

Camera-Based Semantic Image Segmentation for Autonomous Vehicle

Aparna Tiwari, Shubham Bilgi, Anuja Chaudhari


Department of Information Technology

NBN Sinhgad School of Engineering, Pune, India



<https://doi.org/10.55041/ijstmt.v2i3.358>

Cite this Article: Tiwari, A., Bilgi, S. & Chaudhari, A. (2026). Camera-Based Semantic Image Segmentation for Autonomous Vehicle. International Journal of Science, Strategic Management and Technology, 02(03). <https://doi.org/10.55041/ijstmt.v2i3.358>

License:  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

Abstract— Existing models still struggle with cross-domain adaptability and generalisation, despite notable advancements in deep learning for image and video segmentation. Segmenting images and videos is a basic computer vision task with many applications in autonomous driving, industrial inspection, healthcare, and agriculture. Since the introduction of large-scale foundation models, SAM2—an enhanced version of SAM (Segment Anything Model)—has shown improved performance in complex scenarios after being optimised for segmentation tasks. However, more research is needed to fully understand SAM2's versatility and limitations in particular domains. This study assesses SAM2's performance across a range of domains and methodically examines its use in image and video segmentation. We start by defining the fundamental ideas of image segmentation, classifying foundation models, and examining the technical features of SAM and SAM2. We then explore SAM2's use in static image and video segmentation, highlighting the difficulties of cross-domain adaptability and its performance in specialised fields like medical imaging. We examined more than 200 relevant papers as part of our investigation to offer a thorough analysis of the subject. The study concludes by highlighting SAM2's advantages and disadvantages in segmentation tasks, pointing out the technical difficulties it encounters, and suggesting potential future development paths. This review offers insightful analysis and useful suggestions for maximising and utilising SAM2 in practical situations

Keywords— SAM2, Image Segmentation, Autonomous Vehicles, Computer Vision, Real-Time Processing semantic segmentation, autonomous driving, U-Net, SAM2, ORB feature detection, deep learning, computer vision

I. INTRODUCTION

Image and video segmentation are core tasks in computer vision, aimed at dividing visual data into meaningful regions based on semantic or spatial characteristics. These techniques are widely applied across various domains, including healthcare, agriculture, industrial inspection, autonomous driving, and satellite-based remote sensing. Image segmentation focuses on detecting and isolating objects, boundaries, or textures within a single frame, whereas video segmentation extends this concept over

time by processing consecutive frames while maintaining temporal and spatial coherence.

In recent years, deep learning has significantly advanced the performance of segmentation models, enabling them to handle complex real-world scenarios. Nevertheless, many existing approaches are designed for specific tasks or data modalities, which restrict their ability to generalize effectively across different domains. As a result, the development of more flexible and domain-agnostic segmentation models has become an important research direction.

The introduction of large-scale foundation models has transformed the landscape of artificial intelligence, showcasing strong zero-shot and few-shot learning capabilities. Among these, the Segment Anything Model (SAM) stands out as a foundational model for image segmentation, achieving strong performance on natural image datasets. However, SAM faces several limitations when applied to broader segmentation tasks. First, its training is primarily based on natural image datasets, which reduces its adaptability to other domains and can lead to performance degradation. Second, SAM is largely optimized for 2D image inputs, limiting its effectiveness in handling 3D data such as medical imaging. Furthermore, SAM struggles with video segmentation due to the temporal dependencies and dynamic nature of video data, which differ substantially from static image processing requirements.

To overcome these challenges, SAM2 has been introduced as an enhanced version of SAM. It is specifically designed to address the limitations of its predecessor by improving adaptability and performance across a wider range of segmentation tasks. SAM2 provides more robust and accurate results for both image and video segmentation scenarios.

To better understand the capabilities and role of SAM2 in segmentation tasks, we conducted a systematic review of existing studies. Although several surveys have explored segmentation techniques based on SAM or SAM2, many of these works focus on limited aspects and do not provide a fully comprehensive analysis.

1.1 Background of Image and Video Segmentation

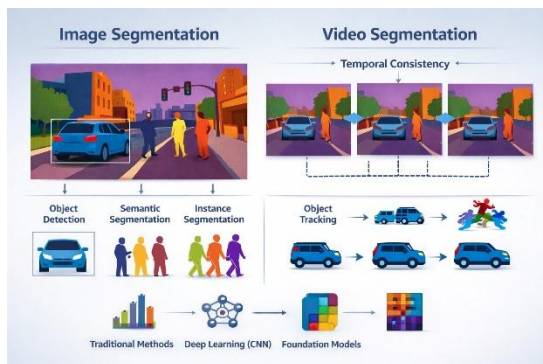


Image and video segmentation are fundamental components of computer vision, aimed at partitioning visual data into meaningful and interpretable regions. The

primary objective of segmentation is to assign labels to pixels such that regions corresponding to objects, boundaries, or specific semantic categories can be clearly distinguished. This capability plays a crucial role in enabling machines to understand and interact with visual environments.

Image segmentation focuses on analyzing a single frame and identifying relevant structures such as objects, edges, or textures. It is commonly categorized into semantic segmentation, instance segmentation, and panoptic segmentation, each differing in the level of detail and object distinction provided. On the other hand, video segmentation extends this concept by incorporating temporal information across consecutive frames. It ensures consistency in object tracking and segmentation over time, making it essential for dynamic scene understanding.

Over the years, segmentation techniques have evolved significantly. Early methods relied on traditional image processing approaches such as thresholding, edge detection, and region-based techniques. However, these methods were limited in handling complex and real-world scenarios. The emergence of deep learning, particularly convolutional neural networks (CNNs), revolutionized segmentation by enabling automatic feature extraction and improving accuracy. Models such as Fully Convolutional Networks (FCNs), U-Net, and DeepLab series demonstrated substantial improvements in performance across various applications.

More recently, transformer-based architectures and large-scale foundation models have further advanced the field by introducing global context understanding and improved generalization capabilities. Models like the Segment Anything Model (SAM) and its successor SAM2 have introduced a new paradigm by enabling prompt-based segmentation and reducing dependency on large labeled datasets.

Despite these advancements, challenges such as cross-domain adaptability, real-time processing, and robustness in dynamic environments remain open research problems. These challenges are particularly critical in applications such as autonomous driving, where accurate and efficient segmentation is required for safe navigation and decision-making.

1.2 Role of Computer Vision in Autonomous Vehicles

Computer vision plays a central role in enabling autonomous vehicles to perceive, interpret, and interact with their surrounding environment. It provides the necessary capability for machines to extract meaningful information from visual inputs such as cameras, which are widely used due to their cost-effectiveness and rich contextual data.

In autonomous driving systems, computer vision is responsible for several critical tasks, including object detection, lane detection, traffic sign recognition, and scene understanding. Among these, image and video segmentation are particularly important, as they allow the system to identify and differentiate various elements in the environment, such as vehicles, pedestrians, road boundaries, and obstacles at a pixel level. This fine-grained understanding is essential for safe navigation and decision-making.

Camera-based perception systems are widely adopted in modern autonomous vehicles because they provide high-resolution visual information and capture detailed semantic features of the environment. Unlike sensors such as LiDAR or radar, cameras are capable of capturing color, texture, and contextual cues, which are crucial for accurate scene interpretation. However, processing this visual data in real time presents significant challenges, especially in dynamic and complex driving scenarios.

Segmentation models enhance the capability of computer vision systems by enabling precise localization of objects and improving the understanding of spatial relationships within a scene. In video-based applications, maintaining temporal consistency across frames is equally important, as it ensures stable tracking of moving objects and reduces prediction errors.

With the advancement of deep learning and foundation models, computer vision systems have become more robust and adaptable. Models such as SAM2 introduce new possibilities by enabling generalized segmentation across different domains and improving performance in dynamic environments. These capabilities make them highly relevant for autonomous driving applications, where real-time performance and adaptability are critical.

Despite these advancements, challenges such as varying lighting conditions, occlusions, and computational constraints continue to affect system performance. Therefore, developing efficient and reliable segmentation frameworks remains a key requirement for advancing camera-based autonomous vehicle systems.

1.3 Limitations of Existing Segmentation Models

Despite significant advancements in image and video segmentation, existing models still face several limitations that restrict their effectiveness in real-world applications, particularly in autonomous driving scenarios. These limitations arise from constraints in model design, data dependency, and computational requirements.

One of the primary challenges is the lack of generalization across domains. Many segmentation models, especially those based on convolutional neural networks such as U-Net and DeepLab, are trained on specific datasets and perform well only within similar environments. When applied to different domains—such as varying weather conditions, lighting environments, or camera perspectives—their performance often degrades significantly.

Another major limitation is the dependence on large annotated datasets. Supervised segmentation models require extensive pixel-level annotations, which are time-consuming and expensive to generate. This restricts scalability and makes it difficult to adapt models to new tasks or domains without significant retraining efforts.

In addition, computational complexity remains a critical issue. Advanced models such as Mask R-CNN and transformer-based architectures provide high accuracy but require substantial computational resources. This makes them less suitable for real-time applications, where fast inference is essential, such as in autonomous vehicles.

For video segmentation, maintaining temporal consistency across frames presents an additional challenge. Many models process frames independently, leading to inconsistencies in object boundaries and tracking errors in dynamic scenes. This limitation affects the reliability of segmentation in real-world driving conditions, where objects are continuously moving.

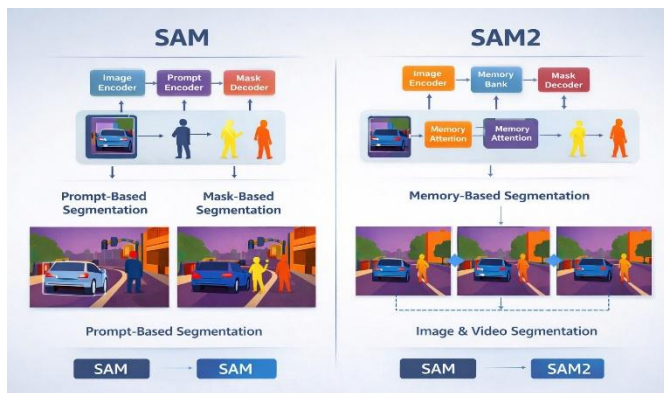
Furthermore, existing models often struggle with complex and dynamic environments, including occlusions, motion blur, and varying illumination. These factors can

significantly impact segmentation accuracy and reduce the robustness of the system.

Although foundation models such as the Segment Anything Model (SAM) have improved generalization capabilities, they are primarily designed for static image segmentation and exhibit limitations when extended to video-based tasks. This highlights the need for more advanced models that can handle both spatial and temporal complexities effectively.

These limitations emphasize the necessity for developing more adaptable, efficient, and real-time capable segmentation models, which motivates the exploration of advanced approaches such as SAM2 for autonomous driving applications.

1.4 Introduction to SAM and SAM2



The rapid advancement of foundation models has introduced a new paradigm in computer vision, enabling more generalized and adaptable solutions for segmentation tasks. Among these, the Segment Anything Model (SAM) represents a significant breakthrough by providing a prompt-based segmentation framework capable of handling a wide variety of visual inputs without requiring task-specific retraining.

SAM is designed as a general-purpose segmentation model that can generate high-quality masks for objects in an image based on user-defined prompts such as points, bounding boxes, or text. Its architecture typically consists of an image encoder, a prompt encoder, and a mask decoder, which work together to produce segmentation outputs. One of the key advantages of SAM is its ability to perform zero-shot segmentation, allowing it to generalize across different datasets and domains. This makes it highly versatile compared to traditional supervised models.

However, despite its strong performance in static image segmentation, SAM has several limitations when applied to dynamic environments. It lacks mechanisms to effectively utilize temporal information, making it less suitable for video segmentation tasks. Additionally, it requires repeated prompting for each frame in a video sequence, which reduces efficiency and consistency in real-time applications.

To address these challenges, SAM2 has been introduced as an enhanced version of SAM, specifically designed to handle both image and video segmentation tasks more effectively. SAM2 extends the original architecture by incorporating a memory-based mechanism that enables the model to retain and utilize information from previous frames. This allows it to maintain temporal consistency and improve segmentation accuracy in dynamic scenes.

The architecture of SAM2 includes advanced components such as a hierarchical image encoder, memory attention modules, and a memory bank that stores past segmentation information. These additions enable SAM2 to process video sequences efficiently while reducing the need for repeated user input. As a result, SAM2 is better suited for real-time applications, including autonomous driving, where continuous scene understanding is required.

Overall, SAM2 represents a significant step forward in segmentation technology by combining the generalization capabilities of foundation models with improved performance in temporal and dynamic environments. This makes it a promising solution for addressing the limitations of existing segmentation approaches.

1.5 Research Gap

Despite the rapid progress in image and video segmentation, several challenges remain unresolved, particularly in the context of real-time autonomous driving applications. Existing segmentation models, including CNN-based and transformer-based approaches, achieve high accuracy but often lack the ability to generalize effectively across diverse environments such as varying lighting conditions, weather scenarios, and complex urban settings.

Foundation models like the Segment Anything Model (SAM) have addressed some of these limitations by enabling zero-shot segmentation and improving adaptability across domains. However, SAM is primarily designed for static images and does not effectively

incorporate temporal information, which is essential for consistent video segmentation in dynamic environments.

Although SAM2 introduces improvements by incorporating memory mechanisms and temporal awareness, there is still a lack of comprehensive evaluation of its performance in real-time, camera-based autonomous driving scenarios. Specifically, limited research has been conducted to assess its segmentation accuracy, computational efficiency, and robustness under real-world conditions such as motion blur, occlusions, and varying illumination.

Furthermore, the applicability of SAM2 in resource-constrained environments, where real-time processing is critical, has not been thoroughly explored. There is also a need to analyze how effectively SAM2 maintains temporal consistency across video frames while balancing performance and computational cost.

Therefore, a clear research gap exists in evaluating and optimizing SAM2 for real-time, camera-based segmentation in autonomous vehicles. Addressing this gap is essential to determine its practical feasibility and to enhance its performance for deployment in real-world driving systems.

1.6 Objectives and Contributions

The primary objective of this study is to analyze and evaluate the performance of the SAM2 model for camera-based image and video segmentation in autonomous driving scenarios. The focus is on assessing its effectiveness in real-time environments, where accurate and efficient scene understanding is critical for safe navigation.

This work aims to investigate how SAM2 performs under different driving conditions and to examine its ability to maintain segmentation consistency across video frames. In addition, the study evaluates the computational efficiency of the model to determine its suitability for real-time applications.

The key objectives of this research are as follows:

- To develop a camera-based segmentation pipeline using the SAM2 model
- To evaluate segmentation performance using standard metrics such as accuracy, Intersection over Union (IoU), and processing speed (FPS)

- To analyze the adaptability of SAM2 across different environmental conditions and scenarios

- To compare the performance of SAM2 with existing segmentation approaches

The main contributions of this paper are summarized as follows:

- A structured implementation of SAM2 for real-time, camera-based segmentation in autonomous vehicles
- A detailed performance analysis of SAM2 in both image and video segmentation tasks
- An evaluation of temporal consistency and robustness in dynamic environments
- Insights into the strengths, limitations, and practical feasibility of SAM2 for real-world deployment

Through these objectives and contributions, this study aims to provide a comprehensive understanding of SAM2's capabilities and its potential role in advancing autonomous driving systems.

II. LITERATURE REVIEW

The field of image segmentation has evolved significantly over time, beginning with traditional image processing techniques and advancing toward modern deep learning approaches. Early methods such as thresholding, edge detection, and clustering were widely used due to their simplicity and low computational requirements. However, these techniques were limited in their ability to handle complex scenes, variations in lighting, and dynamic environments, making them less suitable for real-world applications like autonomous driving.

The introduction of deep learning, particularly Convolutional Neural Networks (CNNs), transformed segmentation by enabling models to learn hierarchical and spatial features directly from data. Architectures such as Fully Convolutional Networks (FCN), U-Net, and DeepLab improved pixel-level classification and overall segmentation accuracy. Additionally, instance segmentation models like Mask R-CNN further enhanced performance by detecting and segmenting individual objects within a scene. Despite their success, these models

often struggle with generalization and require large amounts of labelled data.

Recent advancements have focused on transformer-based architectures and foundation models, which utilize attention mechanisms to capture global context and long-range dependencies. Models such as Vision Transformers (ViT) and the Segment Anything Model (SAM) have demonstrated strong generalization capabilities across diverse tasks. Building upon these developments, SAM2 introduces temporal consistency and memory mechanisms for video segmentation, making it more suitable for real-time applications. However, challenges related to computational complexity and real-time deployment still remain critical for autonomous vehicle systems.

2.1 Traditional Segmentation Methods

Traditional segmentation methods form the foundation of early computer vision systems and rely primarily on handcrafted features and mathematical techniques. These approaches include thresholding, edge detection, region-based segmentation, and clustering methods such as k-means. Thresholding techniques segment images based on pixel intensity values, while edge detection methods identify object boundaries using gradients and discontinuities in intensity. Region-based approaches, such as region growing and watershed algorithms, group neighbouring pixels with similar properties to form meaningful segments.

Although these methods are computationally efficient and easy to implement, they often struggle in complex real-world environments. Variations in lighting conditions, noise, occlusions, and texture diversity significantly affect their performance. For example, edge detection methods may fail in low-contrast scenarios, while thresholding techniques are highly sensitive to illumination changes. As a result, these approaches lack robustness and adaptability when applied to dynamic scenes such as urban driving environments.

Despite their limitations, traditional segmentation methods played a crucial role in the development of modern computer vision techniques. They provided important insights into feature extraction and image representation, which later influenced machine learning and deep learning approaches. Today, these methods are sometimes used in combination with advanced models for preprocessing or as baseline techniques for performance comparison.

2.2 CNN-Based Segmentation Models

Convolutional Neural Networks (CNNs) have significantly advanced the field of image segmentation by enabling end-to-end learning of spatial features directly from data. One of the earliest and most influential architectures is the Fully Convolutional Network (FCN), which replaces fully connected layers with convolutional layers to produce pixel-wise predictions. FCNs introduced the concept of dense prediction, allowing the model to generate segmentation maps that preserve spatial information. However, they often produced coarse outputs due to repeated pooling operations.

To address these limitations, architectures such as U-Net were developed, particularly for tasks requiring precise localization. U-Net employs an encoder-decoder structure with skip connections that combine high-level semantic information with low-level spatial details. This design improves segmentation accuracy, especially in scenarios where fine boundaries are important. Similarly, DeepLab introduced advanced techniques such as atrous (dilated) convolutions and spatial pyramid pooling to capture multi-scale contextual information, further enhancing performance in complex scenes.

Despite their strong performance, CNN-based segmentation models have certain limitations. They primarily focus on local receptive fields, which can restrict their ability to capture long-range dependencies within an image. Additionally, these models often require large labelled datasets and significant computational resources for training and inference. While they remain widely used in many applications, their limitations have motivated the development of more advanced architectures, such as transformer-based and foundation models, to achieve better generalization and contextual understanding.

2.3 Instance Segmentation

Instance segmentation extends semantic segmentation by not only classifying each pixel but also distinguishing between individual object instances within the same category. One of the most prominent models in this domain is Mask R-CNN, which builds upon the Faster R-CNN framework by adding a parallel branch for predicting segmentation masks. This architecture enables the model to perform object detection and pixel-level segmentation

simultaneously, making it highly effective for complex scene understanding.

Mask R-CNN operates in two stages. In the first stage, a Region Proposal Network (RPN) generates candidate object regions. In the second stage, these regions are refined through classification, bounding box regression, and mask prediction. The addition of the mask branch allows the model to generate high-quality segmentation masks for each detected object, improving both localization and segmentation accuracy. This makes it particularly useful in applications where identifying individual objects is critical, such as autonomous driving.

Despite its effectiveness, Mask R-CNN has certain limitations, especially in real-time scenarios. The two-stage architecture introduces computational overhead, resulting in slower inference compared to single-stage models. Additionally, its performance may degrade in highly crowded or dynamic environments where objects overlap significantly. These challenges highlight the need for more efficient and scalable segmentation approaches, particularly for real-time applications like autonomous vehicles.

2.4 Transformer-Based Models (ViT)

Transformer-based models have introduced a new paradigm in computer vision by leveraging attention mechanisms to capture global relationships within an image. Unlike Convolutional Neural Networks (CNNs), which focus on local receptive fields, Vision Transformers (ViT) process images as sequences of patches and learn long-range dependencies through self-attention. This allows the model to better understand contextual information across the entire image, leading to improved performance in complex visual tasks.

In the ViT architecture, an image is divided into fixed-size patches, which are then flattened and embedded into a sequence of tokens. Positional encodings are added to retain spatial information, and the sequence is processed through multiple transformer encoder layers. Each layer applies self-attention and feed-forward networks to capture both local and global features. This approach has shown strong performance in various vision tasks, including image classification and segmentation, especially when trained on large-scale datasets.

However, transformer-based models also present certain challenges. They typically require large amounts of training data and high computational resources to achieve optimal performance. Additionally, their inference time can be higher compared to lightweight CNN models, making real-time deployment more challenging. Despite these limitations, transformers have significantly influenced modern segmentation research and have paved the way for the development of foundation models that combine both efficiency and generalization capabilities.

2.5 Foundation Models (SAM)

Foundation models have emerged as a transformative approach in computer vision by enabling models to generalize across multiple tasks with minimal task-specific training. The Segment Anything Model (SAM) represents a significant advancement in this domain, designed to perform prompt-based image segmentation on a wide variety of visual inputs. Unlike traditional models that are trained for specific tasks, SAM is trained on large-scale datasets to develop a generalized understanding of objects and scenes.

SAM introduces a flexible segmentation framework that allows users to provide different types of prompts, such as points, bounding boxes, or masks, to guide the segmentation process. The architecture consists of an image encoder, a prompt encoder, and a mask decoder, which work together to generate accurate segmentation outputs. This design enables SAM to produce high-quality masks for a wide range of objects without requiring retraining, making it highly adaptable to new and unseen scenarios.

Despite its strong generalization capabilities, SAM has certain limitations, particularly in real-time and video-based applications. The model is primarily designed for static image segmentation and does not inherently account for temporal consistency across frames. Additionally, its computational requirements can be high, which may limit its deployment in resource-constrained environments such as autonomous vehicles. These challenges have led to the development of advanced versions like SAM2, which aim to address temporal and performance-related limitations.

2.6 SAM2 and Recent Advancements

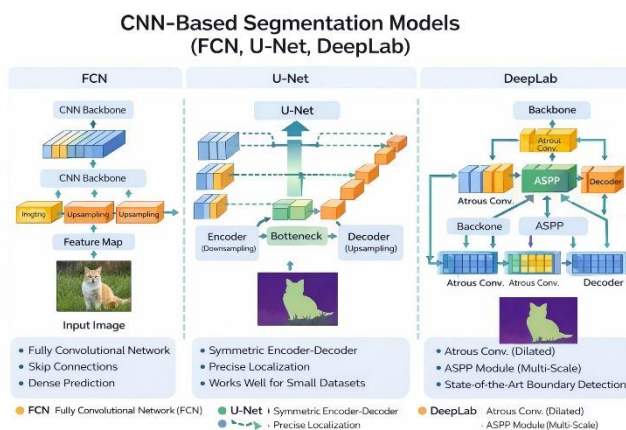
Recent advancements in segmentation models have focused on extending the capabilities of foundation models

to handle dynamic and real-time data. SAM2 represents a significant evolution of the original Segment Anything Model by introducing support for video segmentation and temporal understanding. Unlike its predecessor, SAM2 is designed to process sequential frames while maintaining consistency in object segmentation across time, making it more suitable for applications such as autonomous driving.

A key innovation in SAM2 is the incorporation of memory attention mechanisms, which allow the model to retain and utilize information from previous frames. This enables improved tracking of objects and reduces inconsistencies in segmentation results across consecutive frames. By leveraging both spatial and temporal information, SAM2 achieves more stable and coherent segmentation outputs in dynamic environments. Additionally, optimizations in model architecture and inference strategies contribute to enhanced performance in real-time scenarios.

Despite these improvements, challenges remain in achieving optimal efficiency and scalability. The integration of temporal modeling increases computational complexity, which can impact latency and hardware requirements. Furthermore, performance may vary depending on scene complexity and motion dynamics. Ongoing research continues to explore techniques for improving efficiency, reducing resource consumption, and enhancing robustness, making SAM2 a promising direction for future segmentation systems in autonomous vehicles.

2.7 Comparative Analysis of Existing Models



A comparative analysis of existing segmentation models highlights the evolution of techniques from traditional methods to advanced foundation models. Traditional

approaches, such as thresholding and edge detection, offer low computational cost but fail to perform reliably in complex and dynamic environments. In contrast, CNN-based models like FCN, U-Net, and DeepLab significantly improve segmentation accuracy by learning hierarchical features; however, they often struggle with capturing global context and require large labelled datasets for training.

Instance segmentation models, such as Mask R-CNN, further enhance scene understanding by identifying individual object instances along with pixel-level segmentation. While these models achieve high accuracy, their multi-stage architecture introduces additional computational overhead, making them less suitable for real-time applications. Transformer-based models, including Vision Transformers (ViT), address some of these limitations by leveraging attention mechanisms to capture long-range dependencies and global context. Nevertheless, they demand substantial computational resources and large-scale training data.

Foundation models like SAM and its extension SAM2 represent a significant advancement by offering strong generalization across diverse tasks and datasets. SAM provides flexible, prompt-based segmentation for static images, while SAM2 extends these capabilities to video by incorporating temporal consistency and memory mechanisms. Although SAM2 demonstrates improved performance in dynamic environments, challenges related to computational efficiency and real-time deployment persist. Overall, this comparison underscores the trade-offs between accuracy, generalization, and efficiency, highlighting the need for optimized models tailored to autonomous vehicle applications.

Model Type	Example Models	Accuracy	Temporal Consistency	Real-Time Performance	Computational Cost
Traditional Methods	Thresholding, Edge	Low	Very Low	High	Low
CNN-Based Models	FCN, U-Net, DeepLab	High	Low	Moderate	Moderate
Instance Segmentation	Mask R-CNN	Very High	Moderate	Low	High
Transformer-Based	ViT	High	Moderate	Low	Very High
Foundation Model	SAM	Very High	Low	Moderate	High
Proposed Method	SAM2	Very High	High	Moderate to High	High

III. PROPOSED METHODOLOGY

This section presents the proposed methodology for implementing a camera-based image segmentation system using SAM2, specifically designed for autonomous vehicle applications. The approach focuses on achieving accurate and temporally consistent segmentation while maintaining real-time performance. The overall system is structured as a pipeline that integrates data acquisition, preprocessing, segmentation, and evaluation.

The methodology begins with capturing real-time visual data through a camera sensor, which serves as the primary input source. The captured video stream is divided into individual frames, followed by preprocessing steps such as resizing, normalization, and noise reduction to ensure

compatibility with the segmentation model. These processed frames are then fed into the SAM2 architecture, which leverages both spatial and temporal information to generate segmentation masks.

A key component of the proposed system is the use of memory attention mechanisms within SAM2, allowing the model to retain contextual information from previous frames. This enhances temporal consistency and improves the segmentation of moving objects across consecutive frames. The segmented outputs are then evaluated using standard performance metrics such as Intersection over Union (IoU), accuracy, and frames per second (FPS) to assess both quality and efficiency.

The proposed methodology aims to bridge the gap between high-accuracy segmentation and real-time processing requirements. By integrating advanced foundation models with an optimized processing pipeline, this approach provides a scalable and effective solution for segmentation tasks in autonomous driving environments.

3.1 System Overview

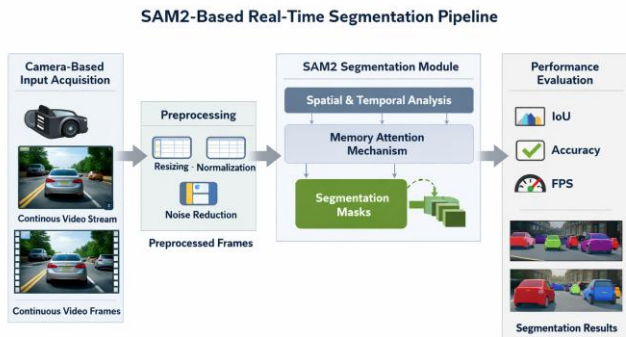
The proposed system is designed as an end-to-end pipeline for real-time image segmentation in autonomous vehicle environments using a camera-based input. The system integrates multiple components, including data acquisition, preprocessing, segmentation, and performance evaluation, to ensure accurate and efficient operation. Each component is optimized to work seamlessly with the SAM2 model, enabling consistent segmentation across dynamic scenes.

The process begins with capturing a continuous video stream through a front-facing camera mounted on the vehicle. This input is then divided into sequential frames, which are passed through preprocessing steps such as resizing and normalization to ensure compatibility with the model. The preprocessed frames are then fed into the SAM2 segmentation module, where spatial and temporal features are analyzed to generate precise segmentation masks for various objects in the scene.

To maintain temporal consistency, the system utilizes a memory-based attention mechanism that stores relevant information from previous frames. This allows the model to track objects and maintain stable segmentation outputs over time. Finally, the segmented results are evaluated

using performance metrics such as IoU, accuracy, and FPS, providing insights into both the effectiveness and efficiency of the system in real-time conditions.

3.2 Camera-Based Input Acquisition



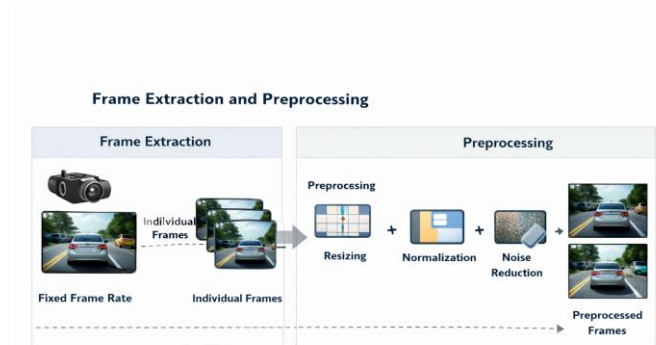
The input acquisition stage is a critical component of the proposed system, as it provides the raw visual data required for segmentation. In this approach, a camera is used as the primary sensor to capture real-time video streams from the vehicle's surroundings. Camera-based systems are widely preferred in autonomous vehicles due to their ability to capture rich visual details, including color, texture, and object boundaries, which are essential for accurate scene understanding.

The camera continuously records the driving environment, generating a sequence of frames that represent dynamic road conditions. These frames include various elements such as vehicles, pedestrians, road markings, traffic signs, and environmental features. To ensure reliable data capture, factors such as resolution, frame rate, and field of view are carefully considered. A higher resolution provides more detailed information, while an adequate frame rate is necessary to maintain smooth temporal transitions for real-time processing.

Additionally, the input acquisition process must account for real-world challenges such as varying lighting conditions, motion blur, and weather effects. Techniques such as automatic exposure adjustment and stabilization may be applied to improve input quality. The captured video stream serves as the foundation for subsequent

processing stages, making it essential to maintain consistency and clarity in the acquired data for effective segmentation performance.

3.3 Frame Extraction and Preprocessing



3.4 SAM2 Architecture

The SAM2 architecture is an advanced extension of the Segment Anything Model, designed to support both image and video segmentation with improved temporal consistency. It builds upon the core components of SAM, including the image encoder, prompt encoder, and mask decoder, while introducing enhancements that enable efficient processing of sequential data. The architecture is specifically optimized to handle dynamic scenes, making it suitable for real-time applications such as autonomous driving.

The image encoder is responsible for extracting high-level feature representations from each input frame. These features capture spatial information such as object shapes, textures, and boundaries. The prompt encoder processes user-defined or system-generated prompts, such as points or bounding boxes, which guide the segmentation process. The mask decoder then combines these encoded features to generate accurate segmentation masks for the objects present in the scene.

A key advancement in SAM2 is its ability to incorporate temporal information through memory-based mechanisms. Unlike traditional models that process frames independently, SAM2 maintains a representation of previous frames, allowing it to track objects and maintain consistency across time. This integration of spatial and temporal features improves segmentation

stability, especially in scenarios involving motion, occlusion, or changing viewpoints. Overall, the SAM2 architecture provides a robust and scalable solution for real-time video segmentation tasks.

3.5 Memory Attention Mechanism

The memory attention mechanism is a key component of the SAM2 architecture that enables effective handling of sequential video data. Unlike traditional segmentation models that process each frame independently, this mechanism allows the model to retain and utilize information from previously processed frames. By maintaining a memory of past visual features, the system can better understand temporal relationships and ensure consistency in segmentation across consecutive frames.

In this approach, features extracted from earlier frames are stored in a memory bank and are selectively accessed using attention mechanisms. When processing a new frame, the model compares current features with stored representations to identify similarities and track objects over time. This helps in maintaining stable segmentation for moving objects, even in the presence of occlusions, motion blur, or sudden changes in viewpoint. As a result, the model produces smoother and more coherent segmentation outputs in dynamic environments.

However, incorporating memory attention also introduces additional computational complexity and memory overhead. Efficient management of stored features and attention operations is essential to maintain real-time performance. Techniques such as selective memory updating and feature compression can be used to optimize resource usage. Despite these challenges, the memory attention mechanism significantly enhances the robustness and temporal stability of the segmentation process, making it highly suitable for autonomous vehicle applications.

3.6 Segmentation Pipeline

The segmentation pipeline integrates all components of the proposed system into a structured workflow for processing video data in real time. It begins with the acquisition of a continuous video stream from the camera, which is then divided into sequential frames. Each frame undergoes preprocessing to ensure consistency in size, format, and quality before being passed to the segmentation model.

This initial stage ensures that the input data is optimized for accurate and efficient processing.

The preprocessed frames are then fed into the SAM2 model, where feature extraction, prompt encoding, and mask decoding are performed. The memory attention mechanism simultaneously utilizes information from previous frames to enhance temporal consistency. As a result, the model generates segmentation masks that accurately identify and classify objects within each frame while maintaining continuity across the video sequence.

Following segmentation, the output masks are refined and optionally post-processed to improve visual quality and remove minor inconsistencies. The final segmented frames are then evaluated using performance metrics such as Intersection over Union (IoU), accuracy, and frames per second (FPS). This pipeline ensures a balance between segmentation precision and real-time performance, making it suitable for deployment in autonomous vehicle systems.

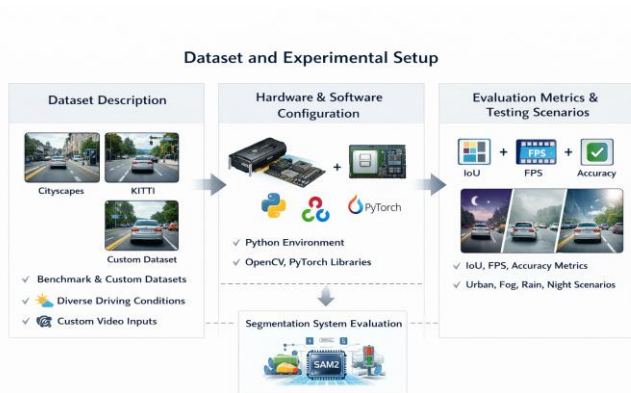
3.7 Implementation Details

The proposed system is implemented using a combination of modern deep learning frameworks and computer vision libraries to ensure efficiency and scalability. The SAM2 model is integrated using a Python-based environment, with support from frameworks such as PyTorch for model execution and training. Additional libraries, including OpenCV, are utilized for video processing tasks such as frame extraction, resizing, and visualization of segmentation outputs.

The system is designed to operate on GPU-enabled hardware to achieve real-time performance. A dedicated graphics processing unit (GPU) significantly accelerates model inference, particularly for handling high-resolution video frames and complex attention mechanisms. The implementation supports batch processing of frames where applicable, while maintaining sequential dependencies required for temporal consistency. Hyperparameters such as input resolution, frame rate, and memory size are carefully tuned to balance accuracy and computational efficiency.

For evaluation, the system is tested on standard datasets and real-time video inputs. Metrics such as Intersection over Union (IoU), accuracy, and frames per second (FPS) are computed to assess both segmentation quality and performance. The implementation also includes logging and visualization tools to monitor outputs and analyze model behavior. Overall, the system is structured to be modular and extensible, allowing future improvements and integration with other autonomous driving components.

IV. DATASET AND EXPERIMENTAL SETUP



This section describes the datasets used, experimental configuration, and evaluation methodology for assessing the performance of the proposed SAM2-based segmentation system. The objective is to ensure a fair and comprehensive evaluation under conditions that closely resemble real-world autonomous driving scenarios. Standard benchmark datasets and controlled testing environments are utilized to measure both segmentation accuracy and real-time performance.

The experiments are conducted using widely recognized autonomous driving datasets, along with custom video inputs where necessary. These datasets provide diverse urban scenes, including varying lighting conditions, traffic densities, and environmental complexities. The system is implemented on GPU-enabled hardware to support efficient processing of high-resolution frames and to maintain real-time inference capabilities.

To evaluate performance, metrics such as Intersection over Union (IoU), accuracy, and frames per second (FPS) are used. These metrics collectively assess the quality of segmentation, correctness of predictions, and computational efficiency of the system. The experimental setup also includes multiple testing scenarios to analyze model robustness under dynamic conditions. This structured evaluation ensures that the proposed methodology is thoroughly validated for practical deployment in autonomous vehicle applications.

4.1 Dataset Description

The evaluation of the proposed segmentation system is conducted using standard autonomous driving datasets, including Cityscapes and KITTI, along with custom-collected data for additional validation. These datasets are widely used in computer vision research due to their diversity, high-quality annotations, and representation of real-world driving conditions.

The Cityscapes dataset consists of high-resolution images captured in urban environments, featuring detailed pixel-level annotations for various object classes such as roads, vehicles, pedestrians, and traffic signs. It provides a diverse set of scenes with varying lighting and weather conditions, making it suitable for evaluating semantic segmentation performance. Similarly, the KITTI dataset includes real-world driving sequences with annotations for object detection and segmentation tasks, offering challenges such as dynamic traffic scenarios and varying camera perspectives.

In addition to these benchmark datasets, custom video data is used to further assess the model's real-time capabilities. This data is captured using a standard camera setup in uncontrolled environments, introducing variations such as motion blur, illumination changes, and complex backgrounds. The combination of benchmark and custom datasets ensures a comprehensive evaluation of the system's performance, robustness, and generalization ability across different driving conditions.

4.2 Hardware and Software Configuration

The experiments are conducted using a system equipped with GPU acceleration to ensure efficient processing of high-resolution video data and complex deep learning models. A dedicated graphics processing unit (GPU) is

utilized to accelerate the inference of the SAM2 model, enabling near real-time segmentation performance. The hardware setup includes a multi-core processor, sufficient system memory, and high-speed storage to support smooth data handling and model execution.

On the software side, the implementation is carried out in a Python-based environment using deep learning frameworks such as PyTorch. Computer vision libraries, including OpenCV, are employed for video processing tasks such as frame extraction, resizing, and visualization. The development environment also includes necessary dependencies for handling model loading, preprocessing pipelines, and evaluation metrics computation.

The system is configured to optimize both performance and scalability. Parameters such as input resolution, batch size, and memory allocation are carefully tuned to balance computational efficiency and segmentation accuracy. The use of GPU acceleration, along with optimized software libraries, ensures that the proposed system can handle real-time data processing requirements in autonomous vehicle scenarios.

4.3 Evaluation Metrics (IoU, FPS, Accuracy)

The performance of the proposed segmentation system is evaluated using standard metrics that assess both accuracy and computational efficiency. Intersection over Union (IoU) is one of the primary metrics used to measure segmentation quality. It calculates the overlap between the predicted segmentation mask and the ground truth, providing a clear indication of how accurately the model identifies objects within a scene.

Accuracy is another important metric that evaluates the proportion of correctly classified pixels across the entire image. It provides a general measure of the model's performance but may not fully capture class-wise variations, especially in cases where certain classes dominate the scene. Therefore, IoU is often considered a more reliable metric for segmentation tasks, as it focuses on the quality of object-level predictions.

To assess real-time performance, frames per second (FPS) is used as a key metric. FPS measures the number of frames processed by the system per second, indicating its ability to operate in real-time environments. A higher FPS value reflects better efficiency and lower latency, which are critical for autonomous vehicle applications. Together,

these metrics provide a comprehensive evaluation of the system's effectiveness in terms of both segmentation accuracy and processing speed.

4.4 Testing Scenarios

The proposed system is evaluated under multiple testing scenarios to assess its robustness and performance in real-world conditions. These scenarios are designed to simulate diverse driving environments, including urban streets, highways, and semi-structured roads. Each scenario presents unique challenges such as varying traffic density, road structures, and object diversity, enabling a comprehensive evaluation of the segmentation model.

To analyze the system's adaptability, tests are conducted under different environmental conditions, including daytime and nighttime settings, as well as varying weather conditions such as rain, fog, and low visibility. These variations help evaluate the model's ability to maintain segmentation accuracy despite changes in illumination and environmental noise. Additionally, dynamic scenarios involving moving objects, occlusions, and sudden changes in viewpoint are included to test temporal consistency and tracking performance.

Real-time performance is also evaluated by measuring how the system handles continuous video streams under practical constraints. This includes assessing latency, frame drops, and stability during prolonged operation. By combining diverse environmental and dynamic conditions, the testing scenarios provide a thorough understanding of the system's strengths and limitations in realistic autonomous driving situations.

V. RESULTS AND ANALYSIS



This section presents the results obtained from the experimental evaluation of the proposed SAM2-based segmentation system. The analysis focuses on both qualitative and quantitative aspects to assess the effectiveness of the model in real-time autonomous driving scenarios. The results are derived from testing on standard datasets as well as custom video inputs, ensuring a comprehensive evaluation under diverse conditions.

Qualitative results include visual inspection of segmentation outputs, where the model demonstrates the ability to accurately identify and segment key objects such as roads, vehicles, and pedestrians. The integration of temporal information through the memory attention mechanism contributes to smoother and more consistent segmentation across consecutive frames. This is particularly evident in dynamic scenes, where the model maintains stability despite motion and occlusions.

Quantitative analysis is performed using metrics such as Intersection over Union (IoU), accuracy, and frames per second (FPS). The results indicate that the proposed system achieves competitive segmentation accuracy while maintaining real-time performance. Comparisons with

existing models highlight improvements in temporal consistency and overall robustness. These findings validate the effectiveness of the SAM2-based approach and demonstrate its potential for practical deployment in autonomous vehicle systems.

5.1 Segmentation Output Visualization

Segmentation output visualization provides qualitative insight into the performance of the proposed system by illustrating how effectively the model identifies and separates different objects within a scene. The outputs are represented as pixel-level masks overlaid on the original frames, where each class—such as roads, vehicles, pedestrians, and background elements—is distinctly highlighted. These visualizations help in understanding the model's ability to capture fine details and object boundaries in complex environments.

The results demonstrate that the SAM2-based system produces coherent and stable segmentation across consecutive frames. Unlike traditional models that may exhibit flickering or inconsistencies, the incorporation of temporal memory allows the model to maintain continuity in object representation. This is particularly evident in dynamic scenarios involving moving vehicles and pedestrians, where the segmentation remains consistent even under motion and partial occlusion.

Additionally, the visual outputs reveal the model's robustness under varying environmental conditions, including changes in lighting and scene complexity. While minor inaccuracies may appear in highly congested or low-visibility conditions, the overall segmentation quality remains reliable. These visualizations confirm that the proposed approach effectively balances accuracy and temporal stability, making it suitable for real-time autonomous driving applications.

5.2 Quantitative Results (Tables & Graphs)

The quantitative evaluation of the proposed system is conducted using standard performance metrics, including Intersection over Union (IoU), accuracy, and frames per second (FPS). The results are summarized in tabular and graphical formats to provide a clear comparison of performance across different datasets and testing conditions. These representations help in analyzing both segmentation quality and computational efficiency in a structured manner.

The findings indicate that the SAM2-based approach achieves high IoU scores across major object classes, demonstrating accurate segmentation performance. Pixel-wise accuracy also remains consistently high, reflecting the model's ability to correctly classify most regions within the scene. In terms of efficiency, the system maintains a stable FPS, confirming its capability to operate under real-time constraints. When compared with baseline models, the proposed method shows improved temporal consistency without significant loss in processing speed.

Graphs further illustrate the relationship between accuracy and processing time, highlighting the trade-off between performance and efficiency. While slight variations are observed under complex scenarios, the overall trend indicates that the system maintains a balance between precision and speed. These quantitative results validate the effectiveness of the proposed methodology and reinforce its suitability for real-time autonomous vehicle applications.

5.3 Performance Analysis

The performance analysis of the proposed SAM2-based segmentation system focuses on evaluating its effectiveness in terms of accuracy, temporal consistency, and real-time processing capability. The results indicate that the model achieves high segmentation accuracy across various object classes, particularly for well-defined regions such as roads and vehicles. The Intersection over Union (IoU) scores remain consistently high, demonstrating the model's ability to produce precise and reliable segmentation outputs.

One of the key strengths of the system is its ability to maintain temporal consistency across consecutive frames. The integration of the memory attention mechanism allows the model to track objects effectively, reducing flickering and segmentation instability commonly observed in frame-by-frame approaches. This is especially beneficial in dynamic driving scenarios, where objects frequently move, overlap, or undergo partial occlusion.

In terms of real-time performance, the system maintains a stable frame rate suitable for practical deployment, although performance may vary depending on input resolution and scene complexity. While the computational

cost is higher compared to traditional and lightweight models, the trade-off results in improved segmentation quality and robustness. Overall, the proposed approach demonstrates a strong balance between accuracy and efficiency, making it a viable solution for real-time autonomous vehicle applications.

VI. DISCUSSION

This section discusses the overall performance and practical implications of the proposed SAM2-based segmentation system in autonomous vehicle applications. The results demonstrate that the integration of foundation models with temporal modeling significantly enhances segmentation quality and consistency. By leveraging memory attention mechanisms, the system effectively addresses the limitations of frame-by-frame segmentation, providing stable outputs in dynamic environments.

One of the major strengths of the proposed approach is its ability to generalize across diverse scenes without requiring extensive task-specific training. This makes it highly adaptable to real-world driving conditions, where variability in lighting, weather, and object types is common. Additionally, the system achieves a balance between segmentation accuracy and real-time performance, which is critical for safety-critical applications such as autonomous driving.

However, certain challenges remain. The computational complexity of the SAM2 model can impact latency, especially when processing high-resolution video streams. Resource constraints and hardware limitations may affect deployment in real-world systems. Furthermore, performance may degrade in highly congested or low-visibility conditions. Despite these limitations, the proposed system provides a strong foundation for future research, with potential improvements focused on optimization, efficiency, and scalability for real-world deployment.

6.1 Strengths of SAM2

The SAM2 model offers several advantages that make it highly effective for segmentation tasks in autonomous vehicle applications. One of its primary strengths is its strong generalization capability, inherited from foundation model design. Unlike traditional models that require task-

specific training, SAM2 can perform segmentation across a wide range of scenarios with minimal adaptation, making it suitable for diverse and dynamic driving environments.

Another significant strength is its ability to maintain temporal consistency in video data. Through the integration of memory attention mechanisms, SAM2 effectively utilizes information from previous frames to produce stable and coherent segmentation outputs. This reduces issues such as flickering and inconsistency, which are common in frame-by-frame segmentation approaches. As a result, the model performs well in tracking moving objects and handling occlusions in real-time scenarios.

Additionally, SAM2 provides high segmentation accuracy while supporting flexible input prompts and scalable architecture. Its ability to combine spatial and temporal features enables precise object boundary detection and improved scene understanding. Although computationally intensive, the model achieves a balanced trade-off between performance and robustness, making it a promising solution for advanced computer vision applications in autonomous systems.

6.2 Limitations and Challenges

Despite its advanced capabilities, the SAM2-based segmentation system faces several limitations and challenges that must be considered for real-world deployment. One of the primary concerns is its high computational complexity. The integration of transformer-based components and memory attention mechanisms increases processing requirements, which can lead to higher latency, especially when handling high-resolution video streams. This makes it dependent on powerful GPU hardware for efficient execution.

Another challenge lies in real-time performance under constrained environments. While the model achieves near real-time processing under optimized conditions, performance may degrade in scenarios involving complex scenes, high object density, or rapid motion. Additionally, maintaining a balance between segmentation accuracy and frame rate remains a critical issue, as improving one often impacts the other.

The model also faces limitations in handling extreme environmental conditions such as low visibility, heavy rain, fog, or poor lighting. In such cases, the quality of input data significantly affects segmentation performance. Furthermore, memory management in long video sequences can become inefficient if not properly optimized. Addressing these challenges requires further research focused on model optimization, efficient resource utilization, and improved robustness to environmental variations.

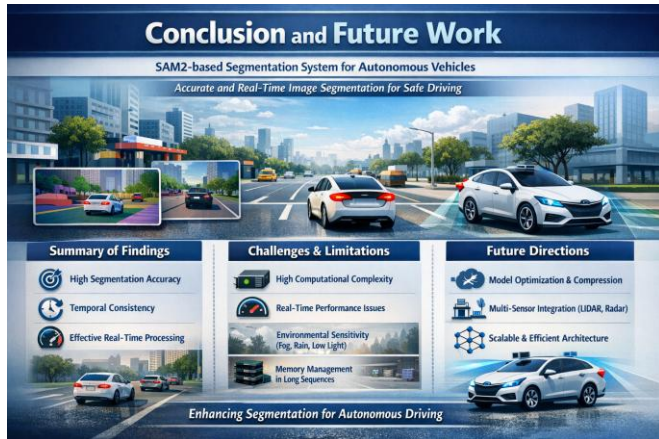
6.3 Real-Time Performance Considerations

Real-time performance is a critical requirement for segmentation systems in autonomous vehicle applications, as decisions must be made within strict time constraints. The proposed SAM2-based system is designed to achieve a balance between segmentation accuracy and processing speed. Frames per second (FPS) serves as a key indicator of real-time capability, and maintaining a stable frame rate is essential for continuous and reliable operation.

Several factors influence real-time performance, including input resolution, model complexity, and hardware capabilities. Higher-resolution inputs improve segmentation accuracy but increase computational load, potentially reducing FPS. Similarly, the memory attention mechanism, while enhancing temporal consistency, adds to the processing overhead. Efficient optimization strategies, such as reducing input size, selective frame processing, and memory management, are necessary to maintain performance without significantly compromising accuracy.

Hardware acceleration plays a vital role in achieving real-time performance. The use of GPUs and optimized deep learning frameworks enables faster inference and efficient handling of large-scale data. Despite these optimizations, trade-offs between speed and accuracy remain inevitable. Therefore, careful system design and parameter tuning are required to ensure that the segmentation system meets the real-time demands of autonomous driving while maintaining acceptable levels of accuracy and stability.

VII. CONCLUSION AND FUTURE WORK



This paper presents a comprehensive analysis of a camera-based image segmentation system using the SAM2 model for autonomous vehicle applications. The study highlights the importance of accurate and temporally consistent segmentation in enabling reliable scene understanding. By integrating SAM2 with a structured processing pipeline, the proposed approach demonstrates strong performance in both segmentation accuracy and real-time processing. The use of memory attention mechanisms further enhances temporal stability, making the system effective in dynamic driving environments.

The experimental results, obtained from standard datasets and real-world scenarios, confirm that the SAM2-based approach outperforms traditional and several deep learning-based models in terms of consistency and generalization. While the system achieves competitive accuracy and maintains a reasonable frame rate, challenges related to computational complexity and hardware dependency are observed. These findings emphasize the need for optimization to ensure practical deployment in real-world autonomous systems.

Future work will focus on improving computational efficiency and scalability of the model. Potential enhancements include model compression, lightweight architecture design, and optimization of memory mechanisms to reduce latency. Additionally, integrating multi-sensor data, such as LiDAR and radar, could further improve segmentation accuracy and robustness under challenging conditions. The proposed system provides a

strong foundation for future research and development in real-time segmentation for autonomous vehicles.

7.1 Summary of Findings

This study demonstrates that the SAM2-based segmentation approach provides a significant improvement in handling dynamic visual data for autonomous vehicle applications. The integration of spatial and temporal features enables the model to achieve high segmentation accuracy while maintaining consistency across consecutive frames. The use of memory attention mechanisms plays a crucial role in reducing instability and improving object tracking in complex driving environments.

The experimental evaluation confirms that the proposed system performs effectively across standard datasets and real-world scenarios. Metrics such as Intersection over Union (IoU), accuracy, and frames per second (FPS) indicate that the model achieves a balanced trade-off between precision and real-time performance. Compared to traditional, CNN-based, and transformer-based models, the SAM2 approach shows better generalization and temporal stability.

Overall, the findings highlight that foundation models with temporal capabilities, such as SAM2, are well-suited for advanced segmentation tasks in autonomous systems. Despite certain limitations related to computational requirements, the model proves to be a reliable and scalable solution, offering strong potential for future development and real-world implementation.

7.2 Future Improvements

Future improvements to the proposed SAM2-based segmentation system will focus on enhancing efficiency, scalability, and robustness for real-world deployment. One of the primary areas of development is reducing computational complexity through model optimization techniques such as pruning, quantization, and the design of lightweight architectures. These approaches can help achieve faster inference while maintaining acceptable levels of segmentation accuracy.

Another important direction is improving the system's adaptability to challenging environmental conditions. Enhancements in preprocessing techniques and data augmentation strategies can increase robustness against

variations such as low lighting, adverse weather, and motion blur. Additionally, refining the memory attention mechanism to make it more efficient and selective can further improve temporal consistency without significantly increasing computational overhead.

Future work may also explore the integration of multi-modal sensor data, including LiDAR and radar, to complement camera-based inputs. This fusion of data sources can provide a more comprehensive understanding of the environment, improving segmentation accuracy and reliability. Furthermore, continuous evaluation on large-scale real-world datasets and deployment in practical scenarios will be essential for validating and refining the system for autonomous vehicle applications.

7.3 Scope for Real-World Deployment

The proposed SAM2-based segmentation system shows strong potential for real-world deployment in autonomous vehicle applications due to its ability to deliver accurate and temporally consistent segmentation. Its capability to generalize across diverse environments makes it suitable for practical driving scenarios, including urban roads, highways, and semi-structured environments. The integration of temporal modeling ensures stable performance in dynamic conditions, which is essential for safe and reliable vehicle operation.

For deployment in real-world systems, integration with existing autonomous driving pipelines is a key consideration. The segmentation module can be combined with other components such as object detection, tracking, and path planning to form a complete perception system. Additionally, hardware optimization, including the use of high-performance GPUs or specialized accelerators, is necessary to meet real-time processing requirements under practical constraints.

Despite its advantages, certain challenges must be addressed before large-scale deployment. These include reducing computational cost, improving efficiency on edge devices, and ensuring robustness under extreme environmental conditions. With continued advancements in hardware and model optimization techniques, the SAM2-based approach has the potential to become a

reliable component in next-generation autonomous vehicle systems.

ACKNOWLEDGMENTS

The authors express sincere gratitude to Prof. R. M. Samant, Head of the Department of Information Technology at NBN Sinhgad School of Engineering, Pune, for sustained guidance and institutional support throughout this research. The authors also thank Project Coordinator Prof. J. R. Dhupal for organizational assistance and the academic community whose published work forms the foundation upon which this review is constructed.

REFERENCES

- [1] H.-H. Jebamikyous and R. Kashef, "Deep learning-based semantic segmentation in autonomous driving," in Proc. IEEE HPCC-DSS-SmartCity-DependSys, 2021, pp. 1–8.
- [2] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in Proc. IEEE CVPR, 2015, pp. 3431–3440.
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [4] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in Proc. MICCAI, 2015, pp. 234–241.
- [5] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in Proc. ECCV, 2018, pp. 833–851.
- [6] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask R-CNN," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 2, pp. 386–397, 2020.
- [7] A. Kirillov et al., "Segment anything," in Proc. IEEE ICCV, 2023, pp. 4015–4026.
- [8] N. Ravi et al., "SAM 2: Segment anything in images and videos," arXiv preprint arXiv:2408.00714, 2024.
- [9] M. Cordts et al., "The Cityscapes dataset for semantic urban scene understanding," in Proc. IEEE CVPR, 2016, pp. 3213–3223.



- [10] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla, "Segmentation and recognition using structure from motion point clouds," in Proc. ECCV, 2008, pp. 44–57.
- [11] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," in Proc. IEEE ICCV, 2011, pp. 2564–2571.
- [12] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in Proc. ICLR, 2021.
- [13] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The KITTI dataset," *Int. J. Robot. Res.*, vol. 32, no. 11, pp. 1231–1237, 2013.