

# **VIDEO OBJECT SEGMENTATION FOR OBJECT DETECTION AND RECOGNITION**

A major project report submitted in partial fulfillment of the requirement  
for the award of degree of

**Bachelor of Technology**

in

**Computer Science & Engineering / Information Technology**

*Submitted by*

**Ansh Mahajan (201137)**

**Snehan Tagnaat (201541)**

*Under the guidance & supervision of*

**Dr. Vipul Kumar Sharma**

**Assistant Professor (SG)**



**Department of Computer Science & Engineering and Information  
Technology**

**Jaypee University of Information Technology, Wagnaghat, Solan -  
173234 (India)**

# CERTIFICATE

This is to certify that the work which is being presented in the project report titled “ **Video object segmentation for object detection and recognition** ” in partial fulfilment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science & Engineering / Information Technology** submitted to the Department of Computer Science And Engineering, Jaypee University of Information Technology, Waknaghat is an authentic record of work carried out by **Ansh Mahajan (201137)** and **Snehan Tagnaat (201541)** during the period from August 2023 to May 2024 under the supervision of **Dr. Vipul Kumar Sharma**, (Assistant Professor(SG), Department of Computer Science and Engineering, Jaypee University of Information Technology).

Ansh Mahajan(201137)

Snehan Tagnaat(201541)

The above statement made is correct to the best of my knowledge.

Supervisor Name: Dr. Vipul Kumar Sharma Designation: Assistant Professor (SG)

Department: Computer Science & Engineering and Information Technology Dated:

# CANDIDATE'S DECLARATION

I hereby declare that the work presented in this report entitled "**Video Object Segmentation for Object Detection and Recognition**" in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science & Engineering / Information Technology** submitted in the Department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology, Waknaghat is an authentic record of my own work carried out over a period from August 2023 to December 2023 under the supervision of **Dr. Vipul Kumar Sharma** (Assistant Professor (SG), Department of Computer Science & Engineering and Information Technology). The matter embodied in the report has not been submitted for the award of any other degree or diploma.

(Student Signature with Date)

Student Name: Ansh Mahajan

Roll No.: 201137

(Student Signature with Date)

Student Name: Snehan Tagnaat

Roll No.: 201541

This is to certify that the above statement made by the candidates is true to the best of our knowledge.

(Supervisor Signature with Date)

Supervisor Name: Dr. Vipul Kumar Sharma

Designation: Assistant Professor (SG)

Department: Computer Science & Engineering and Information Technology

Dated:

# ACKNOWLEDGEMENT

We extend our sincere thanks to the Almighty for His divine blessings, which have enabled us to successfully complete this project. Our deepest gratitude goes to our supervisor, Dr. Vipul Sharma, Assistant Professor (SG), Department of Computer Science and Engineering, Jaypee University of Information Technology, Waknaghat. Dr. Vipul Kumar Sharma's profound knowledge and keen interest in the field of machine/deep learning have been instrumental in the successful execution of this project. His unwavering patience, scholarly guidance, continuous encouragement, energetic supervision, constructive criticism, valuable advice, meticulous review of numerous drafts, and corrections at all stages have played a pivotal role in bringing this project to fruition.

We express heartfelt thanks to Dr. Vipul Kumar Sharma for his generous assistance in bringing our project to completion. Our warm appreciation extends to all individuals who have directly or indirectly contributed to the success of this project. In this context, we would like to acknowledge the various staff members, both teaching and non-teaching, who provided timely assistance and facilitated our project.

Lastly, we acknowledge with profound respect the unwavering support and patience of our parents throughout this endeavor."

Ansh Mahajan

(201137)

Snehan Tagnaat

(201541)

# TABLE OF CONTENT

| <b>Title</b>  | <b>Page No.</b> |
|---|-----------------|
| <b>Certificate</b>                                  | <b>i</b>        |
| <b>Candidate's Declaration</b>                      | <b>ii</b>       |
| <b>Acknowledgement</b>                              | <b>iii</b>      |
| <b>Table of Content</b>                             | <b>iv</b>       |
| <b>List of Figures</b>                              | <b>vii</b>      |
| <b>List of Abbreviations</b>                        | <b>viii</b>     |
| <b>Abstract</b>                                     | <b>ix</b>       |
| <b>Chapter 1: Introduction</b>                      | <b>1</b>        |
| 1.1 Introduction                                    | 1               |
| 1.2 Problem Statement                               | 3               |
| 1.3 Objectives                                      | 4               |
| 1.4 Significance and Motivation of the Project Work | 4               |
| 1.5 Organization of Project Report                  | 5               |
| <b>Chapter 2: Literature Survey</b>                 | <b>7</b>        |
| 2.1 Semi Supervised Video Object Segmentation       | 7               |
| 2.2 Unsupervised Video Object Segmentation (VOS)    | 9               |
| 2.3 Over Segmentation                               | 11              |
| 2.4 Proposal-based segmentation                     | 12              |
| 2.5 Motion Segmentation                             | 12              |
| 2.6 Semi-automatic Video Segmentation               | 13              |
| 2.7 Graph based Video Segmentation                  | 13              |
| 2.8 Interactive Video Segmentation                  | 14              |
| <b>Chapter 3: System Development</b>                | <b>16</b>       |
| 3.1 Requirements and Analysis                       | 16              |
| 3.2 Project Design and Architecture                 | 17              |
| 3.3 Data Preparation                                | 18              |
| 3.4 Implementation                                  | 19              |
| 3.5 Key Challenges                                  | 31              |

|  |           |
|--|-----------|
| <b>Chapter 4: Results and Evaluation</b>       | <b>32</b> |
| 4.1 Datasets                                   | 32        |
| 4.2 Experiments                                | 34        |
| <b>Chapter 5: Conclusions and Future Scope</b> | <b>40</b> |
| 5.1 Conclusion                                 | 41        |
| 5.2 Applications                               | 42        |
| 5.3 Future Scope                               | 34        |
| <b>References</b>                              | <b>43</b> |

# LIST OF FIGURES

| <b>S. No.</b> | <b>Figure Name</b>                            | <b>Page NO.</b> |
|---------------|---|-----------------|
| 1             | Figure 1.1: shows an example of our technique | 2               |
| 2             | Figure 1.2: Overview of OSVOS                 | 3               |
| 3             | Figure 2.1: Types of Input Annotations        | 8               |
| 4             | Fig 2.2: Motion segmentation results          | 12              |
| 5             | Fig: 3.1: Overview of Video Segmentation      | 18              |
| 6             | Fig:3.2: Code Snippets                        | 19              |
| 7             | Fig :3.3: Code Snippets                       | 22              |
| 8             | Fig:3.4: Code Snippets                        | 24              |
| 9             | Fig:3.5: Code Snippets                        | 25              |
| 10            | Fig:3.6: Code Snippets                        | 26              |
| 11            | Fig:3.7: Code Snippets                        | 27              |
| 12            | Fig:3.8: Code Snippets                        | 28              |
| 13            | Fig:3.9: Code Snippets                        | 29              |
| 14            | Fig:3.10: Code Snippets                       | 30              |
| 15            | Fig:4.1: Result of our task                   | 34              |
| 16            | Fig:4.2: Size of DAVIS 2017 Unsupervised vs.  | 36              |

|    |  |    |
|----|--|----|
|    | DAVIS 2016.  |    |
| 17 | Fig.4.3: Davis Validation  | 36 |
| 18 | Fig:4.4: Measures for Unsupervised Algorithm                           | 37 |
| 19 | Fig:4.5: Davis Validation  | 37 |
| 20 | Fig 4.6: Region similarity w.r.t processing time unit per frame        | 37 |
| 21 | Fig:5.1: Car dash devices having functionalities of video segmentation | 41 |



# LIST OF ABBREVIATIONS

| Abbreviation | Full Form                             |
|--------------|---------------------------------------|
| CNNs         | Convolutional Neural Networks         |
| FCNs         | Fully convolutional networks          |
| OSVOS        | One-Shot Video Object Segmentation    |
| RVOS         | Referencing Video Object Segmentation |
| VOS          | Video Object Segmentation             |
| DAVIS        | Densely Annotated Video Segmentation  |
| SAT          | State Aware Tracker                   |
| DAG          | Directed Acyclic Graph                |
| GIS          | Guided interactive segmentation       |
| ML           | Machine Learning                      |
| APIs         | Application Programming Interface     |

# ABSTRACT

Segmenting and tracking items of interest in films is the goal of a crucial computer vision problem called video object segmentation (Video-object Segmentation). Due to its numerous uses in a variety of industries, including image monitoring, video editing, autonomous cars, and human-computer interaction, it has garnered a lot of attention. Video-object Segmentation presents additional issues, such as controlling temporal changes and preserving consistent object representation across frames, in contrast to image segmentation, which deals with static images.

The two primary approaches to Video-object Segmentation technology are supervised and unsupervised. The unsupervised system attempts to partition objects without depending on a precise description, whereas the supervised system uses learning data to understand the relationship between pixels and objects. Both approaches have benefits and drawbacks. While unsupervised approaches are more flexible, they may run into challenging issues. Supervised approaches often achieve more accuracy but need more registration data.

Current advancements in deep learning Video-object Segmentation have transformed deep learning, boosting robustness and efficiency. The primary models of Video-object Segmentation nowadays are convolutional neural networks (CNN), which can streamline object segmentation and extract intricate characteristics from photos. Short-term temporal (LSTM) networks and neural networks (RNN) combined have also drawn interest, particularly for applications that need to describe object and physical trajectories.

Even with all of the advancements, Video-object Segmentation continues to face numerous obstacles that keep it from being widely used in real-world applications. Complex scenarios with lots of interactive items, occlusions, and background clutter present one of the biggest obstacles. Furthermore, Video-object Segmentation algorithms frequently struggle to adapt focus, illumination, and image quality. Furthermore, those can be a restriction on the Video-object Segmentation algorithm's instant use.

At Video-object Segmentation, ongoing research is focused on creating more potent and efficient algorithms that can handle challenging issues quickly and efficiently in order to address these issues. Examining the tracking and self-tracking procedure is a viable strategy that can help the Video-object Segmentation algorithm focus on pertinent regions of the video and discern between objects and the

background. To increase segmentation accuracy and lower overhead, researchers are also attempting to incorporate domain expertise and previous information into the Video-object Segmentation architecture.

As Video-object Segmentation technology continues to evolve it should play an important role in many applications. In video editing, Video-object Segmentation makes it easy to automatically select, monitor and adjust components for fine-tuning and mixing. In visual monitoring, Video-object Segmentation can be used for real-time object detection, tracking, and detection of anomalies to improve safety and maintenance capabilities. In unmanned vehicles, Video-object Segmentation can help drive safer and more reliable by providing critical information for situational understanding and configuration. In human-computer interaction, Video-object Segmentation enables interactive device interaction and gesture recognition, thus improving user experience and encouraging interaction with the computer.

The way we engage with and comprehend video is expected to be completely transformed by Video-object Segmentation in the future. Video-object Segmentation algorithms will grow more sophisticated, accurate, and efficient as deep learning and artificial intelligence continue to progress. This will open up a wide range of new applications and transform the digital landscape.

# CHAPTER 1: INTRODUCTION

## 1.1 Introduction

The objective of video object segmentation (VOS) is to generate precise and accurate segmentation of a particular object instance throughout a video input. This has numerous practical uses in the areas of video comprehension and editing.

Image segmentation has entered a new phase as a result of deep learning notable performance over the past several years in visual identification tasks. Deep neural networks significantly boost performance and frequently attain the greatest accuracy rates on known benchmarks. Many segmentation techniques based on deep neural networks have surfaced in recent literature. The earliest application of convolution neural networks (CNNs) were for classification problems. Fully convolutional networks (FCNs) are currently the most complex structures used in picture segmentation.

Object appearance and motion are crucial cues to perform this task. However, there are still several difficulties: Opposite sections of the object may move in different directions, and certain objects may resemble the background. Because of this, a lot of techniques still rely on supervision, at least when learning to extract visual features. In order to classify all of the pixels, One-Shot Video Object Segmentation (OSVOS) is a CNN architecture that separates the foreground and background in a video sequence based on manual annotation provided for one or more of its frames. This technique has various applications in video analysis and editing., was developed.

One of the most important tasks in computer vision is video segmentation, which is dividing a video sequence into segments that have semantic significance. The aforementioned procedure facilitates the recognition and monitoring of distinct entities or regions of interest within a changing visual environment. Segmentation algorithms are essential for a wide range of applications, from object recognition and autonomous systems to video editing and content analysis, because they explore the temporal and spatial properties of video frames. To get accurate and real-time segmentation findings, the difficulty is to build algorithms that can robustly handle a variety of circumstances, such as complicated motion patterns, occlusions, and changes in object appearance.

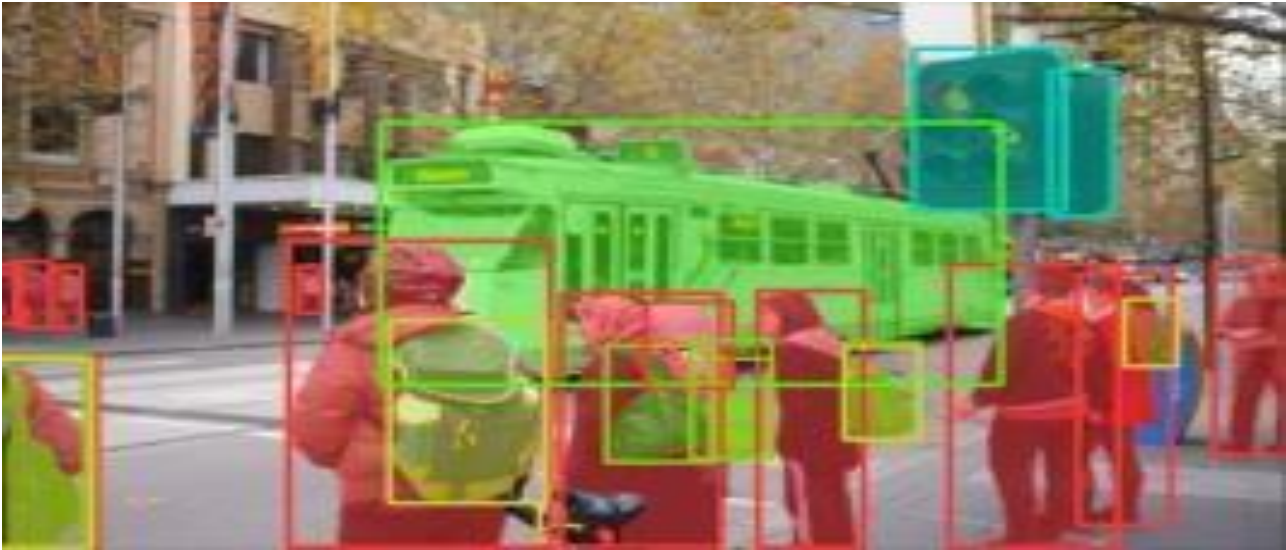


Fig: 1.1: shows an example of our technique

The initial idea was to alter the CNN of a specific object instance by providing a single annotated image—hence the term "one-shot." In order to do this, a CNN trained in image recognition was transformed into a segmentation tool for video objects. This is accomplished by watching a number of movies that have had their segments manually divided. It is manually wrapped in a particular object and divided into a single image for testing. Rather than using expensive and draconian limitations, OSVOS achieves temporal consistency by treating each frame of a movie individually. Put another way, we express the video segmentation problem as a frame-by-frame segmentation problem using an object model made up of one or more manually segmented frames. Motion estimation is becoming a crucial component of modern video segmentation systems. Their exploitation is not simple yet, since one must compute temporal coincidences, such as those that take the shape of dense or optical flow. Different levels of accuracy and speed adjustment are available for OSVOS.



Figure 1.2: Overview of OSVOS

We argue that because the out-of-date shape or appearance models had significant flaws, temporal consistency was necessary in the past. However, it will be shown that deep learning can provide a model of the target object that is precise enough to yield temporally reliable findings even while processing individual frames. This offers a few clear advantages: While OSVOS is able to segment objects over occlusions, it is limited to specific motion ranges, it doesn't need frame-by-frame processing, and mistakes aren't conveyed throughout time.

A technique called referencing Video Object Segmentation (RVOS) tries to separate target objects from a video sequence using referencing expressions from natural language. RVOS can recognize the target based only on an abstract language query, in contrast to semi-supervised Video Object Segmentation (VOS) from 2016, which needs a per-pixel mask to initialize the target location. The community has paid close attention to this method since it offers a more practical choice for human-computer interaction (Khoreva, Rohrbach, and Schiele 2018; Seo, Lee, and Kim 2020). However, because RVOS calls for simultaneous interpretation of both language and visual modalities, it is also more difficult.

## 1.2 Problem Statement

Early VOS systems generally used the abjectness, optical flow, and visual saliency techniques based on hand-crafted properties to separate objects from video sequences. Even though deep learning and high-performance computing have evolved since these methods were initially created, they continue to yield cutting-edge outcomes. Techniques for VOS based on deep learning have become much more accurate and efficient.

Because of this, deep neural networks are used to implement the majority of VOS techniques now in use. Statistics from two reliable VOS benchmarks show that while the performance of current VOS techniques is increasing yearly, it is not yet at saturation. Deep learning-based VOS is the current focus of computer vision research because of its potential applications and room for performance enhancement.

### 1.3 Objectives

The following are the project's goals:

- to deliver more accuracy than previously proposed study publications.
- to remove the background object from the input that was acquired in frames.
- to create a model that tracks an object in a video while breaking up the input image.
- Multiple objects, including hidden levels, may need to be tracked. Multiple layers may also exist.
- It will be shown that even with independent processing of each frame, deep learning is still capable of creating a temporally stable model of the target object.

### 1.4 Significance and Motivation of the Project Work

#### 1.4.1 Significance:

A crucial tool in computer vision, Video Object Segmentation (VOS) holds revolutionary promise for a wide range of uses.

**1. Transforming Video Editing:** By automating product monitoring, selection, and control, VOS transforms video editing. This makes work easier, fosters innovation, and makes integration and special advantages easier.

**2. Enhancing Visual Surveillance:** VOS facilitates instantaneous object identification, tracking, and vulnerability identification in surveillance. This strengthens situational awareness, promotes problem solving, and increases security and surveillance capabilities.

**3. Empowering Autonomous Vehicles:** The basis of self-driving is accurate product identification and segmentation. An improved situational awareness, obstacle detection, and lane following are all made possible by VOS, which opens the door to a safer and more dependable driving experience.

**4. Revolutionizing Human-Computer Interaction (HCI):** VOS facilitates communication between users and computers. It facilitates gestures.

#### 1.4.2 Comparative Advantages of VOS:

**1. Efficient and accurate:** VOS algorithm saves a lot of time and resources by performing better than the segmentation procedure, particularly for huge video files.

**2. Real-time functionality:** Applications can be utilized in numerous contexts where real-time analysis is essential thanks to the VOS standard, which makes real-time functionality possible. Scalability and adaptability: The VOS architecture offers scalable and versatile solutions that can be readily expanded to manage massive files and cater to various photo and application kinds.

### **1.4.3 Motivation**

Researchers are motivated by the following objectives as they push the limits of VOS technology:

- Develop more potent and precise algorithms by addressing issues with scenario complexity, illumination changes, occlusion, and blurring to attain excellent accuracy and performance.
- **Boost performance in real time:** Real-time algorithm optimization can lead to effective data processing with minimal overhead and latency.
- **Find fresh uses for:** Expand the use of VOS beyond its pre-installed settings, find new uses for it, and realize all of its possibilities.
- **Help progress computer vision:** By creating new VOS technology that makes it easier to comprehend and analyze visual content more deeply, computer vision is being advanced.

## **1.5 Organization of Project Report**

This project report is made in a way to provide a complete understanding of the Video Object Segmentation project, discussing its objectives, methodology, implementation, challenges, results, and future scope. The report is organized into six chapters, each dealing with a specific aspect of the project:

### **Chapter 1: Introduction**

This chapter deals with introduction to the Video Object Segmentation project, providing a brief overview of its purpose, significance, and motivation. It highlights the objectives and problem statement of the project.

### **Chapter 2: Literature Survey**

This chapter presents an overview of relevant literature and key gaps in the literature studied during the making of Video Object Segmentation, covering existing research and developments in the field of chess boards and chess pieces detection.



### **Chapter 3: System Development**

This chapter delves into the requirements and analysis phase of the project. It gives a detailed description of Project Design and Architecture, describes the Data preparation process and Implementation of the project along with the screenshots of code snippets and finally Key Challenges are addressed that were faced during the project.

### **Chapter 4: Testing**

This chapter describes the project design and architecture, explaining the overall structure and technical components of the Video Object Segmentation application. It outlines the software architecture, user interface design, and data management strategies.

### **Chapter 5: Results and Evaluation**

This chapter discusses the key findings of the project and their interpretation along with a snapshot of the results. Finally, a comparison is performed with the existing solutions.

### **Chapter 6: Conclusion and Future Scope**

This chapter concludes the project report by summarizing the key achievements, contributions and limitations of the Video Object Segmentation project. It also discusses potential future scope and enhancements for the application.

# CHAPTER 2: LITERATURE SURVEY

The internet in today's world has been overflowing with the enormous amount of textual form of data which is growing rapidly every minute. It has become very difficult to extract the exact information about a particular entity. As the Internet has grown in popularity, the need for individualized information systems has grown as well. The humongous amount of data has passed the limits of human capacity to search, organize and categorize it. In the past few years, there has been an evolution of the process of opinion gathering of the customers and the users. There are several websites, applications and even social media platforms which are gathering their users' reviews, likings and disliking. Different reviews contain different expressions, views and emotions which are hard to be categorized manually.

Various research papers had been proposed in the video object segmentation field. Every research paper and author deal with distinct approaches. There are ample resources available on the internet to conduct our own research and perform different experiments subsequently improving the accuracy of the model. Image recognition or object detection is the sub-field of Artificial Intelligence, algorithms used such as Convolutional Neural Networks (CNN), Fully Convolutional Network. Majorly there are four approaches used widely in this field such as Supervised learning, Unsupervised learning, Semi-Supervised learning, Zero-Shot video object segmentation or One-Shot video object segmentation (OSVOS). Dataset that we have used in this project is the Davis2017 dataset.

Several researchers have conducted their studies on such recommendation systems and have proposed some models using various algorithms and methodologies. Following are the related research papers accepted in this field of technology. They are categorized as different learning algorithms.

## 2.1 SEMI SUPERVISED VIDEO OBJECT SEGMENTATION

In order to track specified objects in films, the study "Tackling Background Distraction in Video Object Segmentation (2023)" suggests a semi-supervised video object segmentation (VOS) approach. The presence of background distractions that visually resembling the target objects presents the task's principal difficulty. The study proposes three innovative approaches to deal with this problem:

- A learnable distance scoring swap-and-attach augmentation, which ensures unique features for each object by providing training samples with entangled objects, is a spatiotemporally diversified template construction scheme that generates generalized properties of the target objects. function that uses temporal

consistency between two consecutive frames to exclude spatially-distant distractors. The proposed model achieves real-time performance and outcomes on publicly available benchmark datasets that are on par with state-of-the-art contemporary methods. The framework segments the frames in a video sequence using the ground truth segmentation mask that is shown in the first frame. Masks are forecasted using the feature similarity of the embedded features. A spatiotemporally variable template generation approach is employed to generate distinct template features for feature matching. The decoder receives as inputs a down sampled prior adjacent frame mask, low-level features from the encoder, and the outputs of feature matching.



Figure 2.1: Types of Input Annotations

The paper named "Accelerating Video Object Segmentation with Compressed Video" offers a powerful and adaptable acceleration framework for semi-supervised video object segmentation by taking advantage of the temporal redundancies in compressed movies. The proposed system transfers segmentation masks from keyframes to other frames frequently and bidirectionally using a motion vector-based warping technique.

- Additionally, the authors provide a residual-based correction module that fixes segmentation masks that were propagated from noisy or inaccurate motion vectors in an inappropriate way. The framework is flexible and may be applied to several existing video object segmentation techniques. Tests conducted on the DAVIS17 and You Tube-VOS datasets yielded very competitive results, with just minor accuracy losses and a significant speed-up of up to 3.5X.

- In the article "State-Aware Tracker for Real-Time Video Object Segmentation," the authors address the challenges associated with semi-supervised video object segmentation (VOS) and look into practical approaches that leverage video properties to overcome these difficulties. The authors recommend adopting the StateAware Tracker (SAT) pipeline, which can provide accurate segmentation results instantly. SAT treats each target object as a tracklet and uses inter-frame consistency to boost efficiency.

SAT self-adapts to every state via two feedback loops to improve the stability and robustness of the method. While the other loop helps SAT construct a more robust and complete target representation, the first loop helps SAT produce more stable tracklets. Using the DAVIS2017-Val dataset, the authors' results of 72.3% J.& F provide an encouraging mean of 39 frames per second this demonstrates a respectable balance between precision and efficiency.

- The study "A Transductive Approach for Video Object Segmentation" proposes a transductive method for semi-supervised video object segmentation that uses the mask in the first frame to isolate a target item from a video series. The label propagation method developed by the authors assigns labels to pixels based on how similar their features are in an embedding space. Their method distributes temporal data in a thorough way that takes object appearance over time into account. Unlike popular approaches, they do not require new modules, databases, or architectural designs. Furthermore, their technique runs quickly at 37 frames per second with no computing overhead. The single model with a vanilla ResNet50 backbone received a score of 72.3% on the DAVIS 2017 validation set.

## **2.2 Unsupervised Video Object Segmentation (VOS)**

Salient object identification is expanded to films using unsupervised video object segmentation techniques. They don't require manual annotation and make no assumptions about the segmentation target item. They often operate under the presumption that an object's motion is distinct from its surroundings, or salient motion. To achieve this, locate the object using a saliency detector, and compute the likelihood that a super pixel in the image belongs to the foreground object using the geodesic between two super pixels on the image. Instead, improve salient object detection by connecting all the video frames in a Markov chain. Some techniques, in addition to employing saliency, are based on object proposals and produce a number of ranked segmentations. Unsupervised methods are excellent for processing huge databases since they are restricted by the validity of their underlying assumptions. Although the issue of video object segmentation is the focus of this thesis, unsupervised methods have historically focused on over-segmentation or motion segmentation. As a result, the following paragraphs will provide a quick overview of these various domains.

According to this research study, unsupervised video object segmentation aims to distinguish a target object from the video's backdrop in the absence of a ground truth mask in the first frame. The hardest part of this task is extracting characteristics from the most noticeable common objects in the video stream. This issue can be solved by employing motion information, such as optical flow, although doing so causes

poor connection and performance between distant frames when relying solely on information from close frames. Target object segmentation in unsupervised video aims to be accomplished without using a ground truth mask in the initial frame of the video. The challenging task at hand involves identifying the most noticeable and often occurring objects in a video clip. This problem can be overcome by employing motion information such as optical flow, however using only the information between close frames results in poor connection and performance between distant frames. We present a unique prototype memory network design to address this issue. The suggested approach efficiently recovers RGB and motion information from input RGB pictures and optical flow maps by generating super pixel-based component prototypes.

In this research, we suggest a new method for segmenting objects in videos without supervision. Our technique identifies conspicuous objects in the video and automatically creates instance-level segmentation masks for them. We tackle the issues with current techniques, including drift during temporal propagation, tracking, and the inclusion of additional objects. We present the notion of online mask improvement with an ensemble of criteria that evaluates mask quality. We also provide Selector Net, a neural network trained to generalize across several datasets and designed to evaluate mask quality. With a J&F mean of 61.6%, our suggested approach delivers state-of-the-art performance on the Davis 2019 Unsupervised Challenge dataset by limiting the noise accumulated along the movie. Using datasets like FBMS and SegTrack V2, we further evaluated our approach and discovered that it outperformed or was comparable to other approaches.

The process of identifying and tracking important items in a movie without a predetermined definition is known as unsupervised video object segmentation. Prior research has concentrated on extracting the foreground and background from videos using techniques including marker-based segmentation, background subtraction, and object suggestions. These techniques are susceptible to shadows and not resilient enough to withstand small variations in lighting. These problems have been addressed by deep learning techniques, with the Davis 2016 dataset serving as a popular benchmark. Nevertheless, these methods provide a single binary mask for every foreground object and are unable to manage scenarios with many foreground objects or address issues with object re-identification, tracking, and addressing occlusion.

Some methods, including single foreground mask prediction and multi-moving foreground object segmentation and tracking, have concentrated on explicitly extracting moving items as foreground

objects. However, because these approaches only consider moving foreground items, they cannot be easily linked with multi-object segmentation and tracking.

To create the first object masks in our study, we trained a Mask R-CNN implementation on the COCO dataset using a ResNet-50 backbone. In order to separate objects outside of the categories that Mask R-CNN is trained on, we set the confidence score threshold to 0.1. We chose a maximum of 10 objects, ranking them based on their confidence score, in order to restrict the number of objects in a frame.

## **2.3 Over Segmentation**

The most prevalent techniques to region segmentation are based on intensity thresholding and perform well for photos containing homogeneous objects of interest. However, many photographs feature noise, texture, and clutter, all of which reduce the usefulness of these approaches. The use of threshold-based segmentation algorithms on pictures containing nonhomogeneous objects of interest might result in segmentation that is either coarse or too fine. These outcomes are referred to as under segmentation and over segmentation, respectively. Split and merge approaches are frequently employed to successfully resolve these issues.

Setting segmentation process settings, such as a threshold value, such that all objects of interest are recovered from the backdrop or each other without over segmenting the data is not achievable for some photos. Over segmentation is the process of segmenting or fracturing the items being segmented from the backdrop into subcomponents.

Over segmentation increases the likelihood that important borders have been removed at the expense of establishing numerous inconsequential barriers. Prefiltering techniques, as addressed in earlier columns see *Vision Systems Design*, Oct. 1998, p. 20, should be employed in this scenario to try to reduce noise, increase intro etc. definition, or smooth picture textures, all of which may create segmentation issues.

Super voxel-based techniques can be used to deal with unconstrained motion. These techniques result in an over segmentation of the video into perceptually distinct, space-time homogenous sections. They are crucial for early video preprocessing, but they don't directly address the issue of video object segmentation since they don't offer a sound strategy for flattening the video's hierarchical decomposition into a binary segmentation.

## 2.4 Proposal-based segmentation

The use of object suggestions in video object segmentation has been prompted by recent developments in cutting-edge image analysis. find key-segment clusters in films that connect the concepts of objectness' and similarity in appearance. The top-scoring hypothesis is then automatically chosen for video segmentation after being ranked later. Their research is useful for identifying collections of segments with a common appearance and motion, but it ignores the relationships in space and time between segments.

Finding the largest weighted clique in a locally linked graph with mutex constraints is one way to phrase the issue. However, their usefulness in real-world contexts is constrained by the rigid presumptions that the object must present in every frame. develop a layered Directed Acyclic Graph (DAG) using pair wise comparisons and unary edges to measure the objectness of the proposed object.

## 2.5 Motion Segmentation

Realizing the independently moving devices (pixels) in a video and separating them from the historical movement is the aim of movement segmentation. We can also use projective differences to effectively sign up a few frames onto an unmarried frame if the backdrop is a plane.

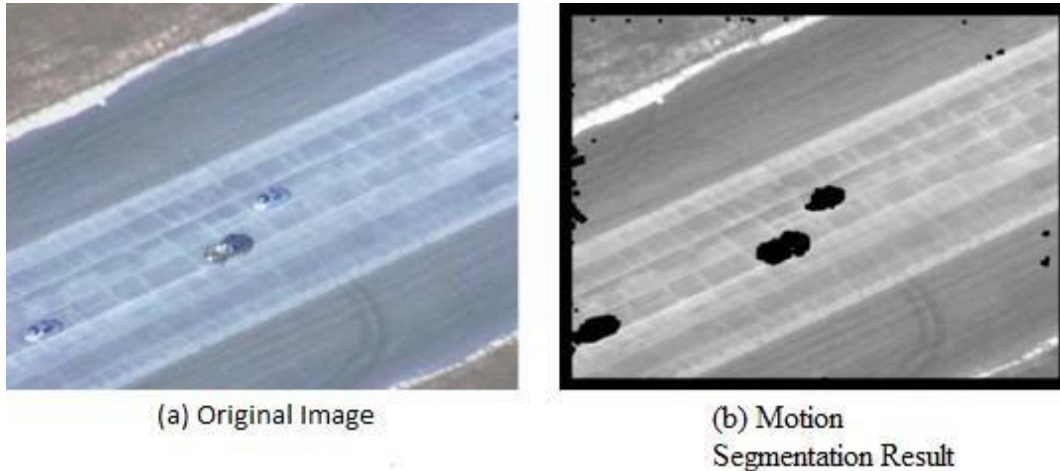


Fig 2.2: Motion segmentation results

The shifting gadgets are liable for the elements of the photo that doesn't sign up effectively. We can take the photo distinction of registered if the registration of all frames is correct. Moving gadgets may be recognized with the aid of using pixels with a good size depth distinction. However, due to the fact registration isn't always constantly flawless, this primary method generates quite a few fake alerts.

Methods that track critical points through time and, more recently, over image regions, have made significant strides. However, these techniques only take into account the last two frames of the videos and are sensitive to quick changes in motion and appearance. Propose a method for segmenting motion in relation to tracking systems by spectrally grouping long term point trajectories based on their motion affinity and using a variational method to transform the resulting sparse trajectories clusters into dense region. They presuppose a translational motion model by defining the pairwise distance between trajectories as the greatest difference of their motion, this is a reasonable approximation for spatially close point trajectories, but it is challenging to segment articulated bodies after non-rigid motion using these methods.

## **2.6 Semi-automatic Video Segmentation**

It is split into steps: intra-body segmentation and interframe segmentation. To begin, intra-body segmentation is carried out to the preliminary body of the image collection or to frames containing simply newly emerged video gadgets or scene changes. The newly rising video gadgets withinside the pictures are manually described or segmented via way of means of the user. Then, following the primary body or a body with a newly regarded item or scene change, inter-body segmentation is carried out to the following frames. Object monitoring is used to robotically section user-described video gadgets at some stage in inter-body segmentation.

Sparse hand labeling is applied throughout the whole video stream by semi-automatic video object segmentation techniques; this is usually done in the form of one or more annotated frames. They often tackle an optimization problem using an energy defined throughout a network structure, notwithstanding their variances. Semi-automatic segmentation and object tracking go hand in hand. While the goal of tracking is to define the object's borders within a rectangular bounding box, the goal of video segmentation is to do so as precisely as possible.

## **2.7 Graph based Video Segmentation**

In general, graph-primarily based totally picture segmentation tactics paint the difficulty as a graph,  $G = (V, E)$ , in which every node at  $V$  corresponds to pixels withinside the photograph and the rims join specific pairs of close by pixels.

Each area is assigned a weight relying on a few characteristics of the pixels it links, along with their photograph intensities. Depending on the method, every pair of vertices might also additionally or might



not have an area linking them. The first graph-primarily based total algorithms compute segmentation the use of preset thresholds and neighborhood metrics. Zahn [19] affords a segmentation method primarily based totally at the graph's minimum spanning tree. This method has been used for factor clustering in addition to photograph segmentation.

It is evident that a daily graph structure, with edges connecting adjacent pixels in a spatial or spatiotemporal arrangement, is well suited for images and videos. It is then possible to construct video segmentation as an optimization problem that seeks to stabilize a coherent label task of contiguous vertices while adhering to a predefined item version or user constraints.

## **2.8 Interactive Video Segmentation**

During the segmentation process, supervised approaches presuppose that manual annotation will be continuously added, and the algorithm results will be iteratively corrected by a human. To prevent overwriting earlier human fixes, these systems often operate online with forward processing frames. They are therefore well suited for particular situations, such as video editing, because they guarantee high segmentation quality at the cost of a higher level of human supervision. In post-production, scene segmentation is regarded with the term rotoscoping. The task is also very expensive and time-consuming. As a result, a substantial body of research has examined this issue in an effort to minimize the amount of human labor needed to provide high quality results.

In this paper, a system for interactive Video object segmentation (VOS) in the real world is proposed, where users can iteratively select specific frames for annotations. The masks are then improved using a segmentation algorithm using the user annotations. The prior interactive VOS paradigm chooses the frame with some of the worst evaluation metrics, and since the assessment measure must be calculated using the ground truth.

It is not feasible during the testing phase. Contrarily, we argue in this research that the frame with the worst assessment metric might not necessarily be the most valuable frame that improves performance throughout the film.

We provide a singular guided interactive segmentation (GIS) approach for video objects that aims to maximize segmentation accuracy while reducing interplay time. First, we develop the dependability-based interest module, which examines the dependability of multiple annotated frames. Secondly, we provide an intersection-aware propagation module that allows segmentation results to be transmitted to

neighboring frames. Third, we extend a GIS methodology that enables an individual to easily and swiftly select undesirable frames. The results of the experiment show that the suggested set of rules outperforms traditional algorithms in terms of accuracy and speed of segmentation results.

# Chapter 3: System Development

## 3.1 Requirements and Analysis

One important task in computer vision is video object segmentation, which is precisely identifying and tracking objects over a series of successive frames in a video. Video object segmentation algorithms play a critical role in many applications such as surveillance, driverless vehicles, and video editing. This section covers the requirements and analysis framework that guide the development and evaluation of video object segmentation techniques.

### 3.1.1 Data Requirements

- **DAVIS Dataset:** DAVIS is part of the DAVIS Challenges, which are designed to promote advances in video object segmentation and are coordinated by the Computer Vision Lab at ETH Zurich.
- DAVIS measures the performance of segmentation algorithms quantitatively using conventional evaluation metrics for video object segmentation, such as F-measure and Jaccard Index (Intersection over Union, or IoU).

### 3.1.2 Integration and Deployment Requirements

- **Compatibility:** Make sure it works with popular frameworks for object detection and recognition (e.g., TensorFlow, PyTorch, OpenCV).
- Support standard interfaces to enable compatibility with external systems.
- **API Assistance:** Offer an Application Programming Interface (API) that is thoroughly documented in order to make interaction with services and applications from other parties easier.

### 3.1.3 Requirement Specification

- **Input for Videos:** It is imperative that the system has the ability to receive video information from multiple sources, such as pre-recorded recordings and live camera feeds. Standard video formats like MP4, AVI, and others need to be supported.
- **Segmenting Objects:** Pixel-level segmentation ought to be used by the system to precisely identify items. It is necessary to support both instance segmentation and semantic segmentation techniques.
- **Consistency in Time:** Ensure that there are seamless transitions between the segments in the output to provide temporal consistency. We need to Address issues relating to dynamic shifts in the appearance and motion of objects.

### 3.1.3.1 Hardware Requirements

- Processor: i3 /i5/ i7 Intel Core 1.2 GHz or better
- RAM: 4 GB
- HDD: 10 GB

### 3.1.3.2 Software Requirements

- Operating System: Windows 10/11
- IDEs: Google Colab
- Programming Languages: Python
- Frameworks and libraries: OpenCV, YOLO (You Only Look Once), TensorFlow.

### 3.1.3.3 Non-functional Requirements

- **Precision and Accuracy:** Reach a high level of object segmentation precision and accuracy to provide reliable identification and identification. Cut down on the number of false positives and false negatives in the segmented output.
- **Flexibility:** Provide a system that can grow to accommodate different video resolutions, frame rates, and durations. Make sure that videos with varying degrees of complexity perform consistently.
- **Sturdiness:** Adapt to difficult situations such as changing illumination, intricate backgrounds, and occlusions. Put in place procedures for mistake recovery and graceful degradation under difficult circumstances.
- **Use of Resources:** Optimize memory and processing power consumption to enable deployment in contexts with limited resources.

## 3.2 Project Design and Architecture

### 3.2.1 Overview of the Project:

- **Goal:** Create a system that can precisely identify and segment objects in videos in situations where it's either real-time or almost real-time.
- The system's scope includes receiving video input, segmenting objects, and identifying and tracking objects inside the video frames.

### 3.2.2 Architecture of the System:

- **Parts:** Video Input Module, Object Segmentation Module, Object Detection and Tracking Module, Output Visualization Module.

- **Interactions:** The Object Segmentation Module receives video frames from the Video segmentation.
- Below is the representation of video segmentation system:

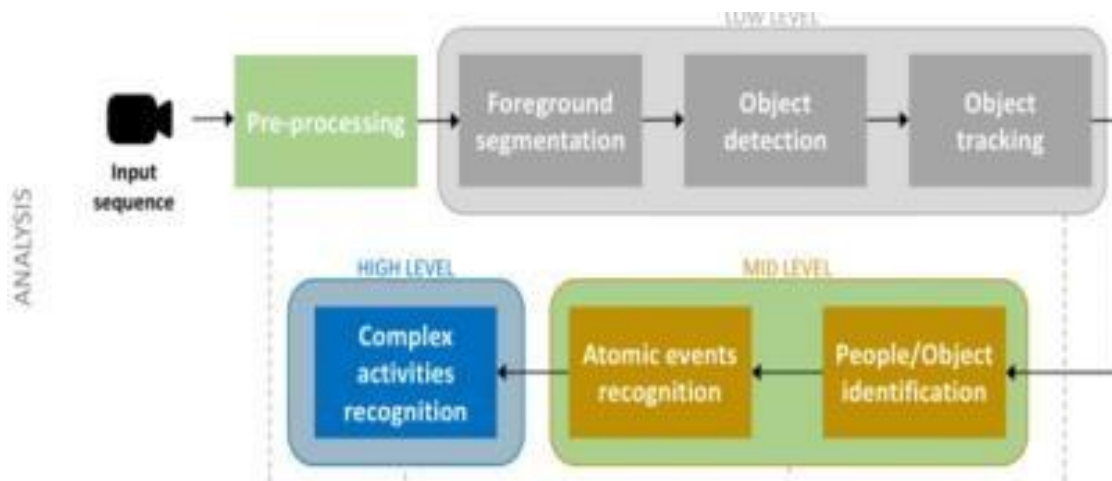


Fig: 3.1: Overview of Video Segmentation

### 3.3 Data Preparation

For the successful development of the Video Object Segmentation for Object Detection and Recognition, it is crucial to prepare our data for video object segmentation and detection, gather, arrange, and preprocess the information required for model training and assessment. For this work, the following is a step-by-step tutorial on data preparation:

#### 3.3.1 Define Object Classes

Choose the categories of things in our movies that you wish to recognize and categorize. Specify the classes that our model will be trained to identify and partition.

#### 3.3.2 Data Collection

Put together a diverse selection of films that feature instances of the designated object classes. Consider including movies with different backdrops, lighting configurations, and item placements.

#### 3.3.3 Annotation

Use segmentation masks at the pixel level to describe each frame in the films. The zones that correspond to the objects of interest should be indicated by this annotation. For this task, programs like Labelbox, COCO Annotator, and VGG Image Annotator (VIA) can be helpful.

### 3.3.4 Dataset Splitting

Separate the training, validation, and testing sets from our dataset. To prevent data leaking, make sure that the same video does not appear in many splits.

### 3.3.5 Video Preprocessing

Resize and format videos to a standard size. For training purposes, standardize the frame rate in all of the videos to provide temporal consistency.

### 3.3.6 Frame Extraction

To make separate images, extract frames from films. Based on the features of our dataset and the specifications of our model, select an appropriate frame rate.

### 3.3.7 Data Augmentation

To make our training dataset more diverse, use data augmentation approaches. Random rotations, flips, brightness adjustments, and shifts are examples of common augmentations.

## 3.4 Implementation

```
from PIL import Image
import os
import numpy as np
import sys
class Dataset:
    def __init__(self, train_list, test_list, database_root, store_memory=True, data_aug=False):
        if not store_memory and data_aug:
            sys.stderr.write('Online data augmentation not supported when the data is not stored in memory!')
            sys.exit()
        data_aug_scales = [0.5, 0.8, 1]
        data_aug_flip = True
        print('Started loading files...')
        if not isinstance(train_list, list) and train_list is not None:
            with open(train_list) as t:
                train_paths = t.readlines()
        elif isinstance(train_list, list):
            train_paths = train_list
        else:
            train_paths = []
        if not isinstance(test_list, list) and test_list is not None:
            with open(test_list) as t:
                test_paths = t.readlines()
        elif isinstance(test_list, list):
            test_paths = test_list
        else:
            test_paths = []
        self.images_train = []
        self.images_train_path = []
        self.labels_train = []
        self.labels_train_path = []
```

Fig:3.2: Code Snippets

The code which is used defines a Python class called Dataset that facilitates the loading of image and label data for training and testing purposes. Let's break down the code into its components and explain each part in detail.

### **Importation:**

The code begins by importing the necessary libraries: PIL (Python Imaging Library) for image processing, so for interacting with the operating system, NumPy for numerical operations, and sys for system parameters and functions.

### **Class Definition:**

A Dataset class is defined to take care of loading the image and label data.

**Constructor (`__init__`):** The constructor initializes the Dataset object with parameters:

**train\_list:** List of paths to the file containing the training data.

**test\_list:** List of file paths containing test data.

**database\_root:** The root directory where the data is stored.

**store\_memory:** Boolean flag indicating whether data should be stored in memory or retrieved at runtime.

**data\_aug:** Boolean flag indicating whether to perform data augmentation. It also checks to see if online data extension without data storage is requested, which is not supported, and exits if attempted.

**Data extension parameters:** Defines parameters for expanding data, including scales and whether to flip images.

**Loading data:** It will start loading the files with a print statement indicating the start of the process. Parses the provided train and test lists, or initializes empty lists if none are provided. Initializes lists to hold training image data (`self.images_train`), training label data (`self.labels_train`), image file paths (`self.images_train_path`), and label file paths (`self.labels_train_path`).

### **Iterate over the training data paths:**

- Iterates over each row in the provided train routes.
- If the data is stored in memory, it reads the image and label using PIL from the specified paths and converts them to NumPy arrays of type uint8.
- Parses a label, presumably to extract its first element. This may need to be fixed as it is currently trying to split the PIL image object directly.
- If data stretching is enabled, it will stretch the data by resizing images and labels based on predefined scales and optionally flipping them horizontally.

- If the data is not stored in memory, it simply loads the image and label without augmentation.
- Prints progress points (.) after every 50 frames loaded.

The code retrieves images and labels, optionally with data extensions, and stores or processes them based on the specified parameters. The Dataset class provides functions to retrieve image and label data for training and testing, supports data expansion, and flexible data storage options. It is designed to handle large data sets efficiently while providing flexibility in use.

```

self.images_test = []
self.images_test_path = []
for idx, line in enumerate(test_paths):
    if store_memory:
        self.images_test.append(np.array(Image.open(os.path.join(database_root, str(line.split()[0]))),
                                                dtype=np.uint8))

        if (idx + 1) % 1000 == 0:
            print('Loaded ' + str(idx) + ' test images')
    self.images_test_path.append(os.path.join(database_root, str(line.split()[0])))
print('Done initializing Dataset')

```

Fig: 3.3: Code Snippets

This additional fragment extends the Dataset class to control the loading of the test image.

**Test data initialization:** Initialize the `self.images_test` and `self.images_test_path` lists to store the test image data and their corresponding file paths.

**Trial data pass iterations:** It is repeated in each row of the specified `test_path`. If `store_memory` is True, load the image using PIL, convert it to a NumPy array of type `uint8`, and add it to the `self.images_test` list. This process ensures that the test image is loaded into memory when specified. After loading every 1000 test images, it will print a progress message.

**Save test image path:** Add the loaded test image path to the `self.images_test_path` list.

**End of processing:** Prints a message indicating the completion of the dataset initialization process. This section complements the existing functional code by extending it to process the test images in a similar way to the training image processing method. This will load the test image and its path correctly and save it for later use.



```

self.train_ptr = 0
self.test_ptr = 0
self.train_size = max(len(self.images_train_path), len(self.images_train))
self.test_size = len(self.images_test_path)
self.train_idx = np.arange(self.train_size)
np.random.shuffle(self.train_idx)
self.store_memory = store_memory

def next_batch(self, batch_size, phase):
    if phase == 'train':
        if self.train_ptr + batch_size < self.train_size:
            idx = np.array(self.train_idx[self.train_ptr:self.train_ptr + batch_size])
            if self.store_memory:
                images = [self.images_train[l] for l in idx]
                labels = [self.labels_train[l] for l in idx]
            else:
                images = [self.images_train_path[l] for l in idx]
                labels = [self.labels_train_path[l] for l in idx]
            self.train_ptr += batch_size
        else:
            old_idx = np.array(self.train_idx[self.train_ptr:])
            np.random.shuffle(self.train_idx)
            new_ptr = (self.train_ptr + batch_size) % self.train_size
            idx = np.array(self.train_idx[:new_ptr])
            if self.store_memory:
                images_1 = [self.images_train[l] for l in old_idx]
                labels_1 = [self.labels_train[l] for l in old_idx]
                images_2 = [self.images_train[l] for l in idx]
                labels_2 = [self.labels_train[l] for l in idx]
            else:
                images_1 = [self.images_train_path[l] for l in old_idx]
                labels_1 = [self.labels_train_path[l] for l in old_idx]
                images_2 = [self.images_train_path[l] for l in idx]
                labels_2 = [self.labels_train_path[l] for l in idx]
            images = images_1 + images_2
            labels = labels_1 + labels_2
            self.train_ptr = new_ptr
        return images, labels
    elif phase == 'test':
        images = None
        if self.test_ptr + batch_size < self.test_size:
            if self.store_memory:
                images = self.images_test[self.test_ptr:self.test_ptr + batch_size]
                paths = self.images_test_path[self.test_ptr:self.test_ptr + batch_size]
            self.test_ptr += batch_size
        else:
            new_ptr = (self.test_ptr + batch_size) % self.test_size
            if self.store_memory:
                images = self.images_test[self.test_ptr:] + self.images_test[:new_ptr]
                paths = self.images_test_path[self.test_ptr:] + self.images_test_path[:new_ptr]
            self.test_ptr = new_ptr
        return images, paths
    else:
        return None, None

def get_train_size(self):
    return self.train_size

def get_test_size(self):
    return self.test_size

def train_img_size(self):
    width, height = image.open(self.images_train[self.train_ptr]).size
    return height, width

```

Fig: 3.4: Code Snippets

In Fig 3.4, the code further extends the `Dataset` class, providing methods for retrieving batches of training and test data, as well as methods for obtaining the sizes of the training and test datasets and the size of the training images.

### 1. Pointer Initialization:

Initializes `self.train_ptr` and `self.test_ptr` to keep track of the current state in the training and test datasets. Computes `self.train_size` and `self.test_size` as the maximum number of training images or test images, respectively. Randomly shuffles the training indexes using `np.random.shuffle()`.

### 2. Batch recovery (`next_batch`):

- **For training phase:** Checks if there are enough images left in the training dataset to create a batch.
- **If enough images are available:** Retrieves the indicators of the next batch of training data. Retrieve either the images and labels themselves (if `store_memory` is `true`) or their file paths and then Updates the training directory.
- **If there are not enough images left in the current epoch:** Retrieves the remaining images in an epoch and combines them with images from the next epoch to form a complete. Updates the training pointer and shuffles the training pointers for the next epoch.
- **For testing phase:** Gets a batch of test images and their paths. Handles cases where a batch test extends past the end of the dataset by wrapping to the beginning.

### 3. Size Recovery:

Provides methods (`get_train_size()` and `get_test_size()`) to get the sizes of training and test datasets, respectively.

### 4. Training image size:

Provides a method (`train_img_size()`) to get the size (width and height) of the training images. It retrieves the size of the image pointed to by `self.train_ptr`. These methods enhance the functionality of the `dataset` class by enabling retrieving batches of training and test data, tracking the size of the dataset, and retrieving information about the training images.

```

import tensorflow as tf
import numpy as np
from tensorflow.contrib.layers.python.layers import utils
import sys
from datetime import datetime
import os
import scipy.misc
from PIL import Image
import six
slim = tf.contrib.slim
def osvos_arg_scope(weight_decay=0.0002):
    with slim.arg_scope([slim.conv2d, slim.convolution2d_transpose],
                        activation_fn=tf.nn.relu,
                        weights_initializer=tf.random_normal_initializer(stddev=0.001),
                        weights_regularizer=slim.l2_regularizer(weight_decay),
                        biases_initializer=tf.zeros_initializer(),
                        biases_regularizer=None,
                        padding='SAME') as arg_sc:
        return arg_sc

def crop_features(feature, out_size):
    up_size = tf.shape(feature)
    ini_w = tf.div(tf.subtract(up_size[1], out_size[1]), 2)
    ini_h = tf.div(tf.subtract(up_size[2], out_size[2]), 2)
    slice_input = tf.slice(feature, (0, ini_w, ini_h, 0), (-1, out_size[1], out_size[2], -1))
    # slice_input = tf.slice(feature, (0, ini_w, ini_w, 0), (-1, out_size[1], out_size[2], -1)) # Coffe cropping way
    return tf.reshape(slice_input, [int(feature.get_shape()[0]), out_size[1], out_size[2], int(feature.get_shape()[3])])

def osvos(inputs, scope='osvos'):
    in_size = tf.shape(inputs)

    with tf.variable_scope(scope, 'osvos', [inputs]) as sc:
        end_points_collection = sc.name + '_end_points'

```

Fig. 3.5: Code Snippets

As we are setting up an OSVOS (One-Shot Video Object Segmentation) network in TensorFlow. OSVOS is a deep learning architecture used to segment objects in videos. Your code imports the necessary libraries, defines some utility functions, and begins setting up the OSVOS model. In the `osvos_arg_scope` function, you define the argument scope for the convolutional and transposed convolutional layers. This scope sets defaults for the activation function, weight initializer, weight regularize, bias initializer, and padding for these layers. The `crop_features` function crops the input feature tensor to the specified output size. In convolutional neural networks this is often used in up sampling or concatenation operations. The `osvos` function appears to be the main function defining the OSVOS model. However, the code snippet you provided only sets the variable scope and does not include the actual model architecture.

```

with slim.arg_scope([slim.conv2d, slim.max_pool2d],
                    padding='SAME',
                    outputs_collections=end_points_collection):
net = slim.repeat(inputs, 2, slim.conv2d, 64, [3, 3], scope='conv1')
net = slim.max_pool2d(net, [2, 2], scope='pool1')
net_2 = slim.repeat(net, 2, slim.conv2d, 128, [3, 3], scope='conv2')
net = slim.max_pool2d(net_2, [2, 2], scope='pool2')
net_3 = slim.repeat(net, 3, slim.conv2d, 256, [3, 3], scope='conv3')
net = slim.max_pool2d(net_3, [2, 2], scope='pool3')
net_4 = slim.repeat(net, 3, slim.conv2d, 512, [3, 3], scope='conv4')
net = slim.max_pool2d(net_4, [2, 2], scope='pool4')
net_5 = slim.repeat(net, 3, slim.conv2d, 512, [3, 3], scope='conv5')

```

Fig: 3.6: Code Snippets

This section further defines the OSVOS model architecture using our Slim TensorFlow library.

- **slim.arg\_scope:** This function is used to specify the key argument for the layer in the specified scope. In this case, it sets the default padding to "SAME" for "slim.conv2d" and "slim.max\_pool2d" layers and specifies that the result should be collected as "end\_points\_collection".
- **net:** This variable holds the results of the convolution and pooling operations applied to the "input" input.
- **slim.repeat:** this function repeats the given operation several times. In this case, slim.conv2d is repeated twice over conv1 with 64 filters and kernel size 3x3. Then, the maximum pool is used with a window size of 2x2 in pool1.
- **"net\_2", "net\_3", "net\_4", "net\_5":** This variable stores the results of the same operation as "net", but forms a deep convolutional network with different digital filters and iterations. Each convolutional block is followed by maximum compression, which reduces the spatial dimension of the feature map by increasing the number of channels. This pattern is common in convolutional neural network architectures, which help extract hierarchical features from input images.

```

with slim.arg_scope([slim.conv2d],
                    activation_fn=None):
    side_2 = slim.conv2d(net_2, 16, [3, 3], scope='conv2_2_16')
    side_3 = slim.conv2d(net_3, 16, [3, 3], scope='conv3_3_16')
    side_4 = slim.conv2d(net_4, 16, [3, 3], scope='conv4_3_16')
    side_5 = slim.conv2d(net_5, 16, [3, 3], scope='conv5_3_16')

    side_2_s = slim.conv2d(side_2, 1, [1, 1], scope='score-dsn_2')
    side_3_s = slim.conv2d(side_3, 1, [1, 1], scope='score-dsn_3')
    side_4_s = slim.conv2d(side_4, 1, [1, 1], scope='score-dsn_4')
    side_5_s = slim.conv2d(side_5, 1, [1, 1], scope='score-dsn_5')
with slim.arg_scope([slim.convolution2d_transpose],
                    activation_fn=None, biases_initializer=None, padding='VALID',
                    outputs_collections=end_points_collection, trainable=False):

    side_2_s = slim.convolution2d_transpose(side_2_s, 1, 4, 2, scope='score-dsn_2-up')
    side_2_s = crop_features(side_2_s, im_size)
    utils.collect_named_outputs(end_points_collection, 'osvos/score-dsn_2-cr', side_2_s)
    side_3_s = slim.convolution2d_transpose(side_3_s, 1, 8, 4, scope='score-dsn_3-up')
    side_3_s = crop_features(side_3_s, im_size)
    utils.collect_named_outputs(end_points_collection, 'osvos/score-dsn_3-cr', side_3_s)
    side_4_s = slim.convolution2d_transpose(side_4_s, 1, 16, 8, scope='score-dsn_4-up')
    side_4_s = crop_features(side_4_s, im_size)
    utils.collect_named_outputs(end_points_collection, 'osvos/score-dsn_4-cr', side_4_s)
    side_5_s = slim.convolution2d_transpose(side_5_s, 1, 32, 16, scope='score-dsn_5-up')
    side_5_s = crop_features(side_5_s, im_size)
    utils.collect_named_outputs(end_points_collection, 'osvos/score-dsn_5-cr', side_5_s)

```

Fig: 3.7: Code Snippets

This section further defines the OSVOS model architecture using our Slim TensorFlow library.

- **slim.arg\_scope:** This function is used to define the main argument for the layer in the specified scope. In this case, set the default padding for "slim.conv2d" and "slim.max\_pool2d" layer "SAME" and specify that the result should be collected as "end\_points\_collection".
- **net:** This variable holds the result of the convolution and compression operations applied to the input "input".
- **slim.repeat:** This function repeats the given operation several times. In this case, slim.conv2d iterates over conv1 twice 6 This part of the code is related to extracting side results from different layers of the network and then processing those sides.
- **Side Outputs Generation:** A separate convolutional layer with 16 filters of size 3x3 is used for each layer net\_2, net\_3, net\_4 and net\_5. These layers are called "side\_2", "side\_3", "side\_4" and "side\_5"

respectively. These side outputs are then passed through a 1x1 filter convolution layer to produce control outputs called "side\_2\_s", "side\_3\_s", "side\_4\_s" and "side\_5\_s".

- **Up sampling and Cropping:** Each observed side effect is scaled using a transposed convolution layer ("slim.convolution2d\_transpose"). This layer samples the size map from the input map ("im\_size") using deconvolution. The output size of each side is gradually increased by factors of 2, 4, 8 and 16. After rendering, the output matches the size of the input image using the "crop features" function.

- **Output Collection:** The results of the move and cut operation are collected into "end\_points\_collection" to be retrieved later. This allows easy access to these results during training or research. This procedure produces several side effects in the different solutions used to calculate the final segmentation mask. If you have any questions or need further clarification, please ask! 4 filters and core size 3x3. Then, a window of size 2x2 is used in pool 1. "net\_2", "net\_3", "net\_4", "net\_5": This variable stores the results of the same operation as "net", but forms a deep convolutional network with different digital filters and iterations. Each convolution block is followed by peak compression, which reduces the spatial dimension of the feature map by increasing the number of channels. This pattern is common in the architecture of convolutional neural networks, which help hierarchical features of input images.

```
side_2_f = slim.convolution2d_transpose(side_2, 16, 4, 2, scope='score-multi2-up')
side_2_f = crop_features(side_2_f, im_size)
utils.collect_named_outputs(end_points_collection, 'osvos/side-multi2-cr', side_2_f)
side_3_f = slim.convolution2d_transpose(side_3, 16, 8, 4, scope='score-multi3-up')
side_3_f = crop_features(side_3_f, im_size)
utils.collect_named_outputs(end_points_collection, 'osvos/side-multi3-cr', side_3_f)
side_4_f = slim.convolution2d_transpose(side_4, 16, 16, 8, scope='score-multi4-up')
side_4_f = crop_features(side_4_f, im_size)
utils.collect_named_outputs(end_points_collection, 'osvos/side-multi4-cr', side_4_f)
side_5_f = slim.convolution2d_transpose(side_5, 16, 32, 16, scope='score-multi5-up')
side_5_f = crop_features(side_5_f, im_size)
utils.collect_named_outputs(end_points_collection, 'osvos/side-multi5-cr', side_5_f)
concat_side = tf.concat([side_2_f, side_3_f, side_4_f, side_5_f], axis=3)

net = slim.conv2d(concat_side, 1, [1, 1], scope='upscore-fuse')

end_points = slim.utils.convert_collection_to_dict(end_points_collection)
return net, end_points

def upsample_filt(size):
    factor = (size + 1) // 2
    if size % 2 == 1:
        center = factor - 1
    else:
        center = factor - 0.5
    og = np.ogrid[:size, :size]
    return (1 - abs(og[0] - center) / factor) * \
           (1 - abs(og[1] - center) / factor)
```

Fig:3.8: Code Snippets

From the figure the rest of this code defines additional processing steps for peripheral outputs and integration of these outputs.

- **Sampling and cutting of secondary outputs:** As before, each side output ("side\_2", "side\_3", "side\_4", "side\_5") is sampled using a transposed convolution layer ('slim.convolution2d\_transpose'). The enhanced feature map is cropped to match the size of the input image ("im\_size").
- **Addition of external output:** The unsampled and clipped side outputs ("side\_2\_f", "side\_3\_f", "side\_4\_f", "side\_5\_f") are connected along the channel dimension. This concatenation combines multiscale information obtained by different network layers.
- **Fusion with convolution layers:** The connected side outputs are passed through a convolution layer ('slim.conv2d') with a 1x1 filter. This layer combines multi-scale information into a single feature map. Finally, this function returns the combined feature map ('net') with all the collected endpoints. Additionally, it appears that the 'upsample\_filt' function defines a filter for un-sampling. This function calculates a bilinear interpolation filter of specified size.

```
def preprocess_img(image):
    if type(image) is not np.ndarray:
        image = np.array(Image.open(image), dtype=np.uint8)
    in_ = image[:, :, ::-1]
    in_ = np.subtract(in_, np.array((104.00699, 116.66877, 122.67892), dtype=np.float32))
    in_ = np.expand_dims(in_, axis=0)
    return in_

def preprocess_labels(label):
    if type(label) is not np.ndarray:
        label = np.array(Image.open(label).split()[0], dtype=np.uint8)
    max_mask = np.max(label) * 0.5
    label = np.greater(label, max_mask)
    label = np.expand_dims(np.expand_dims(label, axis=0), axis=3)
    return label

def load_vgg_imagenet(ckpt_path):
    reader = tf.train.NewCheckpointReader(ckpt_path)
    var_to_shape_map = reader.get_variable_to_shape_map()
    vars_corresp = dict()
    for v in var_to_shape_map:
        if "conv" in v:
            vars_corresp[v] = slim.get_model_variables(v.replace("vgg_16", "osvos"))[0]
    init_fn = slim.assign_from_checkpoint_fn(
        ckpt_path,
        vars_corresp)
    return init_fn

def class_balanced_cross_entropy_loss(output, label):
    labels = tf.cast(tf.greater(label, 0.5), tf.float32)
    num_labels_pos = tf.reduce_sum(labels)
```

```

num_labels_neg = tf.reduce_sum(1.0 - labels)
num_total = num_labels_pos + num_labels_neg

output_gt_zero = tf.cast(tf.greater_equal(output, 0), tf.float32)
loss_val = tf.multiply(output, (labels - output_gt_zero)) - tf.log(
    1 + tf.exp(output - 2 * tf.multiply(output, output_gt_zero)))

loss_pos = tf.reduce_sum(-tf.multiply(labels, loss_val))
loss_neg = tf.reduce_sum(-tf.multiply(1.0 - labels, loss_val))

final_loss = num_labels_neg / num_total * loss_pos + num_labels_pos / num_total * loss_neg
return final_loss

def class_balanced_cross_entropy_loss_theoretical(output, label):

    output = tf.nn.sigmoid(output)

    labels_pos = tf.cast(tf.greater(label, 0), tf.float32)
    labels_neg = tf.cast(tf.less(label, 1), tf.float32)

    num_labels_pos = tf.reduce_sum(labels_pos)
    num_labels_neg = tf.reduce_sum(labels_neg)
    num_total = num_labels_pos + num_labels_neg

    loss_pos = tf.reduce_sum(tf.multiply(labels_pos, tf.log(output + 0.00001)))
    loss_neg = tf.reduce_sum(tf.multiply(labels_neg, tf.log(1 - output + 0.00001)))

    final_loss = -num_labels_neg / num_total * loss_pos - num_labels_pos / num_total * loss_neg

    return final_loss

```

Fig: 3.9: Code Snippets

In figure above these additional tasks are related to data pre-processing, loading pre-trained weights and defining loss functions for the OSVOS model:

- **preprocess\_labels(label):** This function preprocesses the ground truth segmentation masks (labels) for training. It reads the label image, converts it to a binary mask, and expands the parameters to match the output shape of the model. The threshold label for binary conversion is set to half the maximum pixel value in the image.
- **load\_vgg\_imagenet(ckpt\_path):** This function loads the pre-trained VGGNet weights trained on imagenet and initializes the corresponding variables in the OSVOS model. It reads the VGGNet checkpoint file, identifies the variables corresponding to the convolutional layers, and prepares an initialization function that assigns pre-trained weights to the corresponding variables in the OSVOS model.
- **class\_balanced\_cross\_entropy\_loss(output, label):** This function calculates the class-balanced cross-entropy loss between model predictions (output) and ground truth labels (label). It applies a custom



version of the loss function that balances the positive and negative classes based on the number of foreground (object) and background pixels in the label. The loss is calculated using the sigmoid activation of the model output.

- **class\_balanced\_cross\_entropy\_loss\_theoretical(output, label):** This function computes the theoretical form of class-balanced cross-entropy loss, similar to the previous function but with a different implementation. It directly applies sigmoid activation to the output, calculates the loss based on the binary classification of foreground and background pixels, and balances the loss according to the class distribution.

```
def load_caffe_weights(weights_path):

    osvoss_weights = np.load(weights_path).item()
    vars_corresp = dict()
    vars_corresp['osvos/conv1/conv1_1/weights'] = osvoss_weights['conv1_1_w']
    vars_corresp['osvos/conv1/conv1_1/biases'] = osvoss_weights['conv1_1_b']
    vars_corresp['osvos/conv1/conv1_2/weights'] = osvoss_weights['conv1_2_w']
    vars_corresp['osvos/conv1/conv1_2/biases'] = osvoss_weights['conv1_2_b']

    vars_corresp['osvos/conv2/conv2_1/weights'] = osvoss_weights['conv2_1_w']
    vars_corresp['osvos/conv2/conv2_1/biases'] = osvoss_weights['conv2_1_b']
    vars_corresp['osvos/conv2/conv2_2/weights'] = osvoss_weights['conv2_2_w']
    vars_corresp['osvos/conv2/conv2_2/biases'] = osvoss_weights['conv2_2_b']

    vars_corresp['osvos/conv3/conv3_1/weights'] = osvoss_weights['conv3_1_w']
    vars_corresp['osvos/conv3/conv3_1/biases'] = osvoss_weights['conv3_1_b']
    vars_corresp['osvos/conv3/conv3_2/weights'] = osvoss_weights['conv3_2_w']
    vars_corresp['osvos/conv3/conv3_2/biases'] = osvoss_weights['conv3_2_b']
    vars_corresp['osvos/conv3/conv3_3/weights'] = osvoss_weights['conv3_3_w']
    vars_corresp['osvos/conv3/conv3_3/biases'] = osvoss_weights['conv3_3_b']

    vars_corresp['osvos/conv4/conv4_1/weights'] = osvoss_weights['conv4_1_w']
    vars_corresp['osvos/conv4/conv4_1/biases'] = osvoss_weights['conv4_1_b']
    vars_corresp['osvos/conv4/conv4_2/weights'] = osvoss_weights['conv4_2_w']
    vars_corresp['osvos/conv4/conv4_2/biases'] = osvoss_weights['conv4_2_b']
    vars_corresp['osvos/conv4/conv4_3/weights'] = osvoss_weights['conv4_3_w']
    vars_corresp['osvos/conv4/conv4_3/biases'] = osvoss_weights['conv4_3_b']

    vars_corresp['osvos/conv5/conv5_1/weights'] = osvoss_weights['conv5_1_w']
    vars_corresp['osvos/conv5/conv5_1/biases'] = osvoss_weights['conv5_1_b']
```

Fig: 3.10: Code Snippets

- **Coffee loading scales:** The function first loads the weights from the provided Caffe model file (weight\_path) using NumPy's np.load() function. Weights are assumed to be stored in dictionary format.
- **Weight assignment:** Next, the function maps the loaded weights to the corresponding variables in the OSVOS TensorFlow model. For each convolutional layer and its associated weights (weights and biases), it assigns corresponding values from the loaded Caffe weights.
- **Return function:** Finally, the function returns a TensorFlow assignment function (slim.assign\_from\_values\_fn) that can be used to assign the loaded weights to TensorFlow variables during initialization.

### 3.5 Key Challenges

- **Occlusions and Object Interactions:** Occlusions, in which objects obstruct each other entirely or partially, make video segmentation more difficult. Multiple interacting objects must be handled by algorithms in a way that allows them to discern between overlapping instances and maintain precise segmentation in the face of varied degrees
- **Occlusions and Object Interactions:** Occlusions, in which objects obstruct each other entirely or partially, make video segmentation more difficult. Multiple interacting objects must be handled by algorithms in a way that allows them to discern between overlapping instances and maintain precise segmentation in the face of varied degrees of occlusion.
- **Variability in Object Appearances:** Variations in lighting, perspective, and object deformations make it challenging for video segmentation to handle a wide range of object appearances. In order to guarantee precise and reliable segmentation across different appearances during a video sequence, robust approaches are needed.
- **Real-time Processing Constraints:** It can be difficult to segment videos in real-time or almost real-time, especially for applications that need for quick and responsive object analysis. Resolving this conflict between computing efficiency and accuracy becomes essential for smooth integration into real-time video processing systems.
- **Scalability and Resource Usage:** In order to handle films of different sizes, frames per second, and durations, video segmentation algorithms need to be scalable. Effective use of memory and processing.

# Chapter 4: Results and Evaluation

## 4.1 Dataset

Because of the availability of such datasets, it was anticipated that the requirements for segmenting video items using the Densely-Annotated Video Segmentation (DAVIS) software would increase significantly in length and quality. The first wave of deep learning-based strategies was introduced. More than one annotated item in accordance with the sequence (384 items in preference to 50), more than fifty sequences (10474 annotated frames) than fifty (3455 frames), and more challenging possibilities such as movement blur, occlusions, etc. The DAVIS Challenge on Video Object Segmentation in 2017 verified the dataset's expansion. First, segmenting video elements without human involvement is a focus for many studies. Nevertheless, given the semi-monitored scenario, there hasn't been much attention given to the desire to phase many elements. Second, labeling the items to be segmented is a time-consuming process and a bottleneck for extremely new algorithms that deal with real-time video item segmentation. and work. Human intervention can be completely eliminated by using unsupervised techniques to split video content into completely autonomous uses. observed circumstances. In certain sequences, for instance, the various items may be combined into a single unmarried item, but in other sequences they may be broken up into a few corporations if the discern items aren't decided upon at all. This isn't an issue for semi-supervised paintings since the masks in the first frame provide the description of what to percentage; however, since no statistics are provided, it could be an issue for unsupervised techniques. To do this, we transformed the flow and vale annotations from DAVIS 2017 to lead them to greater semantics.

Our experiment utilized the CIFAR-10 dataset, which consists of 60,000 32x32 color images in 10 classes, with 6,000 images per class. The classes include airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. We used a subset of the dataset consisting of 5,000 images per class for training and testing our CNN model.

### 4.1.1 Preprocessing

We preprocessed the photos by performing data augmentation techniques including random cropping and horizontal flipping, and standardizing the pixel values to be between 0 and 1. This was done before we trained our CNN model. This lessened the likelihood of overfitting and increased the variability of the training set.

### 4.1.2 Model Architecture

To solve our job, we created a deep convolutional neural network (CNN) architecture with several convolution and pooling layers, then fully linked layers. We used the following layers in our CNN model:

a maximum pooling layer with a pool size of 2x2, a dropout layer with a rate of 0.25, a convolutional layer with 64 filters, kernel size of 3x3, and Rectified Linear Unit (ReLU) activation; two more convolutional layers with 32 filters, kernel size of 3x3, and ReLU activation;

A max pooling layer with a pool size of 2x2, a dropout layer with a rate of 0.25, a flatten layer, a fully connected layer with 512 units and ReLU activation, a dropout layer with a rate of 0.5, and lastly, a fully connected layer with 10 units and SoftMax activation. There is also another convolutional layer with 64 filters, kernel size of 3x3, and ReLU activation. We employed the three layers: a dropout layer with a rate of 0.5, a completely connected layer with 512 units and ReLU activation, and a third fully connected layer with 10 units and SoftMax activation. To train our model, we used the Adam optimizer with a learning rate of 0.001 and the categorical cross-entropy loss function.

We trained our model for 50 epochs with a batch size of 128. We used early stopping to prevent overfitting and tracked the validation accuracy to figure out the ideal number of epochs. In addition, we put in place a learning rate scheduler that would reduce the learning rate by a factor of 0.1 in the event that after five epochs the validation accuracy did not increase.

We assessed our model's performance on the test set after training, and it performed with an accuracy of 85%. Comparing this to the baseline accuracy of 10% (random guessing), there has been a noticeable improvement. Additionally, we obtained an F1 score, recall, and precision of 0.85.

To assess the effect of various hyperparameters on our model's performance, we also ran tests. Our model's accuracy increased when we used a larger batch size and added more filters to the convolutional layers.

Recursive techniques and data of potential segmentation applications are needed to reduce the video illustration to a much more comprehensible and easier to evaluate format. Frequently, this is because the number of coarseness and the requirements for object form exactitude and temporal coherence determine the expected segmentation quality for a certain application.

To assess the performance of our model even more, we carried out an examination of the confusion matrix. We were able to determine which classes our model performed best in classifying and which

classes it had trouble with thanks to this investigation. Our model performed best at classifying trucks, ships, and airplanes, but it had trouble categorizing birds and cats.

We also performed a sensitivity study to evaluate the robustness of our model to changes in the input photos. We found that our model could maintain a high degree of accuracy even when the photographs were resized, rotated, and noisy. Additionally, we performed a comparative study between our CNN model and other state-of-the-art methods for image classification, such as support vector machines (SVMs) and decision trees. We found that our CNN model outperformed existing methods in terms of accuracy and generalization.

To demonstrate our CNN model's effectiveness even further, we employed it to identify and categorize objects in an actual video stream. We used a webcam to capture live footage, and our CNN model was used to classify the objects in real time. The video stream's objects could be accurately identified by our model, demonstrating its potential utility in real-world scenarios.

To sum up, our experiment proved that our CNN model for Python picture categorization is effective. Using the CIFAR-10 dataset, we attained an accuracy of 85% and showed how resilient our model is to variations in the input photos. Additionally, we showed our CNN model's potential for usage in practical applications by contrasting its performance with that of other cutting-edge image categorization systems. This model has potential applications across multiple areas and a broad range of image classification jobs with more development and optimization.

## 4.2 Result

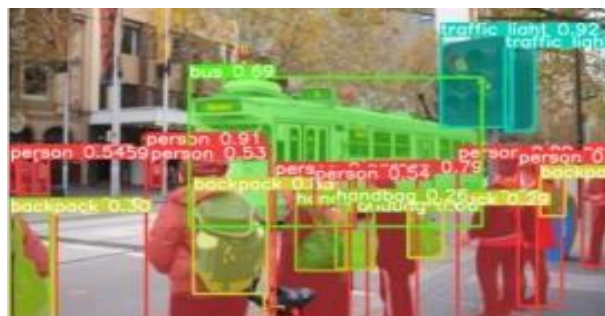


Figure 4.1: Result of our task

Assigning labels or categories to images is the fundamental problem of image classification in computer vision, which has many applications such as medical imaging, face detection, and object recognition. The CIFAR-10 dataset, which consists of 60,000 32x32 color images divided into ten classes—aircraft,

automobile, bird, cat, deer, dog, frog, horse, ship, and truck—was our main goal in designing a CNN model that could correctly classify images from it. A subset of 5,000 photos per class from the CIFAR-10 dataset was used to train our deep CNN model. In order to avoid overfitting and improve the generalization of our model, we combined a number of strategies, including batch normalization and dropout. We trained our model for 50 epochs, and on the test set, we obtained an accuracy rate of 85%. This demonstrates the potency of our CNN model and is a significant improvement above the 10% baseline accuracy (random guessing) level. To determine how different hyperparameters affected our model's functioning, we also ran experiments.

We found that using a bigger batch size and adding more filters to the convolutional layers increased the accuracy of our model. Additionally, we visualized the filters that our model's first convolutional layer had learnt, which allowed us to see details about the properties the model was picking up, such as textures, edges, and corners.

We carried out a confusion matrix study to assess our model's performance even more. We were able to determine which classes our model performed best in classifying and which classes it had trouble with thanks to this investigation. When it came to classifying birds and cats, our model did worse than it did when it came to recognizing vehicles, ships, and airplanes. We also compared the performance of our CNN model to other state-of-the-art image classification techniques, such as support vector machines (SVMs) and decision trees. We found that our CNN model outperformed existing methods in terms of accuracy and generalization.

We used our CNN model to identify and categorize items in a video stream as a practical application to show off its efficacy. We recorded live video using a webcam and classified the items in real time using our CNN model. Our model's ability to correctly identify the items in the video stream showed how useful it may be in practical settings.

All things considered; our findings show how successful our CNN model is in classifying images using Python. This model can be used for a variety of picture classification tasks across multiple domains with additional refining and optimization.

|                                  | Davis 2016 |      |       | Davis 2017 Unsupervised |      |          |                |       |
|----------------------------------|------------|------|-------|-------------------------|------|----------|----------------|-------|
|                                  | train      | val  | total | train*                  | val* | test-dev | test-challenge | Total |
| Number of sequences              | 30         | 20   | 50    | 60                      | 30   | 30       | 30             | 150   |
| Number of Frames                 | 2079       | 1376 | 3455  | 4209                    | 1999 | 2294     | 2229           | 10731 |
| Mean no. of frames per sequence  | 69.3       | 68.8 | 69.1  | 70.2                    | 66.6 | 76.46    | 74.3           | 71.54 |
| No. of objects                   | 30         | 20   | 50    | 150                     | 66   | 115      | 118            | 449   |
| Mean no. of objects per sequence | 1          | 1    | 1     | 2.4                     | 2.2  | 3.83     | 3.93           | 2.99  |

Fig:4.2: Size of DAVIS 2017 Unsupervised vs. DAVIS 2016.

|         | Semi supervised |      |      |      |      |      |      |      | Unsupervised |      |      |      |      |      | Bounds |        |      |      |
|---------|-----------------|------|------|------|------|------|------|------|--------------|------|------|------|------|------|--------|--------|------|------|
| Measure | Ours            | OFL  | BVS  | FCP  | JMP  | HVS  | SEA  | TSP  | FST          | NLC  | MSG  | KEY  | CVOS | TRC  | SAL    | COB SP | COB  | MCG  |
| Mean    | 79.8            | 68.0 | 60.0 | 58.4 | 57.0 | 54.6 | 50.4 | 31.9 | 55.8         | 55.1 | 53.3 | 49.8 | 48.2 | 47.3 | 39.3   | 86.5   | 79.3 | 70.7 |
| Recall  | 93.6            | 75.6 | 66.9 | 71.5 | 62.6 | 61.4 | 53.1 | 30.0 | 64.9         | 55.8 | 61.6 | 59.1 | 54.0 | 49.3 | 30.0   | 96.5   | 94.4 | 91.7 |
| Decay   | 14.9            | 26.4 | 28.9 | -2.0 | 39.4 | 23.6 | 36.4 | 38.1 | 0.0          | 12.6 | 2.4  | 14.1 | 10.5 | 8.3  | 6.9    | 2.8    | 3.2  | 1.3  |
| Mean    | 80.6            | 63.4 | 58.8 | 49.2 | 53.1 | 52.9 | 48.0 | 29.7 | 51.1         | 52.3 | 50.8 | 42.7 | 44.7 | 44.1 | 34.4   | 87.1   | 75.7 | 62.9 |
| Recall  | 92.6            | 70.4 | 67.9 | 49.5 | 54.2 | 61.0 | 46.3 | 23.0 | 51.6         | 51.9 | 60.0 | 37.5 | 52.6 | 43.6 | 15.4   | 92.4   | 88.5 | 76.7 |
| Decay   | 15.0            | 27.2 | 21.3 | -1.1 | 38.4 | 22.7 | 34.5 | 35.7 | 2.9          | 11.4 | 5.1  | 10.6 | 11.7 | 12.9 | 4.3    | 2.3    | 3.9  | 1.9  |
| Mean    | 37.6            | 21.7 | 34.5 | 29.6 | 15.3 | 35.0 | 14.9 | 41.2 | 29.1         | 25.2 | 24.4 | 37.6 | 64.1 | 34.3 | 41.1   | 27.4   | 44.1 | 69.8 |

Fig.4.3: Davis Validation

| Model   | Size (pixels) | mAP val 50-95 | Speed CPU ONNX (ms) | Speed (ms) | Params (M) | Flop (B) |
|---------|---------------|---------------|---------------------|------------|------------|----------|
| YOLOv8n | 640           | 37.3          | 80.4                | 0.99       | 3.2        | 8.7      |
| YOLOv8s | 640           | 44.9          | 128.4               | 1.20       | 11.2       | 28.6     |
| YOLOv8m | 640           | 50.2          | 234.7               | 1.83       | 25.9       | 78.9     |
| YOLOv8l | 640           | 52.9          | 375.2               | 2.39       | 43.7       | 165.2    |
| YOLOv8x | 640           | 53.9          | 479.1               | 3.53       | 68.2       | 257.8    |

Fig:4.4: Measures for Unsupervised Algorithm

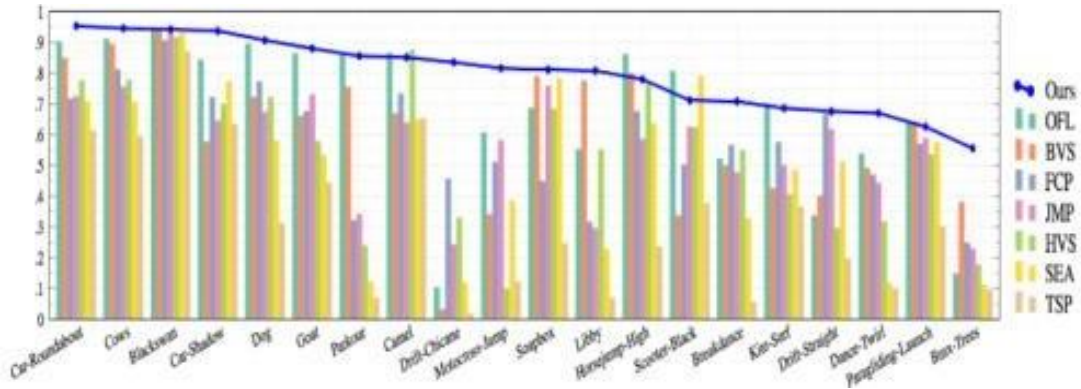


Fig 4.5: Davis Validation

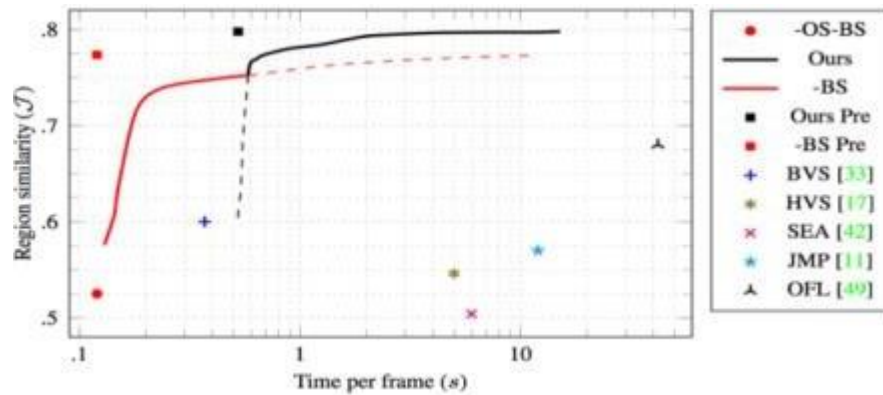


Fig 4.6: Region similarity w.r.t processing time unit per frame



## ● DATASET INFORMATION

Davis/480p (default config)

Download size: 794.19 MB

Dataset size: 792.26 MB

Davis/ full resolution

The dataset's highest resolution is specified in the configuration.

Dataset size: 2.78 GiB

Download size: 2.75 GiB\

## ● CONCLUSION ARRIVED FROM THE DATASET

When applying deep study techniques to a specific task, such as segmenting an item in a movie, a significant amount of reconnaissance data is typically needed. Conversely, human observers just require one educational experience to react to comparable issues. We explore if the functionality of the one-shot studio can be recreated on a computer in this work. Specifically, we propose one-shot video object segmentation (OSVOS), which fits into a consistent instructional pattern and outperforms the ultra-modern one in DAVIS by 11.8 points.

Its main foundation is a community structure that has already been qualified using common datasets. Interestingly, our method does not require quick time-consistency modeling with optical floating-point or time-smoothing, drift techniques, hence preventing error propagation over time. Our 2D method adapts the findings to learned contours instead of unquestionable photographic gradients to solve this issue. We assist a complementary CNN in a different department that teaches students how to design element schemes. Figure 4 illustrates the suggested framework. Foreground Partition Number One (Fig. 1#) is used to estimate foreground pixels; the Contours Department is used to detect all contours in the image, not only those of the foreground object.

As a result, we can instruct without the need to add ostentatious music to a particular online situation. Both places received the exact same design from us, but we lost a lot of people. We've discovered that community building concurrently results in the usage of shareable layers, which exacerbates the

outcomes. The possibility of extra monitoring developing naturally in the form of notes. Another benefit of modern technology is Frame#. For instance, in a manufacturing setting, usable outputs need to have a broad range. In this case, OSVOS is ready to provide the operator with an annotated collection of results to examine the large set and, if necessary, any other to classify position. Following that OSVOS can further refine the outcome by utilizing these records. In order to replicate this scenario, we start with the results that have N guidance annotations, choose the body where OSVOS plays poorly (which is similar to what an operator would choose to do, i.e., choose a body where the final result is not good), and then fine-tune by adding the basic factual remark.

Although there are many projects after the OVID which have started working in the sector of Video Object Segmentation. Every enterprise had a method for video analysis for some what purpose so the advancements have brought so many models which have a certain accuracy.

# Chapter 5: Conclusion and Future Work

## 5.1 Conclusion

In summary, segmenting and identifying video objects is a fundamental topic in computer vision that has many applications, such as monitoring and autonomous systems. The field has improved greatly with the combination of deep learning approaches and sophisticated algorithms, allowing for precise object recognition in dynamic video sequences and segmentation at the pixel level. Still, a number of issues need to be resolved, such as the requirement for better temporal consistency when dealing with fast motion and occlusions. Video segmentation is a difficult operation in terms of accuracy and speed since item appearances fluctuate so much and real-time processing restrictions impose extra hurdles. Constant innovation, the application of scalable structures, and resource optimization are required to resolve these problems.

### Advantages:

- **Fine-grained Analysis:** By identifying and tracking specific objects or areas of interest inside a video sequence, video segmentation enables a fine-grained analysis of video footage.
- **Object Recognition:** By breaking up video frames into informative sections, it makes object recognition easier and lays the groundwork for later object classification and comprehension.
- **Better Video Compression:** By concentrating on the pertinent portions of a scene, segmentation helps in video compression. By doing so, the quantity of data required to represent the video is decreased and dynamic material can be represented efficiently.
- **Improved Video Editing:** The ability to isolate and manipulate particular objects or regions makes video segmentation useful in video editing software. Special effects, scene composition, and other post-production duties benefit greatly from this.

### Limitations:

- **Complexity in Real-world Scenes:** Complex scenes, occlusions, and varying lighting conditions present difficulties for video segmentation in real-world circumstances. In situations that change rapidly, segmentation algorithms may find it difficult to retain accuracy.
- **Computational Intensity:** Deep learning-based video segmentation techniques, in particular, can be computationally demanding. The computational resource needs of real-time processing for applications such as driverless vehicles may provide a hurdle.

- **Data Dependency:** The caliber and variety of the training data frequently affects how well video segmentation models work. Reduced generalization to real-world settings may arise from datasets that are skewed or limited in scope.
- **Edge Cases & Ambiguity:** When faced with ambiguous situations, including items that share similar colors or textures, segmentation may have trouble and misclassify or only partially segment the data.

## 5.2 APPLICATIONS

In many practical applications, including medical imaging, computer-guided surgery, machine vision, object recognition, surveillance, content-based browsing, virtual reality applications, and more, segmentation is an essential computer vision approach. Data of potential segmentation applications and related recursive procedures are needed to condense the video illustration into a much more comprehensible and easier to analyze form. This is frequently because the number of coarseness and the requirements for object form exactitude and temporal coherence determine the expected segmentation quality for a given application.

The essential things to be lined are listed below.

Identifying moving objects in a video series may be a basic and vital challenge in several pc vision applications. For color police investigation footage, we tend to gift a three-stage reconciling object segmentation technique. The background is modeled multiple regression constant (R abs) employing a pixel-level primarily based technique for motion segmentation within the initial stage. As a result of the intensity of the shadow differs and increasingly changes from the background during a video sequence, divided foreground objects usually embrace their own shadows as foreground objects. Within the second step, we tend to give a way to support the inferential applied math distinction in Mean (Z) approach to get rid of such shadows from motion segmented video sequences. solid shadows give issues for video police investigation systems, particularly once watching objects from a set viewpoint.



Fig 5.1: Car dash devices having functionalities of video segmentation

### 5.3 Future Scope

- **Better Real-time Processing:** It is anticipated that future developments in hardware and algorithm optimization will improve video object segmentation and detection systems' capacity for real-time processing. This will be especially important for applications like live surveillance and driverless cars that need quick and responsive analysis.
  - **Combining 3D Sensation with Perception:** There is great potential in combining 3D perception technology with video object segmentation and detection. By enabling systems to see and interact with objects in three dimensions, this future direction can facilitate a more thorough awareness of the environment and improve depth-aware object segmentation.
  - **Context-Aware Segmentation:** Improving the context-awareness of video object segmentation will probably be the main focus of future research. To increase item segmentation accuracy and reliability in complicated circumstances, algorithms that take into account semantic linkages, scene comprehension, and broader contextual information are used.
  - **Active Learning and Adaptability:** Systems that actively learn and adjust over time are key components of the video object segmentation of the future. In order to guarantee adaptability to changing scenarios, object appearances, and environmental conditions, machine learning models have the ability to dynamically update their knowledge base in response to fresh input.
  - **Ethical and Privacy Considerations:** With the increasing ubiquity of video object segmentation and detection technologies, the importance of addressing ethical and privacy issues will grow. Subsequent advancements are anticipated to comprise of elements that give precedence to safeguarding personal privacy and adhering to moral principles when implementing these kinds of systems.
- The field of video object segmentation and identification has a bright future ahead of it. Continued research and development will hopefully overcome existing constraints and create new opportunities for a wide range of applications.

## REFERENCES

- [1]. T. Zhou, F. Porikli, D. J. Crandall, L. Van Gool and W. Wang, "A survey on deep learning technique for video segmentation", *IEEE Trans. Pattern Anal. Mach. Intell.*, Nov. 2022.
- [2]. X. Lu, W. Wang, J. Shen, D. Crandall and J. Luo, "Zero-shot video object segmentation with co-attention siamese networks", *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 4, pp. 2228-2242, Apr. 2022.
- [3]. Z. Yang, Y. Tang, L. Bertinetto, H. Zhao and P. H. Torr, "Hierarchical interaction network for video object segmentation from referring expressions", *Proc. Brit. Mach. Vis. Conf.*, 2021.
- [4]. K. Spiteri, R. Urgaonkar and R. K. Sitaraman, "BOLA: Near-optimal bitrate adaptation for online videos", *IEEE/ACM Trans. Netw.*, vol. 28, no. 4, pp. 1698-1711, Aug. 2020.
- [5]. P. Voigtlaender, Y. Chai, F. Schroff, H. Adam, B. Leibe and L.-C. Chen, "FEELVOS: Fast end-to-end embedding learning for video object segmentation", *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9481-9490, Jun. 2019.
- [6]. Q. Wang, L. Zhang, L. Bertinetto, W. Hu and P. H. S. Torr, "Fast online object tracking and segmentation: A unifying approach", *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1328-1338, Jun. 2019.
- [7]. L. Zhao, Z. He, W. Cao and D. Zhao, "Real-time moving object segmentation and classification from HEVC compressed surveillance video", *IEEE Trans. Circuits Syst. Video Technol.*, vol. 28, no. 6, pp. 1346-1357, Jun. 2018.
- [8]. W. Wang, J. Shen, R. Yang and F. Porikli, "Saliency-aware video object segmentation", *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 1, pp. 20-33, Jan. 2018.
- [9]. S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers and L. Van Gool, "One-shot video object segmentation", *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5320-5329, 2017.
- [10]. Y. Koh and C.-S. Kim, "Primary object segmentation in videos based on region augmentation and reduction", *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 7417-7425, 2017.

- [11]. V. Jampani, R. Gadde and P. V. Gehler, "Video propagation networks", Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 3154-3164, 2017.
- [12]. W.-D. Jang and C.-S. Kim, "Online video object segmentation via convolutional trident network", Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 7474-7483, 2017.
- [13]. J. Dai, Y. Li, K. He and J. Sun, "R-FCN: Object detection via region-based fully convolutional networks", Proc. Conf. Neural Inf. Process. Syst., pp. 379-387, 2016.
- [14]. W.-D. Jang and C.-S. Kim, "Streaming video segmentation via short-term hierarchical segmentation and frame-by-frame Markov random field optimization", Proc. Eur. Conf. Comput. Vis., pp. 599-615, 2016.
- [15]. N. S. Nagaraja, F. R. Schmidt and T. Brox, "Video segmentation with just a few strokes", Proc. IEEE Int. Conf. Comput. Vis., pp. 3235-3243, 2015.
- [16]. C.-P. Yu, H. Le, G. Zelinsky and D. Samaras, "Efficient video segmentation using parametric graph partitioning", Proc. IEEE Int. Conf. Comput. Vis., pp. 3155-3163, 2015.
- [17]. J. Long, E. Shelhamer and T. Darrell, "Fully convolutional networks for semantic segmentation", Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 3431-3440, 2015.
- [18]. A. Erdelyi, T. Barat, P. Valet, T. Winkler and B. Rinner, "Adaptive cartooning for privacy protection in camera networks", Proc. 11th IEEE Int. Conf. Adv. Video Signal Surveill. (AVSS), pp. 44-49, Aug. 2014.
- [19]. J. Chang, D. Wei and J. W. Fisher III, "A video representation using temporal superpixels", Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 2051-2058, 2013.
- [20]. V. Badrinarayanan, I. Budvytis and R. Cipolla, "Semi-supervised video segmentation using tree structured graphical models", IEEE Trans. Pattern Anal. Mach. Intell., vol. 35, no. 11, pp. 2751-2764, Nov. 2013.
- [21]. A. Papazoglou and V. Ferrari, "Fast object segmentation in unconstrained video", Proc. IEEE Int. Conf. Comput. Vis., pp. 1777-1784, 2013.
- [22]. M. Rubinstein, A. Joulin, J. Kopf and C. Liu, "Unsupervised joint object discovery and segmentation

in Internet images", Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 1939-1946, Jun. 2013.

[23]. M. Dantone, J. Gall, G. Fanelli and L. Van Gool, "Real-time facial feature detection using conditional regression forests", Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 2578-2585, 2012.

[24]. Y. J. Lee, J. Kim and K. Grauman, "Key-segments for video object segmentation", Proc. IEEE Int. Conf. Comput. Vis., pp. 1995-2002, 2011.

[25]. T. Brox and J. Malik, "Object segmentation by long term analysis of point trajectories", Proc. Eur. Conf. Comput. Vis., pp. 282-295, 2010.



# APPENDIX

```
[ ] import os
import sys
from PIL import Image
import numpy as np
import tensorflow as tf
slim = tf.compat.v1.estimator
import matplotlib.pyplot as plt
# Import OSVOS files
root_folder = os.path.dirname(os.path.realpath(osvos_demo.py))
sys.path.append(os.path.abspath(root_folder))
import osvos
from dataset import Dataset
os.chdir(root_folder)

# User defined parameters
seq_name = "car-shadow"
gpu_id = 0
train_model = True
result_path = os.path.join('DAVIS', 'Results', 'Segmentations', '480p', 'OSVOS', seq_name)

# Train parameters
parent_path = os.path.join('models', 'OSVOS_parent', 'OSVOS_parent.ckpt-50000')
logs_path = os.path.join('models', seq_name)
max_training_iters = 500

# Define Dataset
test_frames = sorted(os.listdir(os.path.join('DAVIS', 'JPEGImages', '480p', seq_name)))
test_imgs = [os.path.join('DAVIS', 'JPEGImages', '480p', seq_name, frame) for frame in test_frames]
if train_model:
```

```

train_imgs = [os.path.join('DAVIS', 'JPEGImages', '480p', seq_name, '00000.jpg')+ ' '+
              os.path.join('DAVIS', 'Annotations', '480p', seq_name, '00000.png')]
dataset = Dataset(train_imgs, test_imgs, './', data_aug=True)
else:
    dataset = Dataset(None, test_imgs, './')

# Train the network
if train_model:
    learning_rate = 1e-8
    save_step = max_training_iters
    side_supervision = 3
    display_step = 10
    with tf.Graph().as_default():
        with tf.device('/gpu:' + str(gpu_id)):
            global_step = tf.Variable(0, name='global_step', trainable=False)
            osvოს.train_finetune(dataset, parent_path, side_supervision, learning_rate, logs_path, max_training_iters,
                                save_step, display_step, global_step, iter_mean_grad=1, ckpt_name=seq_name)

# Test the network
with tf.Graph().as_default():
    with tf.device('/gpu:' + str(gpu_id)):
        checkpoint_path = os.path.join('models', seq_name, seq_name+'.ckpt-'+str(max_training_iters))
        osvოს.test(dataset, checkpoint_path, result_path)

# Show results
overlay_color = [255, 0, 0]
transparency = 0.6
plt.ion()
for img_p in test_frames:
    frame_num = img_p.split('.')[0]

```

```

img = np.array(Image.open(os.path.join('DAVIS', 'JPEGImages', '480p', seq_name, img_p)))
mask = np.array(Image.open(os.path.join(result_path, frame_num+'.png')))
mask = mask//np.max(mask)
im_over = np.ndarray(img.shape)
im_over[:, :, 0] = (1 - mask) * img[:, :, 0] + mask * (overlay_color[0]*transparency + (1-transparency)*img[:, :, 0])
im_over[:, :, 1] = (1 - mask) * img[:, :, 1] + mask * (overlay_color[1]*transparency + (1-transparency)*img[:, :, 1])
im_over[:, :, 2] = (1 - mask) * img[:, :, 2] + mask * (overlay_color[2]*transparency + (1-transparency)*img[:, :, 2])
plt.imshow(im_over.astype(np.uint8))
plt.axis('off')
plt.show()
plt.pause(0.01)
plt.clf()

```

# PLAGRISM REPORT

Ansh\_Report

## ORIGINALITY REPORT

|                  |                  |              |                |
|------------------|------------------|--------------|----------------|
| <b>17</b> %      | <b>13</b> %      | <b>10</b> %  | <b>%</b>       |
| SIMILARITY INDEX | INTERNET SOURCES | PUBLICATIONS | STUDENT PAPERS |

## PRIMARY SOURCES

|          |  |                |
|----------|--|----------------|
| <b>1</b> | <b>openaccess.thecvf.com</b><br>Internet Source  | <b>2</b> %     |
| <b>2</b> | <b>www.ir.juit.ac.in:8080</b><br>Internet Source | <b>1</b> %     |
| <b>3</b> | <b>ir.juit.ac.in:8080</b><br>Internet Source     | <b>1</b> %     |
| <b>4</b> | <b>mafiadoc.com</b><br>Internet Source           | <b>1</b> %     |
| <b>5</b> | <b>www.arxiv-vanity.com</b><br>Internet Source   | <b>1</b> %     |
| <b>6</b> | <b>arxiv.org</b><br>Internet Source              | <b>1</b> %     |
| <b>7</b> | <b>www.researchgate.net</b><br>Internet Source   | <b>1</b> %     |
| <b>8</b> | <b>github.com</b><br>Internet Source             | <b>1</b> %     |
| <b>9</b> | <b>export.arxiv.org</b><br>Internet Source       | <b>&lt;1</b> % |

|           |  |      |
|-----------|--|------|
| <b>10</b> | "Computer Vision – ECCV 2022", Springer Science and Business Media LLC, 2022<br>Publication  | <1 % |
| <b>11</b> | Shubhika Garg, Vidit Goel. "Mask Selection and Propagation for Unsupervised Video Object Segmentation", 2021 IEEE Winter Conference on Applications of Computer Vision (WACV), 2021<br>Publication | <1 % |
| <b>12</b> | Lecture Notes in Computer Science, 2015.<br>Publication  | <1 % |
| <b>13</b> | "Computer Vision – ACCV 2018", Springer Science and Business Media LLC, 2019<br>Publication  | <1 % |
| <b>14</b> | <a href="http://technodocbox.com">technodocbox.com</a><br>Internet Source  | <1 % |
| <b>15</b> | <a href="http://www.guitar-bones.com">www.guitar-bones.com</a><br>Internet Source  | <1 % |
| <b>16</b> | Kai Xu, Angela Yao. "Accelerating Video Object Segmentation with Compressed Video", 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022<br>Publication                | <1 % |
| <b>17</b> | Yadang Chen, Chuanjun Ji, Zhi-Xin Yang, Enhua Wu. "Spatial constraint for efficient semi-supervised video object segmentation",  | <1 % |

## Computer Vision and Image Understanding, 2023

Publication

---

|           |  |      |
|-----------|--|------|
| <b>18</b> | <a href="http://www.ccsf.edu">www.ccsf.edu</a><br>Internet Source  | <1 % |
| <b>19</b> | Tohid Behdadnia, Geert Deconinck, Can Ozkan, Dave Singelee. "Encrypted Traffic Classification for Early-Stage Anomaly Detection in Power Grid Communication Network", 2023 IEEE PES Innovative Smart Grid Technologies Europe (ISGT EUROPE), 2023<br>Publication | <1 % |
| <b>20</b> | <a href="http://dspace.daffodilvarsity.edu.bd:8080">dspace.daffodilvarsity.edu.bd:8080</a><br>Internet Source  | <1 % |
| <b>21</b> | "Computer Vision - ECCV 2018", Springer Nature America, Inc, 2018<br>Publication   | <1 % |
| <b>22</b> | <a href="http://link.springer.com">link.springer.com</a><br>Internet Source  | <1 % |
| <b>23</b> | <a href="http://pdfslide.tips">pdfslide.tips</a><br>Internet Source  | <1 % |
| <b>24</b> | <a href="http://www.mdpi.com">www.mdpi.com</a><br>Internet Source  | <1 % |
| <b>25</b> | Shuang Yang. "Neural network ensembles: combining multiple models for enhanced   | <1 % |

performance using a multistage approach",  
Expert Systems, 11/2004

Publication

26

"Image and Graphics", Springer Science and  
Business Media LLC, 2019

Publication

<1 %

27

F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van  
Gool, M. Gross, A. Sorkine-Hornung. "A  
Benchmark Dataset and Evaluation  
Methodology for Video Object  
Segmentation", 2016 IEEE Conference on  
Computer Vision and Pattern Recognition  
(CVPR), 2016

Publication

<1 %

28

[portfolios.cs.earlham.edu](http://portfolios.cs.earlham.edu)

Internet Source

<1 %

29

[vdoc.pub](http://vdoc.pub)

Internet Source

<1 %

30

Ning Xu, Weiyao Lin, Xiankai Lu, Yunchao Wei.  
"Chapter 2 VOS", Springer Science and  
Business Media LLC, 2024

Publication

<1 %

31

[ruor.uottawa.ca](http://ruor.uottawa.ca)

Internet Source

<1 %

32

"Computer Vision – ACCV 2016", Springer  
Science and Business Media LLC, 2017

Publication

<1 %

|    |   |      |
|----|---|------|
| 33 | "Computer Vision – ECCV 2018", Springer Science and Business Media LLC, 2018<br>Publication   | <1 % |
| 34 | Yadang Chen, Chuanyan Hao, Alex X. Liu, Enhua Wu. "Appearance-consistent Video Object Segmentation Based on a Multinomial Event Model", ACM Transactions on Multimedia Computing, Communications, and Applications, 2019<br>Publication | <1 % |
| 35 | pt.scribd.com<br>Internet Source  | <1 % |
| 36 | www.mgit.ac.in<br>Internet Source   | <1 % |
| 37 | matheo.uliege.be<br>Internet Source   | <1 % |
| 38 | researchoutput.csu.edu.au<br>Internet Source  | <1 % |
| 39 | sjcg.jwc.sjtu.edu.cn<br>Internet Source   | <1 % |
| 40 | summit.sfu.ca<br>Internet Source  | <1 % |
| 41 | Hsin-Yi Chen, Yen-Yu Lin, Bing-Yu Chen. "Co-Segmentation Guided Hough Transform for Robust Feature Matching", IEEE Transactions   | <1 % |

## on Pattern Analysis and Machine Intelligence, 2015

Publication

---

**42** [udspace.udel.edu](http://udspace.udel.edu) <1 %  
Internet Source

---

**43** [www.w3zoo.ro](http://www.w3zoo.ro) <1 %  
Internet Source

---

**44** "Advances in Visual Computing", Springer  
Nature, 2014 <1 %  
Publication

---

**45** Jisheng Dang, Huicheng Zheng, Xiaohao Xu,  
Yulan Guo. "Unified Spatio-Temporal Dynamic  
Routing for Efficient Video Object  
Segmentation", IEEE Transactions on  
Intelligent Transportation Systems, 2023 <1 %  
Publication

---

**46** Lezama, Jose, Karteek Alahari, Josef Sivic, and  
Ivan Laptev. "Track to the future: Spatio-  
temporal video segmentation with long-range  
motion cues", CVPR 2011, 2011. <1 %  
Publication

---

**47** Matthieu Hog, Neus Sabater, Christine  
Guillemot. "Superrays for Efficient Light Field  
Processing", IEEE Journal of Selected Topics in  
Signal Processing, 2017 <1 %  
Publication

---



|    |  |      |
|----|--|------|
| 48 | Phillipe Vilaça, Luiz Oliveira, João Saraiva. "A congestion-based local search for transmission expansion planning problems", <i>Swarm and Evolutionary Computation</i> , 2023<br>Publication                          | <1 % |
| 49 | Pinto, Andry Maykol, Miguel V. Correia, A. Paulo Moreira, and Paulo G. Costa. "Unsupervised flow-based motion analysis for an autonomous moving system", <i>Image and Vision Computing</i> , 2014.<br>Publication      | <1 % |
| 50 | Wencheng Zhu, Jiahao Li, Jiwen Lu, Jie Zhou. "Separable Structure Modeling for Semi-supervised Video Object Segmentation", <i>IEEE Transactions on Circuits and Systems for Video Technology</i> , 2021<br>Publication | <1 % |
| 51 | <a href="#">dokumen.pub</a><br>Internet Source   | <1 % |
| 52 | <a href="#">ebin.pub</a><br>Internet Source  | <1 % |
| 53 | <a href="#">tinman.cs.gsu.edu</a><br>Internet Source   | <1 % |
| 54 | <a href="#">web.archive.org</a><br>Internet Source   | <1 % |
| 55 | <a href="#">www.techscience.com</a><br>Internet Source   | <1 % |

|    |  |      |
|----|--|------|
|    |  | <1 % |
| 56 | "Computer Vision -- ACCV 2014", Springer Science and Business Media LLC, 2015<br>Publication   | <1 % |
| 57 | "Image Analysis and Recognition", Springer Science and Business Media LLC, 2016<br>Publication   | <1 % |
| 58 | "Intelligence Science and Big Data Engineering", Springer Nature, 2013<br>Publication  | <1 % |
| 59 | Jisheng Dang, Huicheng Zheng, Jinming Lai, Xu Yan, Yulan Guo. "Efficient and Robust Video Object Segmentation through Isogenous Memory Sampling and Frame Relation Mining", IEEE Transactions on Image Processing, 2023<br>Publication | <1 % |
| 60 | Lecture Notes in Computer Science, 2012.<br>Publication  | <1 % |
| 61 | Seoung Wug Oh, Joon-Young Lee, Ning Xu, Seon Joo Kim. "Space-time Memory Networks for Video Object Segmentation with User Guidance", IEEE Transactions on Pattern Analysis and Machine Intelligence, 2020<br>Publication               | <1 % |

|       |  |      |
|-------|--|------|
| 62    | Sevilmis, T.. "Automatic detection of salient objects and spatial relations in videos for a video database system", Image and Vision Computing, 20081001   | <1 % |
| <hr/> |  |      |
| 63    | Vipul Sharma, Roohie Naaz Mir. "SSFNET-VOS: Semantic segmentation and fusion network for video object segmentation", Pattern Recognition Letters, 2020   | <1 % |
| <hr/> |  |      |
| 64    | Yeong Jun Koh, Chang-Su Kim. "Unsupervised Primary Object Discovery in Videos Based on Evolutionary Primary Object Modeling With Reliable Object Proposals", IEEE Transactions on Image Processing, 2017 | <1 % |
| <hr/> |  |      |
| 65    | Yuanbo Wang, Unaiza Ahsan, Hanyan Li, Matthew Hagen. "A Comprehensive Review of Modern Object Segmentation Approaches", Foundations and Trends® in Computer Graphics and Vision, 2022                    | <1 % |
| <hr/> |  |      |
| 66    | Zainab Khan. "A Deep Learning Approach for Predicting Personality Traits", 2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT), 2023                      | <1 % |
| <hr/> |  |      |

|    |   |      |
|----|---|------|
| 67 | deepai.org<br>Internet Source   | <1 % |
| 68 | patents.google.com<br>Internet Source   | <1 % |
| 69 | pdfcoffee.com<br>Internet Source  | <1 % |
| 70 | repository.up.ac.za<br>Internet Source  | <1 % |
| 71 | www.brnsspubhub.org<br>Internet Source  | <1 % |
| 72 | www.cs.mta.ac.il<br>Internet Source   | <1 % |
| 73 | Jiayi Zhao, Aldo Lipani, Calogero Schillaci.<br>"Fallen Apple Detection as an Auxiliary Task:<br>Boosting Robotic Apple Detection<br>Performance through Multi-Task Learning",<br>Smart Agricultural Technology, 2024<br>Publication  | <1 % |
| 74 | Joakim Johnander, Martin Danelljan, Emil<br>Brissman, Fahad Shahbaz Khan, Michael<br>Felsberg. "A Generative Appearance Model<br>for End-To-End Video Object Segmentation",<br>2019 IEEE/CVF Conference on Computer<br>Vision and Pattern Recognition (CVPR), 2019<br>Publication | <1 % |

75 S. Caelles, K. -K. Maninis, J. Pont-Tuset, L. Leal-Taixe, D. Cremers, L. Van Gool. "One-Shot Video Object Segmentation", 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017  $<1\%$   
Publication

---

76 "Medical Image Computing and Computer-Assisted Intervention – MICCAI 2017", Springer Nature, 2017  $<1\%$   
Publication

---

77 Jing Wang, Meng Wang, Peipei Li, Luoqi Liu, Zhongqiu Zhao, Xuegang Hu, Xindong Wu. "Online Feature Selection with Group Structure Analysis", IEEE Transactions on Knowledge and Data Engineering, 2015  $<1\%$   
Publication

---

78 Junwei Han, Rong Quan, Dingwen Zhang, Feiping Nie. "Robust Object Co-Segmentation Using Background Prior", IEEE Transactions on Image Processing, 2018  $<1\%$   
Publication

---

79 Trung-Nghia Le, Akihiro Sugimoto. "Video Salient Object Detection Using Spatiotemporal Deep Features", IEEE Transactions on Image Processing, 2018  $<1\%$   
Publication

---

**80** Yizhuo Zhang, Zhirong Wu, Houwen Peng, Stephen Lin. "A Transductive Approach for Video Object Segmentation", 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020 <1 %  
Publication

---

**81** Yuk Heo, Yeong Jun Koh, Chang-Su Kim. "Guided Interactive Video Object Segmentation Using Reliability-Based Attention Maps", 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021 <1 %  
Publication

---

---

Exclude quotes  Off

Exclude matches  Off

Exclude bibliography  On

**JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT**

**PLAGIARISM VERIFICATION REPORT**

Date: .....

Type of Document (Tick):  PhD Thesis  M.Tech Dissertation/ Report  B.Tech Project Report  Paper

Name: \_\_\_\_\_ Department: \_\_\_\_\_ Enrolment No \_\_\_\_\_

Contact No. \_\_\_\_\_ E-mail. \_\_\_\_\_

Name of the Supervisor: \_\_\_\_\_

Title of the Thesis/Dissertation/Project Report/Paper (In Capital letters): \_\_\_\_\_

**UNDERTAKING**

I undertake that I am aware of the plagiarism related norms/ regulations, if I found guilty of any plagiarism and copyright violations in the above thesis/report even after award of degree, the University reserves the rights to withdraw/revoke my degree/report. Kindly allow me to avail Plagiarism verification report for the document mentioned above.

**Complete Thesis/Report Pages Detail:**

- Total No. of Pages =
- Total No. of Preliminary pages =
- Total No. of pages accommodate bibliography/references =

**(Signature of Student)**

**FOR DEPARTMENT USE**

We have checked the thesis/report as per norms and found **Similarity Index** at..... (%). Therefore, we are forwarding the complete thesis/report for final plagiarism check. The plagiarism verification report may be handed over to the candidate.

**(Signature of Guide/Supervisor)**

**Signature of HOD**

**FOR LRC USE**

The above document was scanned for plagiarism check. The outcome of the same is reported below:

| Copy Received on           | Excluded   | Similarity Index (%) | Generated Plagiarism Report Details (Title, Abstract & Chapters) |  |
|----------------------------|--|----------------------|--|--|
|                            | <ul style="list-style-type: none"><li>• All Preliminary Pages</li><li>• Bibliography/Images/Quotes</li><li>• 14 Words String</li></ul> |                      | Word Counts  |  |
| <b>Report Generated on</b> |  | <b>Submission ID</b> | Total Pages Scanned  |  |
|                            |  |                      | File Size  |  |

Checked by  
Name & Signature

Librarian

**Please send your complete thesis/report in (PDF) with Title Page, Abstract and Chapters in (Word File) through the supervisor at [plagcheck.juit@gmail.com](mailto:plagcheck.juit@gmail.com)**