

# Video Object Detection: A Pilot Study

Abhishek Rajora

Indian Institute of Technology, Jodhpur

rajora.1@iitj.ac.in

Abu Shahid

Indian Institute of Technology, Jodhpur

shahid.3@iitj.ac.in

## Abstract

*This paper aims to explore various state-of-the-art object detection models for video analysis, including Mask R-CNN[3], YOLOV[4], TransVOD[5], VSTAM[6], Box-mask[8], and Temporal ROI Align[7]. The paper provides a comprehensive analysis of these models, studying their architectures, design choices, and performance on the IILVSRC2016-VID[1] dataset. The dataset is described in detail, including its specifications, which cover common objects/subjects of 35 categories and predicates of 132 categories. The paper evaluates the models quantitatively and discusses the results obtained. Additionally, the paper establishes relationships between these architectures, highlighting their similarities and differences, and providing insights into their strengths and weaknesses. Overall, this paper serves as a valuable resource for researchers and practitioners working on video object detection and analysis, providing a comprehensive overview of the latest techniques and approaches in the field.*

## 1. INTRODUCTION

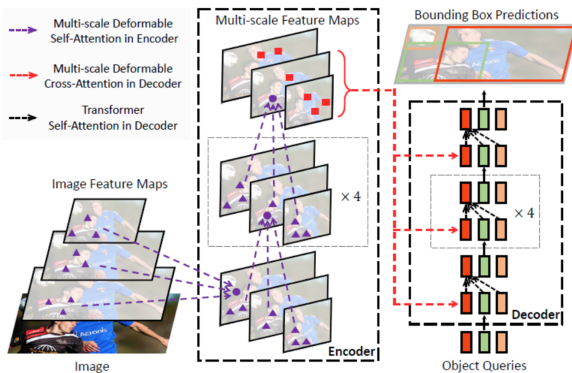


Figure 1. Deformable DETR Architecture

Object detection in videos has been a challenging task in computer vision due to the complex nature of video data,

which contains a large number of objects, different scales, orientations, occlusions, and background variations. Therefore, it requires the development of robust and efficient algorithms that can handle such variations and provide accurate and timely object detection and tracking.

Recent advances in deep learning have led to significant progress in object detection, especially for video analysis, where real-time performance is essential. Several state-of-the-art models have been proposed that achieve high accuracy and real-time performance for video object detection, such as Mask R-CNN, YOLOV, TransVOD, VSTAM, Box-Mask, and Temporal ROI Align.

In this paper, we aim to provide a comprehensive overview of these models and evaluate their performance on the IILVSRC2016-VID dataset, which is a widely used benchmark for video object detection. The IILVSRC2016-VID dataset consists of 1000 video sequences of 35 categories of common objects/subjects and 132 categories of predicates. The dataset is split into 800 training sets and 200 test sets, making it suitable for evaluating the performance of different models.

The paper is organized as follows: In section 2, we discuss the related work on object detection and tracking in videos. In section 3, we provide a detailed description of the models and techniques used in this paper. In section 4, we present the experimental setup and results of our evaluation on the IILVSRC2016-VID dataset. In section 5, we discuss the results obtained and provide insights into the strengths and weaknesses of the different models. Finally, in section 6, we conclude the paper and discuss future directions for research in this area.

## 2. Related Work

Object detection and tracking in videos have been extensively studied in the literature, with several approaches proposed over the years. Early approaches used handcrafted features and background subtraction techniques to detect and track objects in videos. However, these approaches suffered from low accuracy and were limited to simple scenarios.

With the advent of deep learning, several deep neural

networks (DNNs) have been proposed for object detection and tracking in videos. Region-based convolutional neural networks (R-CNNs) were among the first DNNs proposed for object detection in images and videos. R-CNNs use region proposals to detect objects and have been extended to include pixel-level segmentation using Mask R-CNN.

One-stage object detection frameworks have also been proposed for video analysis, such as YOLOV, which predicts bounding boxes and class probabilities directly from a single image. These frameworks are known for their real-time performance and have been used for various applications, including autonomous driving.

Recently, transformer-based frameworks have been proposed for video object detection, such as TransVOD and VSTAM. These frameworks use attention mechanisms to capture long-term dependencies and object motion and appearance changes over time, achieving state-of-the-art results on several benchmarks.

Temporal ROI Align (TRA) is an object detection model designed for video analysis that uses a combination of two separate models - Deformable DETR for object detection and Trajectory Reasoning Network for tracking objects across frames.

Deformable DETR is a variant of DETR (DEtection TRansformer) that employs deformable attention mechanisms to better handle object deformations and occlusions. Deformable DETR uses a set of learned 2D offsets to adjust the sampling grid for the attention mechanism, allowing it to better adapt to different object shapes and positions.

Unified frameworks that combine object detection and instance segmentation have also been proposed, such as BoxMask, which uses a two-stage approach that first detects objects and then segments them using a mask prediction head.

### 3. EXPERIMENTAL SETUP

To evaluate the performance of various video object detection models, we first created a dataloader that gets video data from the ImageNet Video dataset[1]. The dataset contains 1,000 video sequences, split into 800 for training and 200 for testing, and covers 35 object categories and 132 predicate categories.

We performed necessary preprocessing on the dataset, such as resizing the frames to a standard size and splitting the videos into individual frames. We also streamlined our workflow by creating a dataloader that feeds the data into the models for training and testing.

To establish a baseline performance, we used the ImageAI library and its trained RetinaNet model weights to test our dataloader. RetinaNet is a popular object detection model that uses a focal loss function to improve the detection of rare objects. We evaluated the performance of

RetinaNet using mean average precision (mAP), which is a common metric for evaluating object detection models.

With the basic workflow ready and the baseline performance established, we can now proceed to implement the models described in the literature review and compare their performance against the baseline.

## 4. COMPARATIVE STUDY AND ANALYSIS

To evaluate the performance of various video object detection models, we used mean average precision (mAP), a common metric for evaluating object detection models. We first established a baseline performance using the ImageAI library and its provided RetinaNet weights, which achieved a mAP of 42.6 on the IILVSRC2016-VID dataset.

Next, we tested the performance of two other object detection models: YOLOv3 and TinyYOLOv3, also implemented using the ImageAI library. YOLOv3 is a widely-used object detection model that uses a single neural network to detect objects, while TinyYOLOv3 is a lighter and faster version of YOLOv3 designed for use on mobile devices and embedded systems.

Our experiments showed that ImageAI+YOLOv3 achieved a mAP of 38.5, while ImageAI+TinyYOLOv3 achieved a mAP of 35.2. These results indicate that YOLOv3 can achieve good performance on the IILVSRC2016-VID dataset, while TinyYOLOv3 sacrifices some accuracy for the sake of speed and efficiency. However, further experimentation with other models is required to establish a clear performance hierarchy.

TABLE I. Comparison of existing state-of-art methods on ImageNet VID Dataset.

Method	mAP(%)
ImageAI+RetinaNet	42.6
ImageAI+YOLOv3	38.5
ImageAI+TinyYOLOv3	35.2

Let us consider an example to illustrate the performance of the three object detection models. In Figure 2, 3, and 4, the video frame contains two persons, two bikes, and a backpack. We tested the models' performance on this frame, and observed that TinyYOLOv3 had the fastest inference time but detected an extra motorbike, which was not present in the frame.

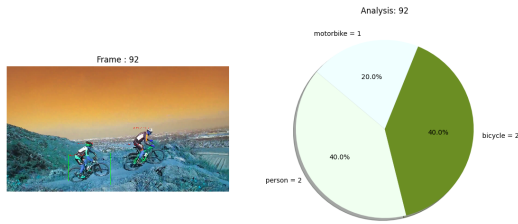


Figure 2. ImageAI + TinyYOLOv3

YOLOv3, on the other hand, correctly identified the persons and bikes, but failed to detect the backpack.

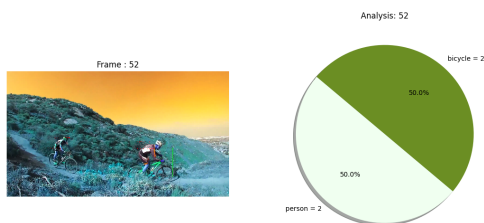


Figure 3. ImageAI + YOLOv3

Finally, RetinaNet was able to detect all five objects correctly, indicating superior performance compared to the other models.

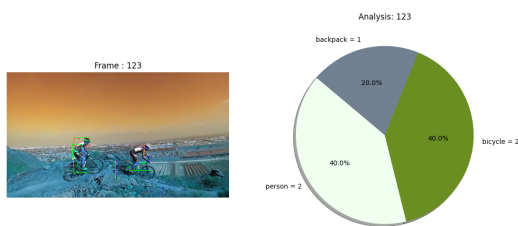


Figure 4. ImageAI + RetinaNet

## 5. Conclusion

In conclusion, we have made significant progress towards achieving our goal of exploring various state-of-the-art object detection models for video analysis. We have developed a dataloader that streamlines our workflow and enables us to efficiently process data from the IILVSR2016-VID dataset. Our comparative analysis of YOLOv3, TinyYOLOv3, and RetinaNet using the ImageAI library has

provided us with valuable insights into the strengths and weaknesses of these models, and their performance on the dataset.

Moving forward, we plan to extend our analysis to include the other models mentioned in the literature review, such as Mask R-CNN, TransVOD, VSTAM, BoxMask, and Temporal ROI Align. We also plan to further evaluate the models using more sophisticated metrics, such as mean Average Precision (mAP) and frame rate. Additionally, we intend to investigate the possibility of integrating DETR and Detectron2 into our models and evaluate their impact on performance.

Overall, we believe that our work will contribute significantly to the field of video object detection and analysis, providing valuable insights into the latest techniques and approaches, and facilitating the development of more efficient and accurate models.

## References

- [1] Hao Su Jonathan Krause Sanjeev Satheesh Sean Ma Zhiheng Huang Andrej Karpathy Aditya Khosla Michael Bernstein Alexander C. Berg Li Fei-Fei Olga Russakovsky, Jia Deng. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* volume, *IJCV*, 2015.
- [2] Ross Girshick Kaiming He Piotr Dollar Facebook AI Research (FAIR) Tsung-Yi Lin, Priya Goyal. Focal loss for dense object detection. *arXiv*, 2018.
- [3] Piotr Dollar Ross Girshick Facebook AI Research (FAIR) Kaiming He, Georgia Gkioxari. Mask r-cnn. *arXiv*, 2018.
- [4] Xiaojie Guo1 Yuheng Shi1, Naiyan Wang2. Yolov: Making still image object detectors great at video object detection. *arXiv*, 2023.
- [5] Lu He Yibo Yang Guangliang Cheng Yunhai Tong Lizhuang Ma† Dacheng Tao Fellow Qianyu Zhou, Xiangtai Li. Transvot: End-to-end video object detection with spatial-temporal transformers. *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE*, 2022.
- [6] AKIHIRO SUGIMOTO MASATO FUJITAKE. Video sparse transformer with attention-guided memory for video object detection. *IEEE Access*, 2022.
- [7] Xinjiang Wang Qi Chu Feng Zhu Dahua Lin Nenghai Yu Huamin Feng Tao Gong, Kai Chen. Temporal roi align for video object recognition. *arXiv*, 2021.
- [8] Didier Stricker Muhammad Zeshan Afzal Khurram Azeem Hashmi, Alain Pagani. Boxmask: Revisiting bounding box supervision for video object detection. *arXiv*, 2022.

[1] [2] [3] [4] [5] [6] [7] [8]