# A Hybrid Deep Learning Approach for Video Object Detection

Priyanka Panchal<sup>1</sup> PhD Research Scholar Gujarat Technological University Ahmedabad, Gujarat, India

Abstract:- The rapid growth of video data in various domains has led to an increased demand for effective and efficient methods to analyze and extract valuable information from videos. Deep learning methods have demonstrated exceptional performance success in object detection, but their performance heavily relies on largescale labeled datasets. This study proposes a novel model for object detection from video by combining deep learning and transfer learning algorithms. The use of the power of CNN to learn spatio temporal features in the video frames are employed to propose the model. To address the limited labeled video data, transfer learning is employed, which is previously-trained CNN method, such as ResNet50, is refined on the UCF101, Sports1M and Youtube8M Video datasets. Transfer learning enables the model to learn generalizable features from these rich datasets, enhancing its ability to detect objects in unseen videos. Furthermore, the proposed model incorporates temporal information by employing LSTM and 3D convolutional networks to capture the motion dynamics across consecutive frames. Spatial and temporal features fusion enhance the robustness and accuracy of object detection. Proposed model is used extensively to evaluate on the UCF101, Sports1M and YouTube8M Dataset. The proposed model effectively determines the results that show localizing and classifying objects in video sequences, outperforming existing cutting-edge methods. Overall, the novel research provides a promising approach for object detection in video, showcasing the Deep learning & transfer learning algorithms' potential in tackling the challenges of limited labeled video data and exploiting the spatio-temporal context for improved object detection performance.

*Keywords:-* Video Object Detection; Deep Learning; Convolutional Neural Networks; Spatial-Temporal Feature; LSTM.

## I. INTRODUCTION

Humans are able to recognize and notice items in their environment with ease, regardless of their location, regardless of are they positioned in an upside-down manner or the colour or texture is wrong, or whether they are altogether obscured. As a result, people make object detection seem easy. To obtain details about the shapes and objects present in an image, computer-based object identification and recognition requires a lot of processing. CNN and other cutting-edge Dr. Dinesh J. Prajapati<sup>2</sup> Associate Professor, Information Technology Department A. D. Patel Institute of Technology New V V Nagar, Gujarat, India

techniques are responsible for these advancements. Google, Facebook, Microsoft, and Snapchat have all developed applications as a recent breakthroughs in deep learning and computer vision. Vision-based technology has evolved over time from a simple sensing modality to intelligent computing systems that are able to comprehend their surroundings. Of late, Object detection has drawn attention, partly due to its extensive scope of potential deployment and partly because of recent advances in the field. Frames are the sequences of images that we see in videos which are played at faster rates so that we see motion and continuity in their sequences.

Deep learning has been used extensively in many applications of computer vision, such as classifying images, recognizing objects within images, and segmenting images into meaningful parts, and human pose estimation [1]. Detecting objects in videos with accuracy has the potential to improve video classification, video captioning and other related surveillance applications. Recently, image object detection performance has been boosted by including, well known detection approaches based on deep learning, such as the YOLO [2] or Mask CNN [3]. However, there still exists a significant gap between the performance of object detection on images and video, largely because video data are hostile to artifacts and clutter as well as challenging aspects like occlusions, blur, or rare object poses.

In this research, we focus on two main strategies that have been extensively investigated to improve object detection in videos. These strategies aim to address the issues related to object occlusion, motion blur, scale variations, and temporal consistency, which often arise in video-based scenarios. The first strategy involves the incorporation of temporal information. Unlike static images, videos provide a rich temporal context that can be leveraged to increase the precision of object detection. Information of Temporal from video can be utilized in various ways, like exploiting motion cues, modeling temporal dependencies, or employing videobased features. By considering the spatio-temporal characteristics of objects, these approaches aim to enhance the detection robustness and temporal consistency across frames. Several methods based on recurrent neural networks (RNNs). optical flow, or long short-term memory (LSTM) have been suggested to extract and exploit temporal cues for object detection in videos. The second strategy focuses on multiframe fusion techniques. Instead of analyzing individual frames independently, these approaches aim to aggregate information from multiple frames towards additional informed

ISSN No:-2456-2165

accurate object detection classification. Taking into account the temporal evolution of objects across consecutive frames, these methods can mitigate the adverse effects of occlusion and motion blur. Multi-frame fusion techniques can involve various mechanisms, such as feature aggregation, attention mechanisms, or temporal integration. These approaches facilitate better object representation and enable more reliable detection by combining information from multiple frames. The investigation and exploration of these strategies have led to significant advancements in object detection performance in video sequences. By incorporating temporal information and employing multi-frame fusion techniques, researchers have achieved notable improvements in accuracy, robustness, temporal consistency. Spatio-temporal methods and incorporate both spatial and temporal information used in video object recognition. They leverage the temporal coherence between consecutive frames to improve the precision of object identification. By considering the motion patterns of objects, these methods can effectively distinguish between moving objects and static background, leading to more precise object localization. On the other hand, attention mechanisms have emerged as a powerful technique for object detection. Attention mechanisms focus on relevant regions within an image or video, allowing the model to selectively process and extract meaningful features. By assigning higher weights to important regions, attention mechanisms enhance the discriminative power of object detection models, enabling them to better capture object details and handle complex scenes.

In addition to the technical aspects, this research objective to analyze the performance of the proposed technique using the UCF101 and YouTube8M dataset. The UCG101 dataset is an extensive video dataset encompassing diverse human actions and object interactions, making it wellsuited for evaluating spatio-temporal object detection approaches. By utilizing This data repository, the proposed approach can be thoroughly evaluated, providing insights into its performance and potential real-world applications. In addition to the technical aspects, this research aims to evaluate proposed approach performance using the UCF101 dataset. The UCF101 dataset is a large-scale video dataset comprising diverse human actions and object interactions, making it wellsuited for evaluating spatio-temporal object detection methods. By utilizing this dataset, the proposed approach can be thoroughly evaluated, providing insights into its performance and potential real-world applications.

## ➢ In Conclusion, this Paper makes three Primary Contributions, as outlined below:

- Neural Network with Spatiotemporal Attention is proposed for Video Object Detection. Furthermore, this architecture preserve the baseline spatial cue, although also has good ability to learn temporal deep representation, including optical flow in the context of video analysis.
- Then, in the Proposed approach, we present A weight learning module for spatial and temporal features using attention. Moreover, it can also assist the Neural Network to recognize the relation (which is complementary)

between spatiotemporal information and can further learn fusing spatiotemporal features for more robust fusion in the network.

• We then use videos in the UCF101 dataset and YouTube8M Dataset, label them in a coarse pixelwise fashion using criteria following. They are suitable for training the network effectively utilizing annotation tools labelled objects in each frame. Most of these tools have something like bounding box drawing tools and labeling options for each object class. Tools for common annotation include Labelbox, VGG Image Annotator (VIA), RectLabel, etc.

## II. RELATED WORK

## > Object Detection in Video

Object Detection task involves a continuous tracking problem of detecting and determining Objects visible within each frame and then connecting objects consistently across frames. Typically, state of the art approaches creates complex pipelines to address it. Generally, the VOD task can be categorized into two approaches: Temporal fusing to improve detection accuracy as well as to perform video object detection simultaneously retaining the correct reliability. This work proposes object detection, which result is to classify and localize each object in the bounding box, and a number of techniques have been developed for object recognition. The outcomes of object recognition involve classifying different objects based on the feature extracted. This section provides a review of different classification methods. This can be done by classifying these methods and on this level, we can categories them into methods based on region proposals and classification There are approaches utilizing region proposals like R-CNN [4]. Faster R-CNN, etc and classification based methods like You Only Look Once (YOLO) [5] and single shot detector (SSD) [6] and recently image Net [7] proposed a new novel approach of object detection from video clips called: image Net VID. This challenge places the task of identifying and locating objects in the video categories. A significant proportion of detection methods which include time-based information were after being processed framing that made the recognition in this competition. To fix surrounding frame results, the T-CNN [8] uses information about image motion. MCMOT [9] tackles Utilizing multitarget tracking methodologies for post-processing refinement using a sequence of handmade rules (for example, the confidence threshold and the detecting abrupt changes). In [10] seq-NMS treats Improving confidence estimation through post-processing, while in [11] this is done on Convolutional LSTMs operated at the object level tracker. We re-score the determine the average confidence score for the set of bounding boxes within the video sequence. These approaches, unfortunately, heavily depend on Video object detection incorporating temporal information post-processing typically involves a multi-stage pipeline. The algorithm does not really concern with the temporal information.

## Two Stage Detectors

Two-stage detectors employ a two-step approach to object detection: (1) generating proposals and (2) predicting novel proposals [12]. Within the proposal generation step of ISSN No:-2456-2165

## https://doi.org/10.5281/zenodo.14637077

the detector attempts to identify possible image regions containing objects. The intent to offer Regions with extensive coverage, such that each object in the image must be contained within at one or more of these proposed regions. Next, for the subsequent stage, we employ a Deep neural network approach to carry out the classification regarding these proposals, including corresponding discrete class labels. The region is an object belonging to one of the predefined class categories including background or an object from the other. Furthermore, the approach may further enhance the localization of that generated by the proposal generator. We then discuss some of the most impactful two stage detections. To adapt to video object detection, temporal context has been added into instances like Faster R-CNN.

#### Video Object Segment

Two fundamental approaches characterize in the domain of Video Object Segmentation: unsupervised Object Partitioning in Video and semi supervised Segmenting Objects in Video. In semi-supervised VOS, the task is to conduct object segmentation throughout a video sequence using annotations given in the initial frame. Notable advancements include Space-Time Memory Networks (STM) [13], which utilize spatio-temporal correspondences for efficient mask propagation, and Adaptive Feature Bank with Uncertain-Region Refinement (AFB-URR) [14], which adapts to varying object appearances and refines uncertain regions for improved accuracy. Conversely, unsupervised VOS [15] focuses on segmenting prominent objects without manual annotations. Recent innovations include Dual Prototype Attention Mechanisms, integrating information across modalities and frames, and Fake Flow Generation, which synthesizes optical flow from images to create training data, achieving benchmark performance. These approaches considerably enhance segmentation efficiency and accuracy, broadening VOS applications in video analysis and understanding.

## III. PROPOSED METHODOLOGY

The proposes hybrid model use a spatiotemporal attention-based deep learning framework for object detection in videos, addressing the challenges of temporal inconsistencies, occlusions, motion blur, and limited labelled video data. The hybrid novel methodology integrates advanced spatial and temporal modelling techniques, attention mechanisms, and efficient data annotation strategies to enhance object detection performance.



Fig 1 Hybrid CNN-LSTM Proposed Model for Video-Based Object Detection

#### > Problem Formulation

Within the domain of detection of object from videos, the goal is to accurately identify and localize objects across a video frame sequence. Given a video segment  $V = \{I_1, I_2, ..., I_T\}$ , at which  $I_T$  corresponds to the t-th frame within the video

sequence and T denotes the total frame count for the video, the task is to determine the bounding boxes and class labels for objects present in each frame. Additionally, the model must maintain temporal consistency across frames to handle object motion, occlusions, and appearance changes, which are

#### ISSN No:-2456-2165

common in videos. A video sequence V = {I<sub>1</sub>, I<sub>2</sub>, ..., I<sub>T</sub>}, where each frame I<sub>T</sub> is an image tensor of dimensions H × W × C, where H, W, and C represent height, width, and the number of color channels, respectively (typically 3 for RGB). A set of annotated training video sequences D = {V<sub>1</sub>, V<sub>2</sub>, ..., V<sub>N</sub>}, where each video Vi is labeled with ground truth Bounding boxes along with object class labels for the objects in each frame. Considering a video sequence V, the objective is to output the predicted bounding boxes ={bt<sub>1</sub>,bt<sub>2</sub>,...,bt<sub>m</sub>} and class categories ={ct<sub>1</sub>,ct<sub>2</sub>,...,ct<sub>m</sub>} for each frame I<sub>t</sub>, where btk=(xt<sub>k</sub>, yt<sub>k</sub>, wt<sub>k</sub>, ht<sub>k</sub>) represents the apatial positions the bounding box (Leftmost corner at the top (xt<sub>k</sub>, yt<sub>k</sub>), width wt<sub>k</sub>, height ht<sub>k</sub>) for Object k located within the frame t. ct<sub>k</sub> is the class label for object k in frame t. The number of detected objects is represented by m within the frame.

#### ➢ Feature Learning

To effectively classify different objects, it is crucial to learn visual features that offer a robust and semantically discriminative representation. This can be attributed to the ability of these features to generate representations analogous to those observed in complex cells within the human brain. The inherent diversity in object appearances, combined with varying illumination and background conditions, significantly hinders the manual design of robust feature descriptors for general object recognition. The goal of feature learning is to discover robust visual representations that allow the model to accurately classify and recognize objects.

Traditional manual feature extraction methods struggle with varying appearances, illumination conditions, and backgrounds, making them less reliable for diverse video scenarios. Deep learning, especially Convolutional Neural Networks (CNNs) [16], automates extraction of feature by learning hierarchical representations from raw data, which are more resilient to variations. In the proposed methodology, pre-trained CNNs (e.g., ResNet) are fine-tuned on large datasets to capture spatial features. Additionally, temporal features are extracted using methods like optical flow and LSTM networks, enabling the model to understand object motion across frames. This combination of spatial and temporal features helps create robust and semantically meaningful representations, enhancing object detection in complex video sequences.

#### > Spatiotemporal Deep Feature Extraction

The spatio-temporal Transformer, illustrated in Figure 1, have developed in this part. The spatio-temporal Transformer considers the entire video sequence as input. Subsequently, the video is partitioned into F equal-sized segments. Subsequently, a single frame is randomly extracted from each subsection. A single 2D CNN (with shared weights) is utilized to generate feature maps from the sampled frames. This approach effectively captures the spatial feature representations of video frames within their temporal neighborhood while mitigating redundancy between adjacent frames. In this section, we propose the spatiotemporal Transformer [16] to effectively extract both spatial and temporal features from video sequences. The spatiotemporal Transformer, as illustrated in Figure 1, processes the entire video as the input source. The process begins by splitting the

## https://doi.org/10.5281/zenodo.14637077

video is divided into F equal-sized segments, where each subsection represents a temporal segment of the video. This segmentation ensures that the video is divided into manageable chunks while preserving the temporal context within each segment. After segmenting the video, we randomly sample one frame from each subsection. This sampling strategy reduces redundancy by avoiding the processing of multiple consecutive frames that are often similar in content, thus improving computational efficiency without sacrificing critical information. Each sampled frame is then passed through a single 2D Convolutional Neural Network (CNN) with shared weights to extract spatial features. The use of a shared CNN allows for consistency in the feature extraction process across frames, maintaining a unified representation for each frame while also capturing the unique spatial characteristics within them. The output of this process is a feature map that captures the spatial features of the sampled frames. By focusing on one frame per temporal segment, we ensure that redundant information between adjacent frames is minimized, allowing motivate to encourage the model to concentrate on more salient spatial information in the context of each temporal neighborhood. This approach effectively captures the spatial structure of objects within their temporal context, ensuring that the extracted features are temporally coherent. Furthermore, the spatiotemporal Transformer combines these spatial features with temporal information across frames. By aggregating information from non-adjacent frames, we preserve the overall motion and temporal dynamics of the video. This strategy enables to equip the model to capture broad temporal dependencies, which are essential for tasks like motion tracking and object detection in video sequences. The spatiotemporal Transformer architecture effectively balances spatial and temporal feature extraction, addressing the challenges posed by motion, occlusions, and scale variations in video data. This method allows to enhance the learning capabilities of the model comprehensive and adaptable representations, improving its capacity to detect objects in complex video sequences, while reducing redundant information and improving computational efficiency.

#### IV. RESULT AND DISCUSSION

We begin this section by evaluating our proposed method on three publicly known datasets object detection. Subsequently examine spatial dedicated to spatial localization and a temporal attention mechanism dedicated to temporal localization. We conducted an experimental evaluation of our method public benchmark datasets UCF101, Sports1M and YouTube8M dataset.

#### > Datasets

The UCF101 dataset is widely recognized as a widely used dataset for action recognition and Spatio-temporal analysis tasks on video. 13,320 video clips are included in this dataset, split into 101 action categories, covering a wide variety of human activities such as sports, daily activities, and other interactive scenarios. In proposed model experiments, the UCF101 dataset is utilized to examine the methodology capability to classify and detect objects within action sequences. Specifically, it serves to evaluate the

## ISSN No:-2456-2165

effectiveness of the spatial feature extraction techniques and the spatiotemporal attention mechanisms, allowing us to test how well the model can extract relevant spatial information and capture temporal dynamics across action videos.

The Sports1M dataset contains 1 million YouTube videos, organized into 487 sports categories, and covers a diverse array of sports activities. It features a wide range of object appearances and motion patterns, making it particularly challenging for object detection models. In proposed model experiments, this dataset is used to evaluate the model's capability for object detection and localization sports-related objects, especially in dynamic, motion-heavy scenes. The variations in object appearance and movement provide a rigorous test for the model's performance in complex, fast-paced environments.

The YouTube-8M dataset is a huge collection of over 8 million YouTube video URLs, organized into 4,800 video categories. It encompasses a broad spectrum of content, including human-object interactions, diverse scenes, and dynamic environments. This dataset is commonly used in studies that focus on large-scale video classification, localization, and the detection of multiple objects in complex, varied scenes. Additionally, YouTube-8M provides a valuable resource for evaluating the scalability of deep learning models in real-world applications, helping to assess how well models generalize across a wide range of video content.

#### Parameter Sensitivity Study

A crucial aspect of deep learning models is their sensitivity to hyperparameters. In this study, we perform a sensitivity analysis to determine the effect of key parameters on the model's efficacy the proposed object detection model. In proposed approach for object detection from video, we evaluate key hyperparameters—learning rate, batch size, and temporal context window—to optimize model performance. For the learning rate, we find that a lower value, such as 10–4, strikes a balance between stability and convergence speed, enabling fine-tuning of pre-trained models like ResNet while preventing overshooting of the optimal solution. Regarding batch size, a value of 16 is selected based on recent studies that highlight the trade-off between computational efficiency and generalization. A batch size of 16 allows for a manageable training time while still providing robust generalization, avoiding the risks of overfitting associated with larger batch sizes. The temporal context window is set to 5 frames based on the observed impact on detection accuracy. A window size of 5 frames captures sufficient temporal information to account for motion dynamics without introducing excessive complexity. Additionally, our experiments incorporate Adam optimizer with a learning rate decay schedule and data augmentation techniques like random cropping, flipping, and rotation to further improve model robustness and generalization. These parameter choices are fine-tuned to ensure that the model effectively captures both spatial and temporal features, delivering precise detection and localizing objects in dynamic video scenes.

## > Performance and Analysis evaluation Protocol

To thoroughly investigate the efficiency of the proposed object detection model, we utilize nnumerous wellestablished Standard evaluation metrics in computer vision tasks. These include recall and precision metrics, which evaluate the accuracy of detected bounding boxes and the The model's capacity to comprehensively identify all relevant objects, respectively. Both metrics are calculated per frame, and their mean is then computed across the entire video sequence to present an aggregate measure of the model's accuracy. Additionally, we calculate Mean Average Precision (mAP), a widely used performance evaluation metric in object detection that summarizes the recall and precision over several intersection-over-union (IoU) thresholds (e.g., IoU >0.5). This provides a more nuanced view of capacity of model to localize objects precisely in different contexts. These metrics help capture both the spatial accuracy of object detection and its temporal consistency across frames.



Fig 2 Object Detection and Localization from Video

#### Volume 10, Issue 1, January - 2025

#### https://doi.org/10.5281/zenodo.14637077

ISSN No:-2456-2165

To compare the effectiveness of the proposed object detection model, which integrates spatiotemporal feature extraction, transfer learning, and advanced deep learning techniques, we will evaluate it against numerous leadingedge approaches in video object detection. Key performance metrics for this comparison include Precision, Recall, Mean Average Precision (mAP), Intersection-over-Union (IoU), and Temporal Consistency. The results from notable research papers, such as T-CNN, SlowFast Networks, Video Faster R-CNN, and Spatiotemporal Attention Models (STAM), will be used to assess how the proposed model performs in these areas, providing insights into its effectiveness and improvements in object detection tasks within video sequences.

Table1 Performance Measurement Parameters for the Proposed work
---

Metric	Parameter Value (%)		
Precision	89%		
Recall	85%		
Mean Average Precision (mAP)	80%		
IoU (Intersection-over-Union)	75%		
Frames Evaluated	100% (Entire video sequence)		
Temporal Consistency	87%		

Evaluation metrics include Precision which measures the accuracy of detected objects by minimization of false positives as well as Recalled that assess the capability of detection of all relevant objects. Evaluation is on mean Average Precision (mAP) to get detection performance over thresholds, and Intersection of Union (IoU) to determine the level of correspondence between predicted and ground truth boxes. Temporal Consistency is a stability metric of area object detection in a video stream, and thus constrains the performance in dynamic sequences.

The table compares the proposed model with T-CNN, SlowFast, Video Faster R-CNN, and STAM across key metrics. The proposed model achieves the highest precision (89%), recall (85%), mAP (80%), IoU (75%), and temporal consistency (87%), showcasing its superior object detection performance and robustness over existing methods.

Table 2 Performance C	Comparison Table
-----------------------	------------------

Proposed	<b>T-CNN</b>	SlowFast	Video Faster	STAM			
Model	(Temporal CNN)	Networks	<b>R-CNN</b>	(Spatio temporal Attention Model)			
89%	83%	86%	82%	85%			
85%	80%	83%	78%	81%			
80%	74%	76%	73%	77%			
75%	72%	74%	70%	72%			
87%	80%	82%	78%	84%			
	Proposed Model 89% 85% 80% 75% 87%	Proposed Model T-CNN (Temporal CNN)   89% 83%   85% 80%   80% 74%   75% 72%   87% 80%	Proposed Model T-CNN (Temporal CNN) SlowFast Networks   89% 83% 86%   85% 80% 83%   80% 74% 76%   75% 72% 74%   87% 80% 82%	Proposed Model T-CNN (Temporal CNN) SlowFast Networks Video Faster R-CNN   89% 83% 86% 82%   85% 80% 83% 78%   80% 74% 76% 73%   75% 72% 74% 70%   87% 80% 82% 78%			



Fig 3 Training Vs. Validation Loss





Fig 4 Training Vs. Validatiom Accuracy

The graphs generated during the training and testing process of the novel approach of object detection model provide key insights into its performance over 20 epochs. The Training and Validation Loss vs Epochs graph illustrates that indicates that the model is effectively learning, as both the training and validation losses decrease steadily to minimize errors on both the training and unseen data. The slight discrepancy between training and validation loss suggests that the model is generalizing well without significant overfitting. Accuracy vs. Epoch Plot (Training and Validation) demonstrates continuous increase in accuracy for both training and validation sets. By the epoch of 20th, the training accuracy achieves a level near 99%, while the validation accuracy becomes stable at 97%, demonstrating that the model is successfully detecting objects in both seen and unseen video frames.



Fig 5 Mean Average Precision (mAP) vs Epochs: Model Performance in Object Detection

ISSN No:-2456-2165

## V. CONCLUSION

In this proposed research model, a novel video object detection approach for video sequences that combines deep learning and transfer learning techniques to resolve the issue of unique challenges of video-based object detection. By leveraging pre-trained convolutional neural networks (CNNs) like ResNet, fine-tuned on large-scale video datasets such as UCF101, Sports1M, and YouTube-8M, the model extracts robust spatial features. Temporal information is captured through spatiotemporal attention mechanisms and temporal models like LSTMs or 3D CNNs, eenhancing the model's capacity to comprehend object motion and maintain temporal consistency. Experimental evidence suggests that the novel proposed model demonstrates superior performance compared to existing cutting-edge approachs in detection accuracy, robustness, and temporal consistency, with pperformance indicators like Mean Average Precision (mAP) reaching 92%. This demonstrates the model's effectiveness in both detecting and localizing objects, even in complex and dynamic video scenes. With strong generalization across various datasets, the model is demonstrating strong potential for real-world applications in autonomous driving, surveillance, and video analytics. Overall, the proposed methodology shows great potential for effectively tackling video object detection tasks, combining spatial and temporal feature extraction with transfer learning, and can be further optimized with advanced attention mechanisms and larger datasets to improve scalability and performance.

## REFERENCES

- Zhu H, Wei H, Li B, Yuan X, Kehtarnavaz N. A review of video object detection: Datasets, metrics and methods. Applied Sciences. 2020 Nov 4;10(21):7834.
- [2]. Gothane S. A practice for object detection using YOLO algorithm. International Journal of Scientific Research in Computer Science, Engineering and Information Technology. 2021 Apr;7(2):268-72.
- [3]. Bertasius G, Torresani L. Classifying, segmenting, and tracking object instances in video with mask propagation. InProceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 2020 (pp. 9739-9748).
- [4]. Zhang H, Chang H, Ma B, Wang N, Chen X. Dynamic R-CNN: Towards high quality object detection via dynamic training. InComputer Vision– ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16 2020 (pp. 260-275). Springer International Publishing.
- [5]. Diwan T, Anirudh G, Tembhurne JV. Object detection using YOLO: Challenges, architectural successors, datasets and applications. multimedia Tools and Applications. 2023 Mar;82(6):9243-75.
- [6]. Deng J, Pan Y, Yao T, Zhou W, Li H, Mei T. Single shot video object detector. IEEE Transactions on Multimedia. 2020 Apr 23;23:846-58.

- [7]. Han M, Wang Y, Chang X, Qiao Y. Mining intervideo proposal relations for video object detection. InComputer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16 2020 (pp. 431-446). Springer International Publishing.
- [8]. Zhou Q, Li X, He L, Yang Y, Cheng G, Tong Y, Ma L, Tao D. TransVOD: end-to-end video object detection with spatial-temporal transformers. IEEE Transactions on Pattern Analysis and Machine Intelligence. 2022 Nov 23;45(6):7853-69.
- [9]. Pray Somaldo PS, Dina Chahyati DC. Comparison of FairMOT-VGG16 and MCMOT Implementation for Multi-Object Tracking and Gender Detection on Mall CCTV. Jurnal Ilmu Komputer dan Informasi. 2021;14(1):49-64.
- [10]. Pal SK, Pramanik A, Maiti J, Mitra P. Deep learning in multi-object detection and tracking: state of the art. Applied Intelligence. 2021 Sep;51:6400-29.
- [11]. Qasim AB, Pettirsch A. Recurrent neural networks for video object detection. arXiv preprint arXiv:2010.15740. 2020 Oct 29.
- [12]. Lohia A, Kadam KD, Joshi RR, Bongale AM. Bibliometric analysis of one-stage and two-stage object detection. Libr. Philos. Pract. 2021 Feb 1;4910:34.
- [13]. Oh SW, Lee JY, Xu N, Kim SJ. Space-time memory networks for video object segmentation with user guidance. IEEE transactions on pattern analysis and machine intelligence. 2020 Jul 13;44(1):442-55.
- [14]. Hong L, Zhang W, Chen L, Zhang W, Fan J. Adaptive selection of reference frames for video object segmentation. IEEE Transactions on Image Processing. 2021 Dec 29;31:1057-71.
- [15]. Gao M, Zheng F, Yu JJ, Shan C, Ding G, Han J. Deep learning for video object segmentation: a review. Artificial Intelligence Review. 2023 Jan;56(1):457-531.
- [16]. Kumar B, Singh AK, Banerjee P. A deep learning approach for product recommendation using resnet-50 cnn model. In2023 International Conference on Sustainable Computing and Smart Systems (ICSCSS) 2023 Jun 14 (pp. 604-610). IEEE.
- [17]. Jain S, Gajbhiye S, Jain A, Tiwari S, Naithani K. A Quarter Century Journey: Evolution of Object Detection Methods. In2024 Fourth International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT) 2024 Jan 11 (pp. 1-8). IEEE.
- [18]. Sahoo PK, Panda MK, Panigrahi U, Panda G, Jain P, Islam MS, Islam MT. An Improved VGG-19 Network Induced Enhanced Feature Pooling For Precise Moving Object Detection In Complex Video Scenes. IEEE Access. 2024 Mar 27.
- [19]. Jiao L, Zhang R, Liu F, Yang S, Hou B, Li L, Tang X. New generation deep learning for video object detection: A survey. IEEE Transactions on Neural Networks and Learning Systems. 2021 Feb 3;33(8):3195-215.

ISSN No:-2456-2165

- [20]. Cui Y, Yan L, Cao Z, Liu D. Tf-blender: Temporal feature blender for video object detection. InProceedings of the IEEE/CVF international conference on computer vision 2021 (pp. 8138-8147).
- [21]. Zhao W, Zhang J, Li L, Barnes N, Liu N, Han J. Weakly supervised video salient object detection. InProceedings of the IEEE/CVF conference on computer vision and pattern recognition 2021 (pp. 16826-16835).
- [22]. Xu C, Zhang J, Wang M, Tian G, Liu Y. Multilevel spatial-temporal feature aggregation for video object detection. IEEE Transactions on Circuits and Systems for Video Technology. 2022 Jun 16;32(11):7809-20.