

Efficient Object Detection and Classification Approach Using an Enhanced Moving Object Detection Algorithm in Motion Videos

K. Madhan¹ and N. Shanmugapriya^{2*}

¹Research Scholar, ²Associate Professor,

Department of Computer Science and Engineering, Dhanalakshmi Srinivasan University, Tamil Nadu, India

E-mail: madhank.phd2022@dsuniversity.ac.in, *shanmugapriyan.set@dsuniversity.ac.in

(Received 30 November 2023; Revised 10 January 2024, Accepted 5 February 2024; Available online 15 February 2024)

Abstract - Object detection and classification have become prominent research topics in computer vision due to their applications in areas such as visual tracking. Despite advancements, vision-based methods for detecting smaller targets and densely packed objects with high accuracy in complex dynamic environments still encounter challenges. This paper introduces a novel and enhanced approach for hyperbolic shadow detection and object classification based on the Enhanced Moving Object Detection (EMOD) algorithm and an improved manta ray-based convolutional neural network optimized for search. In the preprocessing phase, the video data transforms into a sequence of frames, with polynomial adaptive antialiasing applied to maintain frame size and reduce noise. Additionally, an enhanced boundary area preservation algorithm improves the contrast of noise-free edited image sequences. To achieve high-precision detection of smaller objects, the Grib profile of each detected object is also tracked. Finally, a convolutional neural network method employing an enhanced Manta search optimization is deployed for target detection and classification. Comparative experiments conducted across diverse datasets and benchmark methods demonstrate significantly improved accuracy and expanded capabilities in detection and classification.

Keywords: Object Detection, Video Processing, Convolution Neural Network

I. INTRODUCTION

Object detection (OD) plays a pivotal role in photo processing, serving as a crucial component across various domains, including robotic navigation, industrial sensing, intelligent video surveillance, and spatial video surveillance (Wang, M., *et al.*, 2019). Particularly in remote sensing applications, both OD and monitoring hold fundamental importance (Zhan, T. 2019). Among these applications, OD leverages drones for aerial video surveillance (Su, Y. 2021). The primary challenge of OD lies in identifying objects belonging to specific classes and determining their respective locations in images or videos. Machines typically invest substantial time in training and testing to achieve object detection in videos (Fortino, G. 2021). Nonetheless, machines encounter challenges in distinguishing between various objects. To accomplish this objective, a process of continuous knowledge enhancement and the adoption of efficient algorithms are essential (Wang, Y. 2020). Notably, OD algorithms have recently undergone significant advancements (Srivastava, S. 2020).

In computer vision, two significant challenges exist: locating and classifying objects in videos or images (Subban, R. 2018). The primary objective of object localization is to determine the precise positions of objects by delineating specific bounding boxes (BB) around them (Choi, J. W. 2020). Automated object classification represents a crucial challenge with diverse applications. Traditionally, computer vision systems first detect objects and then concentrate on high-quality image processing (Purwanto, D. 2018), aiming to enhance the organizational development and performance of the classifier (W. 2020). Overall, Object Detection (OD) employs various structural models as reference data to identify exciting objects within images (Ro, Y. M. 2019).

To effectively describe objects within images, the selection and extraction of distinctive vital factors play a crucial role in object recognition applications (Sun, J., *et al.*, 2020). Organizational advancement is a critical and dynamic domain within information technology, particularly emphasizing temporal data in various media, including videos (Cao, J. 2019). In video processing, individuals are tracked across frames and consolidated into a unified frame for coherent processing (Kwon, H. 2019). The real-time detection and tracking of moving objects or activities through video surveillance represent tasks of immense significance and complexity (Kumar, R. 2019). The escalating demand for intelligent surveillance systems has intensified the exploration of target tracking as a subject of rigorous research. Object detection and tracking have extensive applications in crowd flow estimation and behavioural analysis (Patan, R., *et al.*, 2020)

The contemporary challenges in OD techniques primarily revolve around computational complexity and accuracy. This study proposes an innovative, eco-friendly classification and OD method aimed at tackling these challenges. The approach harnesses the hyperbolic tangent within the framework of the Enhanced Moving Object Detection (EMOD) algorithm, coupled with a convolutional neural network.

The paper is structured into several sections to address various aspects of the proposed methodology comprehensively. Section 2 conducts an in-depth

exploration of the current landscape in object discovery methods. Following this, Section 3 delineates the proposed system architecture, highlighting preprocessing strategies, algorithms employed, and contrast enhancement techniques. Further elaboration is provided in Section 4, focusing on the utilization of EMOD for object detection. Section 5 delves into the description of the database used, followed by Section 6, which presents a meticulous performance analysis of the proposed method. Finally, in Section 7, the conclusions drawn are based on the empirical findings of the study. This section encapsulates insights gained from the research and offers valuable recommendations for future enhancements in the field of object detection and classification.

II. RELATED WORK

Object Detection (OD) with binary classification is a two-layered method that emphasizes deep learning (DL) to address the challenge of identifying smaller objects within an image frame. At the first level, this method analyzes these proposed smaller objects. The primary emphasis is on individual or character Convolutional Neural Networks (CNNs), which incorporate binarization techniques. This approach effectively reduces the overall number of false positives compared to simple multi-class detection. Moreover, preprocessing techniques are employed to filter out noisy conditions, albeit at the cost of reducing CNNs' detection accuracy. A proposed adjustment replaces the increased Margin pattern with proportional styles.

Additionally, the introduction of a Composite Spatial Pyramid layer is recommended. In the original design, a collapsible core was introduced, leading to a reduction in the overall weight parameters of the layers. This alteration was projected to enhance performance. However, the revised approach may not be suitable for detecting smaller objects. To address this limitation, a multiscale deformable convolutional OD network was proposed by Rao *et al.*, in 2021 (Sugumaran, M. 2021). This multiscale approach employs deep convolutional networks to capture features at various scales. Folding deformations are incorporated to add structural information to mitigate the impact of geometric transformations. Furthermore, attention and region regression techniques are utilized for effective object detection, combining multiscale features through resampling.

The enhancement in detecting smaller goal objects with geometric distortion has been notable. There has been substantial improvement in achieving a balance between velocity and accuracy. However, this scientific approach is unable to track objects within video sequences. Kanimozhi *et al.*, in 2019, presented the advancements in real-time video organization for You Only Look Once (YOLO) networks (Mala, T. 2019). The YOLO Fast OD model is trained to extract information from objects, improving the

Google Inception Net (GoogLeNet) architecture by substituting small convolution operations with local convolution operations. This reduction in configuration parameters also decreases OD time in videos. This method outperforms the original YOLO approach and other conventional methods. Challenges such as high computational load and a low detection rate have been identified.

Ray *et al.*, (2019) presented the hybrid deep learning models, Faster-Region-based Convolutional Neural Network (RCNN) and Mask-RCNN, which consist primarily of two components. The region proposal network is the initial phase in generating a list of region proposals for analyzing input images. Classification makes it possible to determine the return on investment or its absence for each proposal. The ROI pooling layer then accepts proposals and refines the bounding boxes of objects, improving the resulting image. Aerial images of objects are captured within these bounding boxes. The second method for detecting visible sky objects was designed using an R-CNN mask. This model meticulously identifies the actual pixels of each object and detects the bounding box. While this algorithm demonstrates improved discriminative power for specific tasks, the overall increase in accuracy with this method is somewhat limited.

The novel approach introduced in this paper significantly advances the state-of-the-art in object detection and classification. Comparative experiments conducted on diverse datasets and against benchmark methods consistently demonstrate the superior accuracy and expanded capabilities of the proposed method in detecting and classifying objects within challenging contexts. The presented approach not only outperforms existing methods in accuracy but also effectively addresses the challenges related to smaller target detection and noise reduction, making it a promising advancement in the field of computer vision-based object detection.

III. PROPOSED METHOD

The challenge in OD lies in effectively detecting small objects within busy macro scenes, as well as in microscopic modeling, especially for Small Object Detection (SOD) within more significant objects. Enhancing organizational efficiency involves addressing two specific tasks: considering object areas and bounding boxes (BBs) for various purposes like object classification and identification. A novel approach for regularization and classification is proposed using a hyperbolic tangent based on the EMOD within contemporary deep learning models. The video data is initially converted into a sequence of images. In the preprocessing stage, resizing and denoising techniques are applied through the EMOD algorithm to enhance the contrast of noise-free, rescaled image sequences.

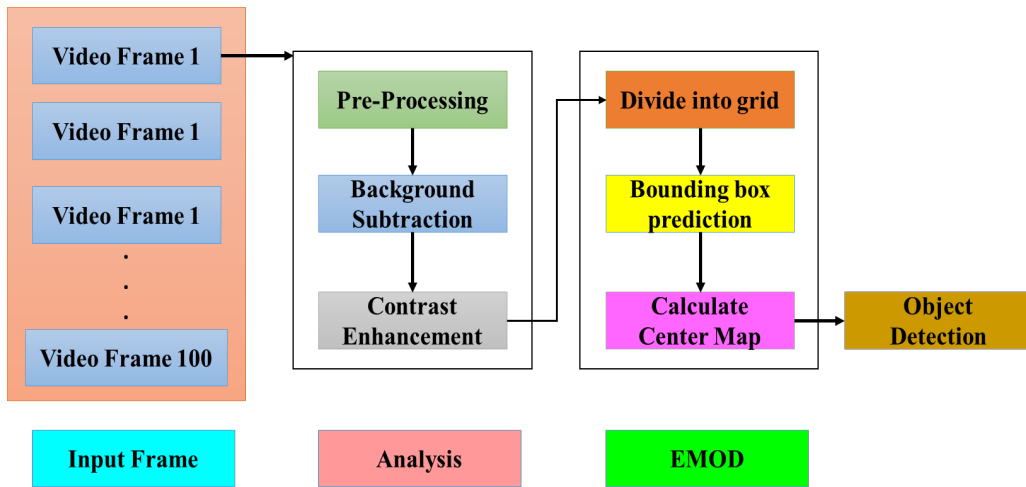


Fig. 1 System Architecture

The resulting contrast-enhanced image sequence is then utilized to train the EMOD algorithm for object detection. This process enables the detection of all objects and their respective bounding boxes, allowing for loss calculation. Consequently, the threshold setting of Grasp is adjusted to detect smaller objects effectively. The EMOD algorithm executes steps within the proposed system architecture, as depicted in Figure 1.

A. Preprocessing

In the initial phase, video data is converted into multiple frames. During preprocessing, the existing noise is eliminated from the input frame photograph using Polynomial Adaptive Edge Perspective Algorithm (PAEPA). To improve clarity, photo resizing is conducted before noise removal. Typically, altering pixel data may diminish image clarity when the image size is increased. Polynomial interpolation was employed to address this issue and enhance image quality. This method ensures equivalent pixel intensities are increased, resolving the clarity concern in the image.

B. Background Subtraction

The input video G_t is first amassed from the Multiple Object Tracking (MOT) dataset. The amassed rec is initially compiled from the MOT dataset. The compiled records, V_f are then converted to G_t as follows.

$$G_t = \{G_1, G_2, G_3, \dots, G_N\} \quad (1)$$

Here, N represents the number of G_t instances. Subsequently, an EMO is utilized to remove unwanted noise from G_t .

The combination of EMOD with the Knowledge Mining Accelerator (KMA) is referred to as EMOD MA, which follows this process.

Algorithm: EMOD MA

1. Consider the input $V_f(G_t)$, and the cluster centroids are represented and start with the value $i = 1, 2, \dots, K$ represents the centroids.
2. Identify the total no of clusters k with start cluster centroids C_i .
3. Use the EDK function, assign (G_t) to centroid C_i , It is represented by

$$E(G_t, C_i) = \sum_{t=1, i=1}^{N, k} e^{\left(\frac{d_e(G_t, C_i)}{2\sigma_e^2}\right)} \quad (2)$$

$$d_e(G_t, C_i) = \sum_{t=1, i=1}^{N, k} \frac{G_t \log G_t + C_i \log C_i}{2} - \frac{G_t + C_i}{2} \log \frac{G_t + C_i}{2} \quad (3)$$

Where $2\sigma_e^2$ represents the ED Function, $d_e(G_t, C_i)$ signifies the entropy-like divergence.

C. Contrast Enhancement

The input for the RCNN consists of the reduced features (R_f), where $f = 1, 2, \dots, n$ represents the number of capabilities. At each time step (τ), the relationship between the input and output of the neural network can be expressed as:

$$H_\tau = \chi(w_{iH}R_f + w_{iH}R_f + B_H) \quad (4)$$

$$O_\tau = \delta(w_{Ho}H_\tau + B_o) \quad (5)$$

Here, $\tau = 1, 2, \dots, m$, O_τ represents the output of the network at time step τ , H_τ denotes the hidden layer (HL), B_H and B_o represents the bias of the HL and the output layer, respectively. w_{iH} and w_{io} correspondingly denote the weight matrices between the input HL and output hidden layer w_{HH} recursion.

The bounded ReLU activation function χ, δ is equation (6)

$$\chi_{R_f} = \delta_{R_f} = \min(\max(0, R_f), Z) = \begin{cases} 0 & R_f \leq 0 \\ R_f & 0 < R_f \leq Z \\ Z & R_f > Z \end{cases} \quad (6)$$

Here, Z signifies the maximum output value.

IV. OBJECT DETECTION USING EMOD

EMOD was originally developed to enhance object detection after contrast enhancement. YOLOv4 is considered an upgraded version of YOLOv3. The backbone, neck, and head are three crucial components or stages. The pre-trained neural network is referred to as the backbone, responsible for extracting essential features from the input image.

A. Divide into Grid

The neck encapsulates multiple top-down and bottom-up pathways to acquire function maps (FMs) at distinct stages. The identification of object classes and bounding boxes has been achieved by the maintainer. Images provided as input are segmented into grids across the picture, which might occur at higher grid levels. To address this issue, reducing the mesh level resolves the problem, leading to more accurate bounding box predictions and reduced loss

functions. Hence, the proposed approach is referred to as EMOD.

B. Bounding Box Prediction

Predefined bounding boxes offer convenient shapes and sizes that adapt during training. The determination of image bounding boxes is preprocessed by utilizing K-means analysis on the training dataset. Crucially, the predictive representation of the bounding boxes is adjusted to minimize the impact of minor changes on predictions. This adjustment enhances the model’s stability. Instead of directly predicting the position and size, a logistic function is employed for prediction, reducing the influence of offsets. This method reshapes predefined anchor boxes for grid cells.

C. Object Detection and Classification

Multiple top-down and bottom-up pathways are integrated within the framework to obtain functional maps (FMs) at distinct stages. The algorithmic framework encompasses the understanding of object features and bounding containers. The input images are segmented into grids across the image, potentially occurring when the grid scale is larger. To address this issue, reducing the grid scale resolves the problem, providing accurate bounding box predictions and minimizing the loss function. Hence, the proposed methodology is referred to as EMOD.

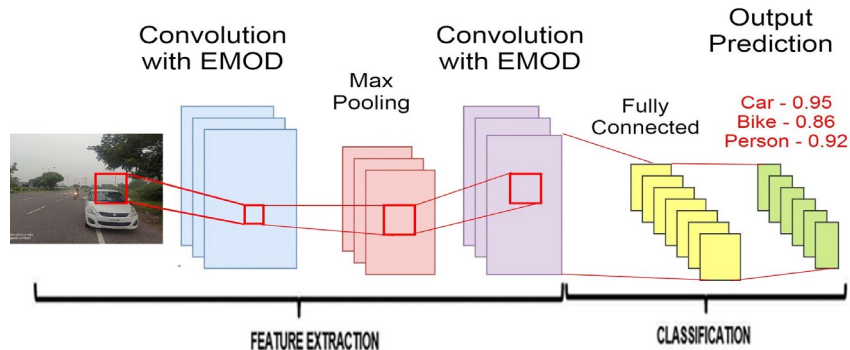


Fig. 2 Structure of EMOD Object Classifier

EMOD tackles the challenge of object detection and classification. Artificial neural networks, particularly CNNs, are widely employed for this purpose. To differentiate between objects, an image is inputted, and adaptable weights are allocated to the objects within the image. The CNN’s ability to detect high-quality features is validated by a Loss Function (LF), which measures the difference between the predicted output and the actual label. Performance prediction is influenced by these losses. The primary goal is to optimize functionality by reducing LF, minimizing training time, and improving model accuracy and structure, as depicted in Figure 2.

EMOD is utilized to optimize the weights of the CNN. Within EMOD, specific layers are exclusively targeted during the optimization process to yield the best solution

amidst the search iterations, preventing the Manta Ray Optimization Algorithm (MRFOA) from obtaining solutions beyond the defined range and upper limit of the problem area. This modification aims to address this issue by introducing EMOD-CNN. The structure of EMOD-CNN is illustrated accordingly.

V. DATABASE DESCRIPTION

The Unmanned Aerial Vehicle (UAV) datasets were utilized for performance analysis. These datasets encompass approximately 10 hours of 720 x 240 40 Hz video footage recorded during typical city traffic conditions involving a car. This footage comprises approximately 150,000 frames with 250,000 bounding boxes and annotations for around 3,500 unique pedestrians. The UAV dataset is divided into

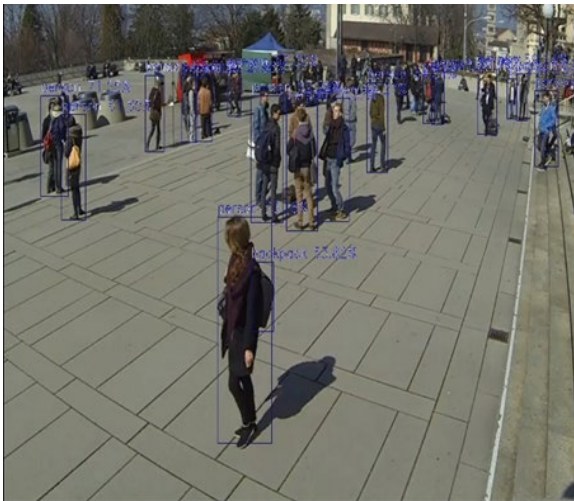
separate training, validation, and test sets, constituting three distinct subsets. Example images from the dataset, along with subsequent image processing, are depicted in the figure. Figure 4a displays the original image from the dataset, while Figure 4b exhibits the preprocessed image using the EMOD algorithm. Additionally, Figures 3a, 3b, 3c, and 3d illustrate the input and output frames of the motion video.



(a) Input Frame 1



(b) Input Frame 2



(c) Output Frame 1



(d) Output Frame 2

Fig. 3 Input and output frames of motion video

VI. PERFORMANCE ANALYSIS

Various ML techniques are evaluated, and their performance metrics are estimated to define the malicious RDP sessions.

$$Accuracy = \frac{TP + TN}{Total\ subjects} \times 100\% \quad (7)$$

$$Precision = \frac{TP}{TP + FP} \times 100\% \quad (8)$$

$$F1\ score = 2 \times \frac{TP}{TP + FN} \quad (9)$$

$$Sensitivity/Recall = \frac{TP}{TP + FN} \times 100\% \quad (10)$$

$$AP\ score = \sum_n (Recall_n - Recall_{n-1}) \times Precision \quad (11)$$

$$Specificity = \frac{TN}{FP + TN} \times 100\% \quad (12)$$

$$GMean = \sqrt{Sensitivity + Specificity} \quad (13)$$

TP is True Positive, TN is True Negative, FP is False Positive, and FN is False Negative values.

TABLE I PERFORMANCE ANALYSIS

Methods	Sensitivity	Specificity	AP Score
Proposed EMOD	99	94	97
Deep Belief Network	94	92	93
Recurrent Neural Network	91	88	89
Convolution Neural Network	92	86	88
Deep Neural Network	88	80	86

The various methods for ML techniques are executed to compare the performance based on Proposed EMOD, DBN, RNN, CNN, DNN with Sensitivity 99,94,91,92,88 and Specificity 94,92,88,86,80 and AP Score 97,93,89,88,86 are shown in Table I and figure 4.

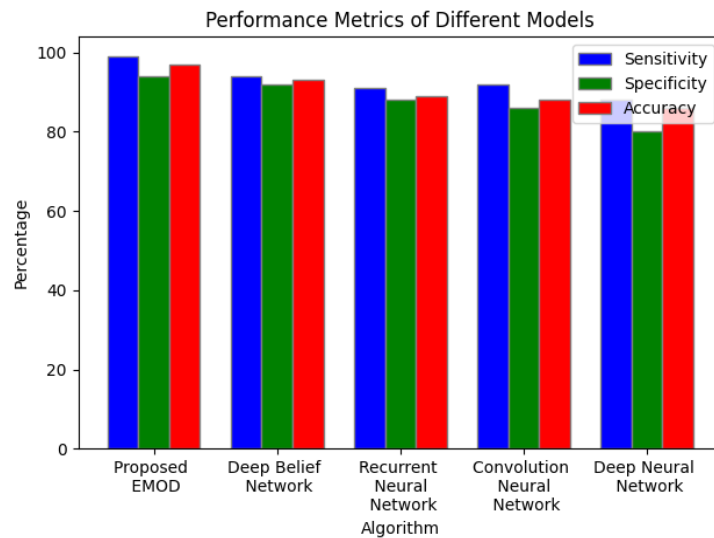


Fig. 4 Comparison of sensitivity, specificity, and AP score

Sensitivity assesses the model’s ability to predict true positive outcomes across all available categories. The capacity of a model to accurately identify true negatives from each available class is evaluated in terms of specificity. Accuracy measures the proximity of the model’s predictions to the actual values. Overall, it has been reported that the

performance outcomes, including sensitivity, specificity, and accuracy, in detecting objects using the proposed technique are observed to be high. Similarly, particularly, Confidence Network deep information such as DBN, RNN, CNN, and DNN are illustrated in Table II and Figure 5.

TABLE II COMPARISON OF ML PERFORMANCE

Methods	Precision	Recall	F-Measure
Proposed EMOD	96	99	97
Deep Belief Network	92	92	92
Recurrent Neural Network	89	89	88
Convolution Neural Network	86	92	86
Deep Neural Network	85	85	81

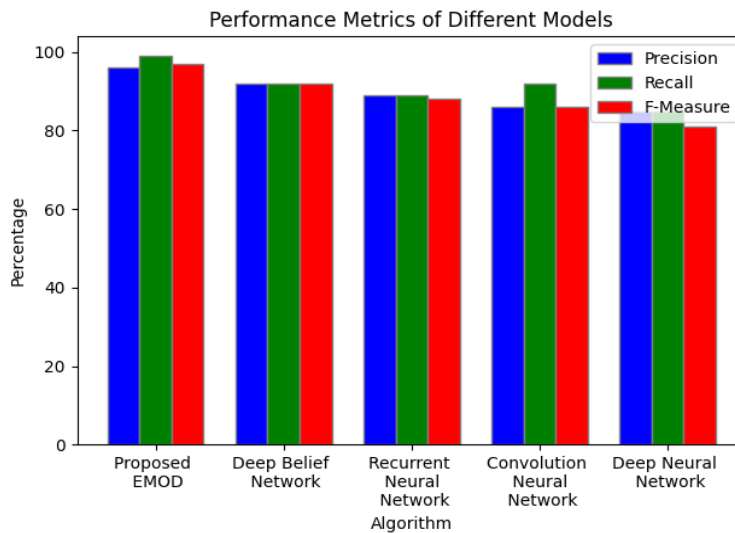


Fig. 5 Comparison of ML parameters

The entire range of predictions for each specific type is assessed using the Accuracy metric. The total range of category predictions generated from each instance in the

dataset defines the recall. In comparison to other outstanding techniques, the proposed approach demonstrates a higher level of accuracy. However, the

recovery and F-measure values of the leading techniques yielded significantly lower performance. Consequently, across all metrics, the proposed approach has proven to be

more effective than both the proposed and existing methods. The accuracy results are presented in Table III and Figure 6.

TABLE III COMPARISON OF ACCURACY

Methods	Accuracy
Proposed EMOD	98
Deep Belief Network	85
Recurrent Neural Network	79
Convolution Neural Network	88
Deep Neural Network	76

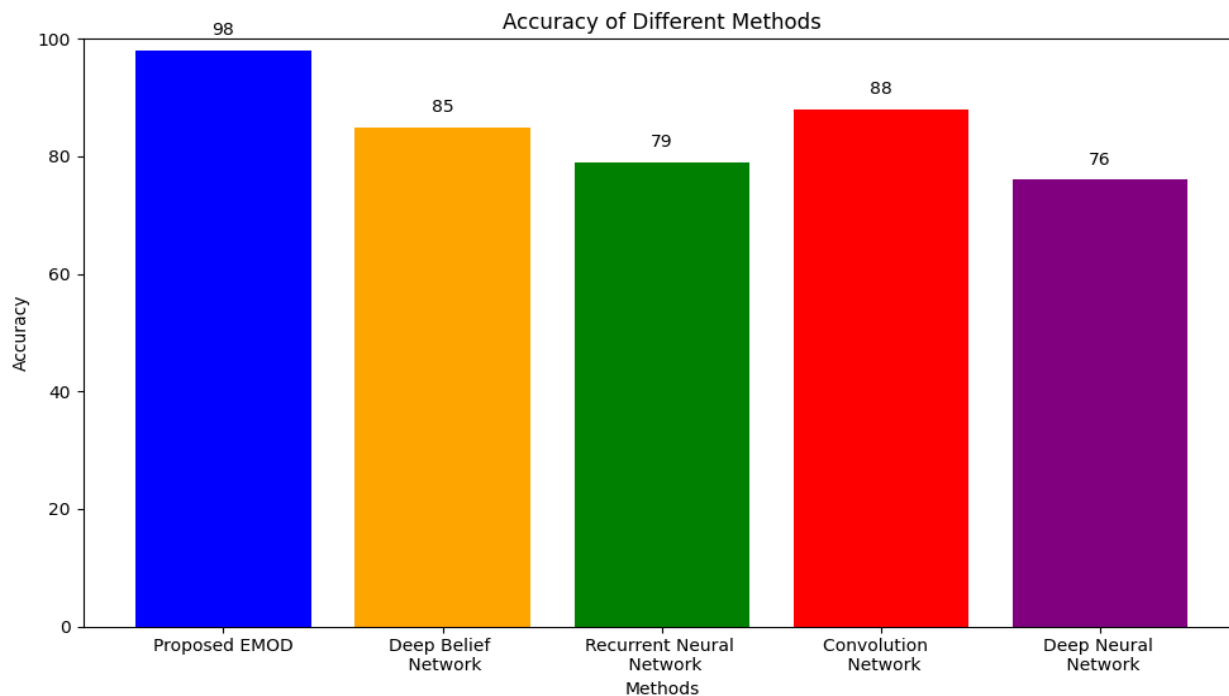


Fig. 6 Comparison of accuracy

VII. CONCLUSION

This study introduces a robust method for object detection and classification using the EMOD algorithm. The efficacy of the proposed method was evaluated across diverse benchmark datasets, notably the UAV datasets, where it underwent rigorous comparison against conventional techniques. The results reveal that EMOD achieves a remarkable accuracy rate of 97.43% in performance analysis, showcasing its superiority in both object detection and classification tasks. The EMOD-based approach demonstrates remarkable improvements in precision, recall, F-measure, sensitivity, and FPR compared to mainstream techniques. Training EMOD with contrast-enhanced images notably enhances the quality of image sequences, positively impacting detection accuracy. Evaluation across benchmark datasets, particularly the UAV datasets, validates the robustness and effectiveness of the proposed approach in diverse real-world scenarios. Moving forward, further advancements can be explored through the integration of more sophisticated algorithms, potentially enhancing the

overall performance and applicability of object detection and classification methodologies. The significant advancements and promising results achieved in this study underscore the potential of EMOD-based techniques in pushing the boundaries of accuracy and efficiency in computer vision tasks, paving the way for enhanced applications across various domains.

ACKNOWLEDGMENT

The authors express their gratitude to the HoD of CSE and the management of Dhanalakshmi Srinivasan University for their invaluable support throughout this research endeavor.

REFERENCES

- [1] Ahmed, I., Din, S., Jeon, G., Piccialli, F., & Fortino, G. (2021). Towards collaborative robotics in top view surveillance: A framework for multiple object tracking by detection using deep learning. *IEEE Chinese Association of Automation Journal of Automatica Sinica*, 8(7), 1253-1270. doi: 10.1109/JAS.2020.1003453.

- [2] Attamimi, M., Nagai, T., & Purwanto, D. (2018). Object detection based on particle filter and integration of multiple features. *Procedia Computer Science*, 144, 214-218. doi:10.1016/j.procs.2018.10.521.
- [3] Bhuvanawari, R., & Subban, R. (2018). Novel object detection and recognition system based on points of interest selection and SVM classification. *Cognitive Systems Research*, 58(1), 1-18. doi: 10.1016/j.cogsys.2018.09.022.
- [4] Hou, B., Li, J., Zhang, X., Wang, S., & Jiao, L. (2019). Object detection and tracking based on convolutional neural networks for high-resolution optical remote sensing video. In *IEEE Int. Geoscience and Remote Sensing Sym.*, Yokohama, Japan, 5433-5436. doi: 10.1109/IGARSS.2019.8898173.
- [5] Kanimozhi, S., Gayathri, G., & Mala, T. (2019). Multiple real-time object identification using single-shot multi-box detection. In *Second Int. Conf. on Computational Intelligence in Data Science*, Chennai, India, 1-5. doi: 10.1109/ICCIDS.2019.8862041.
- [6] Kim, J. U., & Ro, Y. M. (2019). Attentive layer separation for object classification and object localization in object detection. In *IEEE Int. Conf. on Image Processing (ICIP)*, Taipei, Taiwan, 3995-3999. doi: 10.1109/ICIP.2019.8803439.
- [7] Kim, J., Koh, J., & Choi, J. W. (2020). Video object detection using motion context and feature aggregation. In *Int. Conf. on Information and Communication Technology Convergence (ICTC)*, Jeju, Korea (South), 269-272. doi: 10.1109/ICTC49870.2020.9289386.
- [8] Kousik, N. V., Natarajan, Y., Raja, A. R., Kallam, S., Patan, R., et al. (2020). Improved salient object detection using hybrid convolution recurrent neural network. *Expert Systems with Applications*, 166(3), 114064. doi: 10.1016/j.eswa.2020.114064.
- [9] Kumar, A., & Srivastava, S. (2020). Object detection system based on convolution neural networks using single-shot multi-box detector. *Procedia Computer Science*, 171(1), 2610-2617. doi: 10.1016/j.procs.2020.04.283.
- [10] Lee, H., Eum, S., & Kwon, H. (2019). ME R-CNN multi-expert R-CNN for object detection. *IEEE Transactions on Image Processing*, 29, 1030-1044. doi: 10.1109/TIP.2019.2938879.
- [11] Lu, Z., Lu, J., Ge, Q., & Zhan, T. (2019). Multi-object detection method based on YOLO and resnet hybrid networks. In *4th Int. Conf. on Advanced Robotics and Mechatronics (ICARM)*, Toyonaka, Japan, 827-832. doi: 10.1109/ICARM.2019.8833671.
- [12] Pang, Y., & Cao, J. (2019). Deep learning in object detection. In *Deep Learning in Object Detection and Recognition*, 1st ed. Singapore: Springer, 19-57 [Online]. doi: 10.1007/978-981-10-5152-4_2.
- [13] Rani, R., Singh, A. P., & Kumar, R. (2019). Impact of reduction in descriptor size on object detection and classification. *Multimedia Tools and Applications*, 78(7), 8965-8979. doi: 10.1007/s11042-018-6911-7.
- [14] Rao, N. V., Prasad, D. V., & Sugumaran, M. (2021). Real-time video object detection and classification using hybrid texture feature extraction. *International Journal of Computers and Applications*, 43(2), 119-126. doi: 10.1080/1206212X.2018.1525929
- [15] Ray, K. S., & Chakraborty, S. (2019). Object detection by spatio-temporal analysis and tracking of the detected objects in a video with variable background. *Journal of Visual Communication and Image Representation*, 58, 662-674. doi: 10.1016/j.jvcir.2018.
- [16] Yi, S., Ma, H., Li, X., & Wang, Y. (2020). WSODPB weakly supervised object detection with PCS net and box regression module. *Neurocomputing*, 418(12), 232-240. doi: 10.1016/j.neucom.2020.08.028.
- [17] Yin, Y., Li, H., & Fu, W. (2020). Faster-YOLO an accurate and faster object detection method. *Digital Signal Processing*, 102(6), 1-11. doi:10.1016/j.dsp.2020.102756.
- [18] Yu, Q., Wang, B., & Su, Y. (2021). Object detection-tracking algorithm for unmanned surface vehicles based on a radar-photoelectric system. *IEEE Access*, 9, 57529-57541. doi: 10.1109/ACCESS.2021.3072897.
- [19] Yuan, J., Xiong, H. C., Xiao, Y., Guan, W., Wang, M., et al., (2019). Gated CNN integrating multi-scale feature layers for object detection. *Pattern Recognition*, 105(6), 1-33. doi: 10.1016/j.patcog.2019.107131.
- [20] Zhu, Y., Wu, J. S., Liu, X., Zeng, G., Sun, J., et al., (2020). Photon-limited non-imaging object detection and classification based on single-pixel imaging system. *Applied Physics B*, 126(1), 1-8. doi: 10.1007/s00340-019-7373-y.