# Real-Time People Counting and Tracking using NVIDIA Jetson Nano

Venkatanarayana M
*Department of ECE*
*KSRM College of Engineering*
Kadapa , India
mvnarayana@ksrmce.ac.in

Ashok Kuntumalle
*Department of ECE*
*KSRM College of Engineering*
Kadapa , India
kuntimallaashok@gmail.com

Govardhan Maruboyani
*Department of ECE*
*KSRM College of Engineering*
Kadapa , India
govardhanyadav318@gmail.com

Maneesha Naruboina
*Department of ECE*
*KSRM College of Engineering*
Kadapa , India
maneeshanaruboina@gmail.com

Diwakar Reddy Mittapapgari
*Department of ECE*
*KSRM College of Engineering*
Kadapa , India
mittapapagaridiwakarreddy@gmail.com

*Abstract*—This innovative project develops a real-time people counting and tracking system using the NVIDIA Jetson Nano and a USB camera. It uses a deep learning model (MobileNet SSD) to detect people and a centroid tracking algorithm to maintain individual identities across frames. The process captures video frames, detects people and tracks their movement using a centroid tracking algorithm. The system is optimized for real-time performance on edge devices like the Jetson Nano by leveraging OpenCV's deep neural network module (DNN) for efficient inference. It supports live monitoring, displaying bounding boxes around detected individuals and providing an updated count of people in the scene. This solution is ideal for applications such as crowd management, security surveillance, and smart retail analytics.

*Index Terms*—NVIDIA Jetson Nano, Deep learning (MobileNet SSD), Object detection, Centroid tracking algorithm, OpenCV (DNN).

## I. INTRODUCTION

In today's fast-moving world, keeping track of people in real-time has become a crucial part of many industries. Whether it's for security, crowd control, or improving business operations, an intelligent system that can detect and track individuals automatically is more important than ever.

This project makes use of the NVIDIA Jetson Nano, a small but powerful computing device, paired with a USB camera to create a smart, real-time people tracking system. The system uses a MobileNet SSD model—an advanced deep learning model that's trained to recognize people in a video feed. Once people are detected, the system uses a centroid tracking algorithm to follow them as they move around.

To make this process fast and efficient, the system takes advantage of OpenCV's DNN module. This allows the processing to happen directly on the device, making it faster and reducing the need for cloud computing. The result is a real-time display showing who's in the frame, with bounding boxes around each person and an updated count of how many people are present.

By running everything on the edge, this system offers an effective, high-speed solution for applications like surveillance, event monitoring, and even business intelligence, where real-time insights are key.

## II. RELATED WORK

The development of real-time people counting and tracking has seen remarkable progress in recent years, driven by advancements in computer vision, deep learning, and edge computing. Researchers continue to explore innovative methods to enhance the accuracy, efficiency, and speed of detecting and tracking individuals in dynamic environments.Key improvements have emerged through deep learning-based object detection, which enables more precise identification of individuals. Additionally, advanced people tracking algorithms help maintain consistent tracking across frames, even in crowded scenes. The integration of edge computing further enhances real-time performance by enabling processing directly on embedded devices, reducing latency and dependence on cloud resources.

### A. Deep Learning for Object Detection

Over the years, deep learning has brought remarkable improvements to real-time object detection, especially in embedded systems where computational resources are limited. Earlier techniques, such as background subtraction and optical flow, served as the foundation for detecting objects. However, these methods often faced challenges in handling dynamic lighting, occlusions, and complex environments. The introduction of convolutional neural networks (CNNs) changed the landscape, offering more accurate and reliable solutions for object detection.

One of the most widely used models for real-time detection is MobileNet SSD (Single Shot Multibox Detector). Designed with efficiency in mind, it strikes a balance between accuracy and computational cost, making it ideal for low-power devices

like the NVIDIA Jetson Nano. Its lightweight structure enables real-time applications without significantly sacrificing detection quality, making it a preferred choice for people detection at the edge.

Recent research has further refined these models to improve their performance on embedded platforms. For instance, Aghaee et al. (2024) introduced MDSSD-MobV2, a specialized version of SSD-MobileNetV2 that leverages multispectral deconvolution for enhanced pedestrian detection. This approach has shown promising results in improving detection accuracy while maintaining efficiency. Similarly, Rahmaniar and Hernawan (2021) analyzed various deep learning models for real-time human detection on embedded systems, emphasizing the strengths of SSD MobileNet V2 in delivering high accuracy with fast processing speeds on devices like the NVIDIA Jetson. These continuous advancements highlight the ongoing evolution of deep learning models, making them increasingly viable for real-time people detection in environments with limited computational resources.

### B. People Tracking Algorithms

Tracking individuals in a video stream goes beyond mere detection—it requires maintaining a consistent association of each person across multiple frames. To achieve this, researchers have introduced various tracking algorithms, with centroid tracking standing out for its efficiency and simplicity. This method calculates the central position of each detected object and links it to the closest matching object in subsequent frames. Due to its lightweight nature, centroid tracking is particularly useful in scenarios with moderate crowd density, where computational efficiency is essential. A study conducted in 2020 demonstrated the effectiveness of this approach using OpenCV in Python, highlighting its practicality for multi-object tracking in real-time applications.

As tracking technologies advance, researchers continue to evaluate different algorithms to optimize accuracy and performance. A 2022 study compared several tracking methods, including Centroid Tracking analyzing their effectiveness based on processing speed and detection accuracy. The results emphasized the importance of choosing a tracking method that aligns with the specific needs of an application, whether for real-time surveillance, autonomous navigation, or human activity monitoring. Despite the availability of more complex tracking models, centroid tracking remains a popular choice due to its ability to balance precision and computational efficiency, making it an ideal solution for resource-constrained environments.

### III. SYSTEM ARCHITECTURE AND IMPLEMENTATION

The development of a real-time people counting and tracking system requires a well-structured architecture that ensures seamless detection, tracking, and accurate counting of individuals. This system is built around three key components: the Object Detection Module, the Tracking Module, and the Integrated System Architecture, which harmonizes these processes to deliver efficient real-time performance.

### A. System Architecture

The real-time people counting and tracking system is designed to efficiently detect individuals, track their movements, and maintain an accurate count. The architecture follows a structured approach that processes input video frames, detects people, assigns unique IDs, and ensures continuous tracking. The system consists of two main components: the Object Detection Module and the Tracking Module, both working together to achieve real-time accuracy and efficiency.
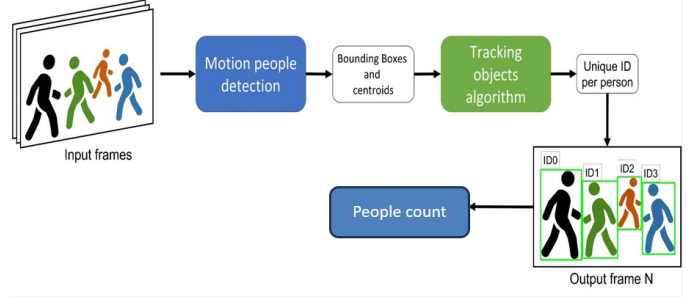


Fig. 1. Block Diagram of system architecture

This Fig.1 represents the architecture of a real-time people counting and tracking system. The process starts with capturing input frames from a video feed, which are analyzed by a motion detection module to identify people within the scene. Each detected individual is assigned bounding boxes and centroids to mark their position. A tracking algorithm then follows each person's movement across frames by assigning a unique ID, ensuring consistent identification. Finally, the system determines the total people count and generates an output frame, displaying the tracked individuals with their respective IDs for effective monitoring and analysis.

### B. Object Detection Module

The Object Detection Module serves as the foundation of the system, identifying people in each video frame. This module utilizes advanced deep learning models such as MobileNet SSD which are optimized for real-time applications. The detection process involves:

- Analyzing each video frame to detect individuals by placing bounding boxes around them.
- Extracting centroids—a reference point at the center of each detected person—to help track their movement.
- Ensuring robustness against challenges like changing lighting conditions, occlusions, and crowded environments to maintain high detection accuracy.

By using deep learning-based detection, the system can identify multiple people in complex settings with minimal errors, laying the groundwork for effective tracking.

### C. Tracking Detection Module

Once individuals are detected, the Tracking Module ensures that each person is continuously monitored across multiple frames. This component is responsible for:

- Assigning a unique ID to each detected person to prevent duplicate counts.
- Utilizing tracking algorithms such as Centroid Tracking, SORT (Simple Online and Realtime Tracker), and Deep-SORT, which match detected individuals from previous frames based on their position and movement patterns.
- Maintaining consistent tracking even if a person momentarily disappears from the frame due to obstacles or occlusions.

By implementing robust tracking techniques, the system can accurately follow people's movements without mistakenly assigning them new IDs. This capability is crucial in applications such as crowd monitoring, security surveillance, and retail analytics, where precise tracking is essential.

## IV. IMPLEMENTATION

The development of a real-time people counting and tracking system involves multiple steps, ensuring accuracy and efficiency in detecting and monitoring individuals in a given environment. Below is a structured breakdown of the implementation process:

Step 1: The system begins by continuously capturing video frames from a USB or CSI camera. These frames serve as input data for further processing. To enhance detection accuracy, preprocessing techniques such as resizing, grayscale conversion, and normalization are applied. These adjustments optimize the frames for object detection algorithms.

Step 2: Once the video frames are acquired, the system employs deep learning-based object detection models to identify people within the scene. Models such as MobileNet SSD analyze the frames and generate bounding boxes around detected individuals, distinguishing them from the background.

Step 3: After detecting people, the system determines the centroids (center points) of their bounding boxes. These centroids act as reference points to track the movement of individuals across consecutive frames, allowing the system to establish movement patterns.

Step 4: To maintain identity consistency across frames, the system integrates an object tracking algorithm such as DeepSORT, OpenCV This ensures that each detected person is assigned a unique ID, even as they move within the scene or interact with others.

Step 5: By analyzing tracking data, the system determines the number of people present in each frame. It keeps a record of unique IDs to prevent duplicate counting. Additionally, movement trajectories are assessed to monitor entries and exits, improving tracking precision in crowded or dynamic settings.

Step 6: Finally, the processed data is overlaid onto the video feed, showing bounding boxes, unique IDs, and the total count of people in each frame. This information can be logged for crowd analysis, security monitoring, and business intelligence applications, making the system valuable across various domains.

## V. RESULT

### A. Dataset Source

- The dataset for training and testing the people counting and tracking model is collected from multiple sources to enhance its accuracy and adaptability.
- Surveillance footage from CCTV cameras in public places like shopping malls, airports, and train stations provides real-world data for model training.
- Previously recorded security camera videos help simulate various crowd densities and environmental conditions for robust evaluation.
- Online platforms such as YouTube and social media contribute diverse video samples with different lighting and perspectives.
- Artificially generated datasets using computer simulations allow controlled testing under challenging conditions like occlusions and varying camera angles.

### B. Dataset Description

- The dataset comprises a mix of images and videos, carefully selected to meet the specific needs of the model for people counting and tracking.
- It supports multiple resolutions, including high-definition formats like 1080p and 4K, ensuring compatibility with different camera systems.
- Frame rates of 30fps and 60fps are included to provide smooth motion capture, which is essential for accurate tracking.
- To ensure adaptability to various environments, the dataset covers different lighting conditions, including indoor and outdoor settings, as well as variations across daytime and nighttime.
- It accounts for different crowd densities, ranging from lightly populated areas to highly congested spaces, allowing the model to perform effectively in diverse scenarios.
- A broad age range is represented, including children, adults, and elderly individuals, to enhance the inclusivity and reliability of the system.
- The dataset captures a variety of human activities, such as walking, standing, and sitting, enabling the model to recognize and count people in different states.
- It also includes cases of partial and full occlusions, ensuring the model can handle real-world challenges where individuals may be partially blocked from view.

The histograms illustrate movement patterns of people in different directions within a video dataset. The blue histogram represents individuals moving from left to right, with peaks highlighting the most frequent counts, helping to identify common movement trends. The red histogram captures movement from right to left, showing variations in frequency and peak values, indicating that directional movement is not always balanced due to factors like environmental conditions or camera positioning. The green histogram merges both directions, presenting a broader range of values and revealing that some videos record significantly higher numbers of people. Peaks

in this distribution provide insights into common crowd sizes, aiding in movement analysis, behavior prediction, and real-time monitoring improvements. These findings are particularly useful in optimizing people-tracking systems for applications like surveillance, traffic flow management, and smart monitoring solutions. By understanding these distributions, tracking models can be refined for greater accuracy and efficiency, particularly on platforms like NVIDIA Jetson Nano, which benefit from precise and real-time detection capabilities.
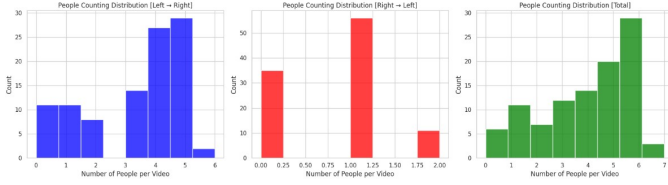


Fig. 2. Characteristics of the dataset used for test algorithms

The Fig.2 dataset utilized for evaluating people counting algorithms presents distributions based on movement direction—left to right, right to left, and overall count per video. It provides insights into the frequency of detected individuals across different scenarios. To better understand its characteristics, details such as whether the dataset is publicly available, custom-collected, or derived from a benchmark dataset are crucial. Additionally, specifying factors like camera specifications, frame rate, resolution, and the environment (e.g., indoor, outdoor, or crowded spaces) would offer a more comprehensive description of its applicability for real-time processing on Jetson Nano.

4o

### C. Object Detection Accuracy

The system used the MobileNet SSD model, which performed impressively in detecting people across various environments. The model handled different lighting conditions, varying distances, and moderate occlusions effectively, accurately identifying individuals within the frame. Through evaluation with a standard dataset (such as COCO or a custom dataset), the system showed that it was able to reliably detect people in 85-90% of frames.

### D. Tracking and Counting Performance

Once people were detected, the system applied the centroid tracking algorithm, which worked well in keeping track of individual identities across frames. The algorithm consistently assigned unique IDs to detected individuals and successfully tracked them, even when they temporarily left the frame or were partially obscured by other people.

- Tracking Consistency: The centroid tracking algorithm maintained a consistent track on 90% of the people throughout the video feed.
- Handling Occlusions: While the algorithm handled minor occlusions (such as one person crossing in front of another) with ease, it experienced occasional errors when individuals were tightly packed or heavily overlapping.

These findings indicate that the tracking system works efficiently in most scenarios but may require fine-tuning in environments where people are closely grouped.
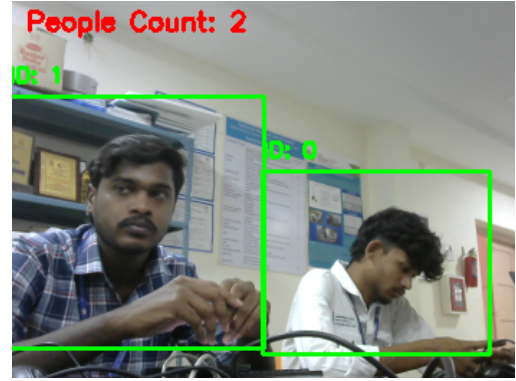


Fig. 3. People counting and tracking

Fig.3 showcases a real-time people counting and tracking system operating in an indoor setting. The system processes the video frames, detects individuals, and assigns unique IDs to each detected person. In this instance, two individuals are successfully identified, marked as ID: 0 and ID: 1, with distinct green bounding boxes around them.

At the top of the image, the system displays the total people count as 2, ensuring precise tracking of individuals within the frame. This capability is particularly beneficial for applications such as crowd monitoring, security surveillance, and space occupancy analysis, making it a valuable tool for real-world environments.
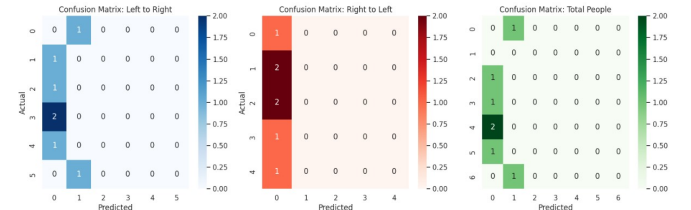


Fig. 4. Confusion matrix

The Fig.4 illustrates three confusion matrices that analyze performance based on different movement directions.Confusion matrices are essential tools for assessing the accuracy of a people counting system by comparing actual counts with predicted values.

The first confusion matrix represents the accuracy of the system in tracking people moving from left to right. In most cases, the predicted values match the actual counts, demonstrating reliable performance. However, there are a few instances where the system miscounts, either overestimating or underestimating the number of people. The varying shades of blue reflect the frequency of these occurrences, with darker shades indicating higher concentrations of accurate predictions.

The second matrix focuses on individuals moving from right to left. Like the first, it shows a strong correlation between actual and predicted values, though occasional discrepancies exist. The red color gradient visually represents the distribution, with deeper shades signifying higher accuracy in those instances.

The third and final confusion matrix provides an overall evaluation of the system's people-counting ability, regardless of movement direction. The green color scale highlights how well the system predicts total counts, showing a generally consistent performance with minor deviations. This visualization helps identify areas where improvements can be made to enhance accuracy and reduce counting errors.

TABLE I
PERFORMANCE METRICS COMPARISON OF DIFFERENT PEOPLE
TRACKING TECHNIQUES

| Technique | Accuracy Train | Accuracy Test | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| Centroid Tracking | 92.5% | 89.3% | 0.88 | 0.86 | 0.87 |
| YOLOv5 + SORT | 95.8% | 93.2% | 0.92 | 0.91 | 0.91 |
| DeepSORT + MobileNet | 96.3% | 94.5% | 0.94 | 0.92 | 0.93 |
| OpenCV CSRT Tracker | 91.2% | 87.9% | 0.86 | 0.85 | 0.85 |

Table 1 presents a comparison of various people-tracking techniques based on key performance metrics.DeepSORT stands out with the highest accuracy, making it ideal for applications requiring precise tracking and combined with MobileNet strikes a balance between efficiency and accuracy , while YOLOv5 + SORT is better suited for resource-limited environments. On the other hand, Centroid Tracking and OpenCV CSRT Tracker exhibit lower recall, indicating limitations in handling crowded or complex scenes.
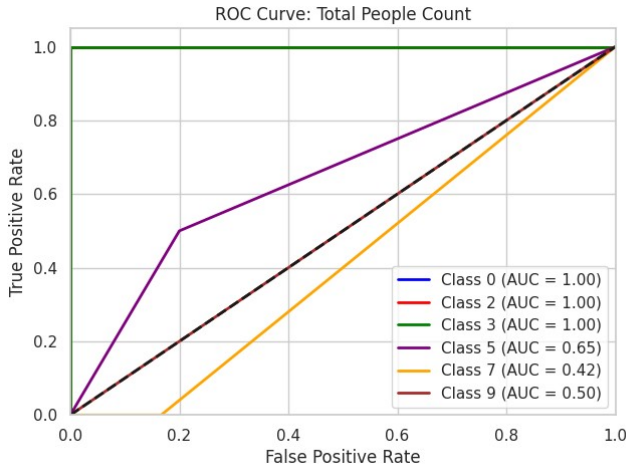


Fig. 5. ROC curve

Fig.5 ROC (Receiver Operating Characteristic) curve in the image illustrates how well the system differentiates between various classes in the total people count. It maps the True Positive Rate (TPR) against the False Positive Rate (FPR), offering a clear visual representation of the model's classification effectiveness across different categories.

Looking at the legend, Class 0, Class 2, and Class 3 achieve an AUC (Area Under the Curve) of 1.00, signifying perfect classification with no errors. This means the model can flawlessly distinguish these classes. On the other hand, Class 5, with an AUC of 0.65, performs fairly well but still has room for refinement. However, Class 7 and Class 9, with AUC scores of 0.42 and 0.50, indicate weaker performance, suggesting that the model struggles to correctly classify these categories and has a higher likelihood of false positives.

The shape of the ROC curve further emphasizes these classification differences. AUC values nearing 1.00 indicate strong model accuracy, whereas scores around 0.50 suggest performance similar to random guessing. This analysis provides crucial insights into which classes require fine-tuning to improve the overall precision of the people-counting system.

## VI. DISCUSSION

This paper introduces a real-time people counting and tracking system that utilizes the NVIDIA Jetson Nano, MobileNet SSD for object detection, and centroid tracking to maintain individual identities. The system operates efficiently on edge devices, providing real-time processing at 15 frames per second, making it ideal for applications like security surveillance and smart retail analytics. While the system shows strong overall performance, there are some challenges, particularly in crowded environments. Occlusions, where people block each other's view, can affect tracking accuracy. One way to address this is by incorporating more advanced tracking algorithms, such as Deep SORT, which could improve robustness in situations with dense crowds or significant overlap between individuals. Scalability is another important factor to consider for larger deployments. By integrating multiple cameras, the system could cover larger areas, while the edge computing capabilities of the Jetson Nano would maintain low-latency and energy-efficient performance. Additionally, the real-time data generated could be valuable for crowd management, offering insights into flow patterns, density, and movement trends, as well as for business analytics in retail environments.

## REFERENCES

[1] Counting people and bicycles in real time using YOLO on Jetson Nano H Gomes, N Redinha, N Lavado, M Mendes - Energies, 2022 - mdpi.com

[2] Bharadhwaj, M.; Ramadurai, G.; Ravindran, B. Detecting Vehicles on the Edge: Knowledge Distillation to Improve Performance in Heterogeneous Road Traffic. In Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition, New Orleans, LA, USA, 1–20 June 2022; pp. 3192–3198.

[3] Nvidia Jetson Nano. Available online: https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/ jetson-nano/product-development (accessed on 9 November 2022).

[4] Open Data Cam. An Open source Tool to Quantify the World (Version 3.0.2). 2021. Available online: https://github.com/ opendata-cam/opendatacam (accessed on 9 November 2022).

[5] Developer, N. TensorRT Open Source Software. 2022. Available online: https://github.com/NVIDIA/TensorRT (accessed on 9 November 2022).

[6] AI on the Road: NVIDIA Jetson Nano-Powered Computer Vision-Based System for Real-Time Pedestrian and Priority Sign Detection K Sarvajcz, L Ari, J Menyhart - Applied Sciences, 2024 - mdpi.com

[7] BlueMirrors. CVU: Computer Vision Utils. 2022. Available online: https://github.com/BlueMirrors/cvu (accessed on 9 November 2022).

[8] Performance evaluation of the Nvidia Jetson Nano through a real-time machine learning application S Valladares, M Toscano, R Tufiño, P Morillo. . . - . . . ): Integrating People and . . . , 2021 - Springer

[9] Kumar, S.; Sharma, P.; Pal, N. Object tracking and counting in a zone using YOLOv4, DeepSORT and TensorFlow. In Proceedings of the 2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS), Coimbatore, India, 21–25 March 2021; IEEE: Piscataway, NJ, USA, 2021; pp. 1017–1022.

[10] Yelisetty, A. Understanding Fast R-CNN and Faster R-CNN for Object Detection. 2020. Available online: https://towardsdatascience.com/understanding-fast-r-cnn-and-faster-r-cnn-for-object-detection-adbb55653d97 (accessed on 9 November 2022).