# Machine Learning for Autonomous Vehicle Environment Perception and Analysis

*By Dr. Carlos Murillo*

*Professor of Industrial Engineering, Universidad Nacional Autónoma de México (UNAM)*

## 1. Introduction to Autonomous Vehicles and Environment Perception

Basically, this paper will cover the perception tasks from both sensor and algorithm perspectives and we will also discuss the current industry and the open problems in AV perception from both research and application perspectives. Based on this study, important views of possible future research topics in camera-radar perception for AVs will be pointed out. This survey is structured in six main sections. The first section is the introduction, the second section includes the surve literature on environmental perception. The surveys are done, in detail, for each sensor such as Cameras, LiDARs, Radars, and. More specifically, this section includes LiDAR surveys and compares the ultimate approach applied for perception. Then, the Radar surveys are presented in detail according to their perception algorithms. Then, Camera-Radar fusion surveys are included. Another sub-section is a very active and beginning trend in this survey. Basic topics are discussed, such as background subtraction for static background especially due to the illumination changes, and the second topic is illumination invariant moving object detection and tracking. Finally, related works and the discussion section is in section four. A very detailed future research direction section is also provided in this section. Section five is the conclusion.

[1] [2] [3]The development of autonomous vehicles (AVs) is facing numerous challenges, including legal, regulatory, cost, and technological issues. On the technological side, the most important task for driving and functioning within the traffic environment is Environment Perception (EP). EP is responsible for identifying all objects around the vehicle and analyzing, identifying, and understanding the environment according to these objects. This forms the basis for subsequent environmental planning and decision-making according to intelligent control. This paper aims to provide readers with a state-of-the-art survey of the perception and analysis task for AVs at a very detailed level. Instead of trying to make an extensive and

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

brief survey covering all perception tasks, this survey mainly focuses on the main perception scenarios including obstacle (e.g., vehicle, pedestrian, cyclist) detection and road segmentation.

## 1.1. Overview of Autonomous Vehicle Technology

The perception module (also known as the Sensing & Perception module) supplies the decision-making and planning modules of an AV with intelligent data, which are required for making decisions and performing planning and control actions. An enhancement in the performance and capabilities of the perception module will introduce a notable improvement in the overall AV performance. Currently, multiple sensors and communication-based mechanisms are jointly used in the perception module to enhance the perception results. This module plays an important role in ensuring the general availability of intelligent data for the overall autonomous vehicle (AV) architecture in different scenarios, including complex on-road and off-road environments, where the physical constraints and dynamics of the environment may require further intelligence in the AV perception and decision-making.

[4] [5] Autonomous vehicles have revolutionized future mobility by attracting notable attention from manufacturers and researchers. Autonomous vehicles (AVs) aim to help people reach their destination safely with minimal human interaction. In an AV, the traditional architecture has been moved from a human-centric system to an automated one. Based on the capabilities of different components, the AV can perform specific tasks. The environment perception and analysis module, which pertains to all the systems involved in the scene understanding, is crucial for autonomous driving. It includes all the information from a wide range of sensors and connectivity components on the vehicle. The perception module supplies the brain of self-driving cars with decision-making information. This involves converting data collected from the perception sensors into intelligent information. The acquired data can include sensed data from various sensors, such as cameras, LiDARs (Light Detection and Ranging), and radar as well as designated communication-based sensors, which take the data from Vehicle-to-Vehicle (V2V) and Vehicle-to-Infrastructure (V2I).

## 1.2. Importance of Environment Perception in Autonomous Vehicles

[6] Autonomous ground vehicle navigation, as the focus of robotics research and real-world applications, has attracted considerable attention and evolved rapidly in recent years. In different scenarios and environments involving road segments, unstructured outdoor areas, and highly dynamic environments, researchers have developed many machine learning-based navigation strategies. Although lower-level sensing and control challenges emerge in different terrains, researchers have generally focused on learning-based algorithms for perception, state estimation, and decision-making problems for vehicle navigation. Automated and semi-automated solutions for robot navigation as well as generic robotics navigation have also been addressed. Human mobility has been enhanced by combining robotics advancements, including advanced tracking and sensor developments, as with these technologies outdoor robots can follow humans and/or carry out useful tasks for humans. In an urban environment, robots are useful for multi-tasking and multi-human support, including elderly and disabled individuals' needs, construction, management, and repair tasks. The result is several surveys of different mobile robot tasks.[7] A traffic scene contains a wealth of dynamic information for vehicles to perceive in the environment. Learning agents for environment perception in traffic-scenes should be aware of the physical properties of vehicles and unavoidable uncertainties in sensing-based perception and decision-making. Combining sensor data with probabilistic and dynamic behavior models of the observed agents can generally reduce the effect of such noise effects. Our frame of choice has been a probabilistic scenario-based traffic scene representation, which includes unpredictable disturbances (steering-amplitude noise) and prediction noise. Ours' learning agents are aware of and act on the probabilistic representations of these artifacts. We have been utilizing this representation and the perceptive and predictive abilities in the acquired policies to analyze and improve the decision-making strategy of the learning agent acting in the traffic.

## 1.3. Role of Machine Learning in Environment Perception

The integrated fusion of visual fields and radar signal data is realized through the common input layer of the neural network in the study dealing with scene biometric recognition of different vehicles, backgrounds, and obstacles. Thus, the combined data is effectively used to form an integrated awareness environment model, which improves the pertinence, convergence, and reliability metrics. A new Cropping Block (CB) and a new basic structure are introduced, which shows that the fusion of the two-modal ice and Lidar data for training

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

capability improves the previous fusion methods in the study presenting three methods of autonomous environment perception based on deep learning [8].

Machine learning used in autonomous vehicle systems is inevitable, and new practical instances are seen now and then. Considering the environmental conditions, the system is supposed to perceive and display an enabling effect for the vehicle to prevent any possible threaten that might be caused by expected and unexpected conditions. Machine learning algorithms in general are involved in different stages of this perception such as data collection, feature extraction, and intelligent decision making [9]. There is a process from collecting visual-and LiDAR-based environmental information throughout the day on the test track to autonomous learning based on convnet learning and joint feature extraction to multi-modal data processing, which can recognize the 360-degree surrounding element of the vehicle including pedestrians, vehicles, and geometrical structures. Section 1.3 is focused on studies emphasizing the application of machine learning in environment perception [10].

## 2. Sensors and Data Collection

The variety of sensors used to collect information from the environment in which autonomous vehicles operate suggests that it impractical to write a unique perception system at the level of the information collected, for each type of sensor [11]. Each type of sensor provides information in a different way and involves different strengths and constraints that are implicit. For instance, camera data give colorized visual representations of the environment, but is limited by poor visibility conditions. LiDAR provides an accurate 3D map of the environment, but it remains sensitive to the weather. RADAR implies a slower frame rate but operates correctly under various harsh conditions; it offers advanced imagery of other moving vehicles. by crossing data across different sensors, therefore, one can benefit from the benefits of each sensor. A salient research choice, addressed by a wide range of researchers, and also addressed in the literature, is to use a multi-sensors architecture and to validate this alternative through the use of an evaluation protocol serving as a reference. Ultimately, the choice on the essential sensor and data processing combination to improve autonomous vehicle perception remains an open issue in the field of machine learning and autonomous vehicles in particular.

We can record and sense the environment in which the autonomous vehicle operates using a wide range of sensory devices, such as cameras, LiDAR, RADAR and ultrasonic sensors [12].

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

Thanks to machine learning, it also became possible to integrate this sensed information with contextual information for a better understanding of the environment and for a better anticipation of possible future events, such as the detection of moving objects, such as pedestrians, cyclists and other vehicles, intentions of vehicles and the history of their movements, the detection of lane markings, traffic signs, road geometry and various other attributes of the environment that are beneficial for safe and efficient driving. The sensors used not only provide detailed information about the environment in which the autonomous vehicle operates, but also about the vehicle itself [13]. Barzel et al. Proposed an architecture for interpreting the dynamics of autonomous vehicle events based on one-state-hidden Markov models, supported by longterm memories. This system is inspired by human reasoning mechanisms with the main ability to predict and understand the future of other drivers. The static and repeating routes of the vehicles are learned and used by the system as a reference point to identify the abnormal driving behaviour of the other vehicles. The system further uses a prediction mechanism based on Gaussian processes, giving rise to the ability to perform short-term anticipations.

## 2.1. Types of Sensors Used in Autonomous Vehicles

The advantage of ultrasonic sensors is the fact that cars have been equipped with these sensors for a long time period, which results in our society being used to them. Unfortunately, they can be used to present information only in the form of sound and sometimes also as a visual source of information. They are also characterized by low resolution and low output data rate, which makes them perfect for detecting obstacles during parking but insufficient for driving on a road [11]. A similar situation occurs with mono-camera based systems. Objects can be recognized based on the intensity of the light reflected from the objects' surface. A comparative analysis of the level of grey in an image of the road situated in front of a vehicle and the level of red in a similar image will allow for the detection of positioned vehicles. The issue is more complicated — the overlapping of diffusive surface ALA a jet engine of a cargo plane and the view suitable for detection, such as windows on a suburban bus, can differ. This house algorithm of object recognition is not capable of solving this problem. An algorithm based on 2D segmentation of objects will enable the visualization of supplemented objects. Even systems in which two cameras are located near infrared also provide 3D features [14]. Due to this, they are perfect for detecting objects. The limitation of the method is the loss of information about the number of detected surfaces, and thus the object geometrical mass.

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

Systems based on mono-camera have as well low resolution and low output data rate, and this result in poor detection of the object. Therefore, it may be relatively good for car detection but does not work for object detection in practice. Only in objects detection can compensate for the limitations of the remaining methods.

Environment perception systems of autonomous vehicles have to analyze the road scenes, avoid possible collisions, localize the vehicle, and plan the most relevant trajectory in real-time. To address these issues, every advanced driver assistance system (ADAS) and autonomous driving platform must be equipped with sensors, which are capable of active or passive perception of the chaotic and complex surrounding [15]. In the currently available ADAS and semi-autonomous vehicles, ultrasonic sensors are mainly used. However, they are not sufficient for the operation of fully autonomous vehicles. They should be equipped with other sensors, such as LiDAR, RADAR or digital cameras, which are capable of active measurement of the surrounding or a very detailed passive sensing of the surrounding.

## 2.2. Data Collection and Preprocessing Techniques

The leader in 3D LIDAR and third behind GPS-INS-Inertia sensors in the global company pattern is the preferred sensor. The data-set used on the autonomous vehicle often uses it as basic environmental data. The UAV123ND dataset is constructed in co-operation with Beijing and Pyongyang under the same working pattern. We installed an OUSTER OS1-64 LIDAR sensor on top of the car's roof and performed multiple korekm-laps around the track for the data set. During the experimental process, two vehicles with permission of law were tested, and the test got the results through them. The dataset was logically divided into three categories based on applications such as scene segmentation, scene classification, and point cloud classification [16].

Various environmental data are required for the developmental proceedings of autonomous vehicle technologies. Vision data, in some cases referred to as RGB-D camera data, is one of the most varied and commonly used vehicle sensor data in recent years. With the advent of deep learning and specifically CNNs, the development of the perception technology that sometimes process vision data to process the environment through sensor fusion has largely flourished [17]. However, there have been a number of obstacles in using vision data from moving vehicles during the driving phase, and therefore only a few data sets could be created. The data collection of the Vision Transformer (ViT) project was completed by 200-hour

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

driving by providing society with the first 3D Depth Map Autonomous Vehicle (DamAV) RGB-D dataset.

## 3. Machine Learning Algorithms for Environment Perception

In the machine learning community, several methods to perform environmental perception in autonomous vehicles have been studied in the literature. A widely used method to perform environmental perception is to use well-established computer vision algorithms that rely on the extraction and classification of local features or segmentation [18]. For instance, the detection of obstacles on the road could be achieved by extracting relevant features like parallel lines or the vanishing point in an image [13,14]. A similar proposal has been extended to perform lane detection by using line segments. The simple design and real time performance but significant robustness makes classic computer vision methods good candidates in the autonomous driving industy. On the other hand, an one-size-fits-all computer vision model lacks the adaptability for different scenarios so as to discriminate the context better. With the fast development of deep learning these years, CNN, the most famous deep learning model, has recently been proposed to address this issue, which adapts to the environments and nails down the perception results by providing higher feature extraction quality.

The introduction of machine learning algorithms has revolutionized the field of autonomous driving. These algorithms have been used for several vital components of the systems, including environment perception, trajectory prediction, and path planning [6]. Among all these components, probably the most critical one is environment perception, where the autonomous vehicle obtains an understanding of its environment by collecting, filtering, and processing raw data to estimate the attributes of key road components such as obstacles, traffic signs, lane markings, and roads themselves. In particular, by perceiving the environment, the autonomous vehicle could drive in a self-supervised manner, knowing where and how to drive to reach the desired destination while not causing damage to both itself and nearby objects [19]. Most environmental perception systems consist of the following general steps: 1)sensing, 2)perception, and 3)prediction. In the first step, the perception of the environment requires the fusion of inputs from multiple sensors, such as camera, light detection and ranging (LiDAR), and radar. Next, the binary (occupied/free) representation of the perception results will be used to generate the contour of the objects and the estimation of their poses.

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

The output of the previous step will provide the basis for the last step where the prediction of the trajectories ofboth dynamic and static obstacles is made. Machine learning has shown its potential on such perception and prediction tasks, thanks to its capacity to learn latent knowledge from large scale data [ref: c6b5335f-f0f8-48e1-bb76-1eb3696d1fa2

### 3.1. Supervised Learning Algorithms

The recent advances in the supervised learning method, like Convolu- tional Neural Network (CNN) in Autonomous Vehicles (AVs) have improved the object detection, classi cation, and instance segmentation tasks from the point of view of environmental perception for the AV systems [20]. CNN is sophisticated classi cation and representation method while the fully connected layers of CNN makes it highly susceptible to the over tting. In addition, the CNN model is sensitive to the different types of noise, light, and environmental disturbance. The situation gets more critical with the non-policy heavy softmax classi ers in the supervised learning. Pixels from the same object category may have quite different features and the classi cation results will depend on the external conditions like the environment, shadow, light, and noise. There are also efficient researches on domain adaptation. The main goal of domain adaptation is to transform a pre-trained model to a different distribution directly; the computer vision community deals with the Longest Zero-shot Learning (ZSL) task where a model used to learn from di- rectly relevant data [21]. In recent years, with the development of deep learning, self-supervised learning has been addressed with unsupervised learning algo- rithms that require less supervision. Speci cally, learning the relationship of the symbols with diagnoses, or simulating. Context, where there isn't supervision, is usually available. To improve the generalization capacity, a self-supervised learning approach has been applied to LiDAR point cloud data for motion segmentation in autonomousdriving scenarios. Even though some un- supervised algorithms closely learned the distribution of theenvironment, limited work in literature focused on creating self- supervised datasets for the environment at di erent wavelengths in the deep reinforcement net formalism [22].

### 3.2. Unsupervised Learning Algorithms

In recent years, the rise of self-taught feature learning and other multi-layer dimensionality reduction algorithms such as t-distributed Stochastic Neighborhood Embedding (t-SNE) have attracted increasing attention. The paper Tao et al. [23] provide a good survey on the utilization of various deep learning hybrid models with recommendation on some popular

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

deep learning architectures such as WaterNet, temporal convolutional neural network, ResNet and etc. The Autoencoder-based dimensionality reduction algorithm is considered as unsupervised learning, in which the training data consists of input samples x = {x(1),x(2),...,x(T)} in which x(t) is a 1 × L vector. They involve two major components: encoding from input layer to hidden layer, i.e., z = g(x; θ) and decoding from hidden layer to output layer or obtained samples, i.e., x = f (z ; θ). Unsupervised algorithms are often used to discover hidden patterns and low-level features in the input data.

Various unsupervised machine learning algorithms can be employed to perform exploratory analysis of unclassified structured data or to learn low-level features from unstructured data. With unsupervised learning, no specific rules are predefined by which the algorithm can be trained to recognize the data, rather the algorithm is left on its own to discover structures in the data [24]. The unsupervised algorithms fall into two main categories, clustering and association algorithms or frequent itemset mining. Clustering is a technique to partition a given set of patterns into clusters in the way that patterns in the same cluster are similar to one another in some respect, rather different from patterns in other clusters. In contrast, association mining generally attempts to discover strong rules for association, frequent itemsets, or closed frequent itemsets in a transaction dataset. In these terms, clustering may be viewed as a process for creating unsupervised labels while association mining may be viewed as a process for discovering relationships among the attributes [25].

### 3.3. Reinforcement Learning Algorithms

To tackle the cumulative cost-to-pgo problem, this article proposes an efficient approach that integrates neural network architectures into the discrete and continuous-control reinforcement learning algorithms [26]. Each algorithm is tested on the CARLA simulator, running three different test scenarios in an open urban environment along a trajectory, and consequently in MERGE scenario. The discrete exploration task exhibits obvious stochastic performance, mainly during the affordance recognition based approach. The two exploration frustrating issues which stochastically arises are: 1) the driving model becomes stuck due to high A/V ratio in the left or right lane and it semaphores inside an A/V vehicle platoon; 2) the driving model runs into the back of an A/V vehicle platoon. Conversely, the continuous infinitesimal knob exploraïtion can effectively mitigate the two issues.

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

To ensure the efficient and safe decision-making of autonomous vehicles, reinforcement learning can be employed [27]. Compared to classic navigation methods, reinforcement learning can better approximate the optimal policy to allow the vehicle to safely navigate through complications and challenging dynamic conditions. We used a deep Q-network (DQN) and prioritized experience replay to jointly learn lane-keeping and vehicle-following tasks. The objective of our work is to minimize the pink region, which is defined as the spatial difference between the max dynamically feasible distance from the lane centerline and the current vehicle position [28]. By considering this definition, the vehicle may always behave like a human before hitting the lane boundary, which is considered a mix of keeping tasks. The DQN-based controller was employed for longitudinal vehicle control tasks and all the observations from visual perceptions were provided as inputs. However, direct adoption of DQN-based algorithms may not be applicable for large-scale autonomous vehicle navigation, particularly when complex topology is present.

## 4. Deep Learning Techniques for Environment Perception

In summary, object detection consists of associating distinct 3D and 2D points in the scene to an object, classifying the object, and estimating its pose. There are a variety of modalities used to perceive these three characteristics: lidar, sonar, ultrasound, radar, vision, and even no modalities (using raw acceleration or position data). Object tracking preserves these associations but repeatedly receives and updates them as new sensor data arrives. Sensor fusion contextualizes this new data to maintain a consistent semantics across both time and space, including data from multiple sensors in the system and providing different levels of reasoning in terms of its position, orientation, and motion. Additionally, this model is the main contributor to the overall perception robustness in both spatial and temporal dimensions [13].

The need for intelligent, context-aware systems has led to an increase in interest in sensor data fusion with deep learning models. Vehicles outfitted with a myriad of sensors can interact with their environments and gather information to assist drivers in making decisions or to control the vehicle themselves [29]. There are three main deep learning tasks for environmental perception: object detection, object tracking, and sensor fusion. These tasks differ from most pure perception tasks in that they require algorithms with the speed and robustness necessary for real-time implementation; it is not sufficient to simply detect an

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

object with high accuracy because the state of the vehicle is constantly changing, moving further from or closer to the object. Moreover, the object itself could be moving, also requiring the perception algorithm to track it. The final significant difference in sensor fusion is the necessity of having robust models that can integrate various types of sensory and contextual information, aggregating disparate pieces of information to maintain scene understanding in dynamic scenarios . This task is also relevant to enhancing the spatial and temporal robustness of the different levels of perception by contextualizing the object detected at multiple points by analyzing the sequence of sensor data and Vision.

**4.1. Convolutional Neural Networks (CNNs)**

In all of the abovementioned models, a CNN is used for feature learning, followed by one or more fully connected layers to produce the final output. Each of these networks uses different strategies for combining information from convolutional and fully connected layers, with the final objective, being to generate a fixed number of high confidence detection bounding boxes, corresponding to the objects present in the input image. For autonomous navigation purposes, semantic segmentation of the entire frame to identify drivable and non-drivable regions can aid in obstacle avoidance and planning. The main challenge in developing high accuracy and reliable CNN architectures is the large number of parameters. Thus, various CNN architectures have been proposed that modify the basic architecture to utilize fewer parameters, while maintaining similar accuracies [30].

Convolutional Neural Networks (CNNs, or ConvNets) are biologically inspired models that have gained immense popularity in computer vision tasks [31]. CNNs have shown state-of-the-art performance in a variety of tasks such as image classification, object detection, semantic segmentation, and instance segmentation. Object detection, which is the task of localizing and classifying objects in a scene, is particularly important in driving applications. Popular CNNs for detection purposes include Faster R-CNN, Single Shot Detector (SSD), RetinaNet, and You Only Look Once (YOLO) detector [22]. Typically, CNNs contain convolutional layers, pooling layers, and fully connected layers. The main characteristics of CNNs include parameter sharing and spatial hierarchy, which increase the effective receptive field and aid in learning patterns within images.

**4.2. Recurrent Neural Networks (RNNs)**

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

Exemplar studies show that RNN hybrids outperform so called handcrafted and tailored designs at vehicular data processing tasks like, event detection, football action recognition, traffic signalization and longitudinal maneuvers classification [32].Charsets0290615149721NoneNatural_Language Processing (NLP) [_] text processing tasks like, Language learning, Voice Recognition, Game Chats, Image to text conversions cleaner text processing has always lacked vivid establishment and it has been believed, Maximum Likelihood Estimation (MLE) hints. Letter, word, and software embeddings in NLP are the precedimg hybrid candidate w.r.t. hybrids. Deep learning already has shown its magnum opus heart warming performance on comprehensive range of NLP processing tasks. Text Data kmsare sequential and is hence taken under the jury of RNNs for long. RNNs in text processing primarily started as a solution in an abandon for bounded and unbounded LSTM. As far as the sequence complexity is concerned, the vanishing and the exploding gradient problems are under the same family of problems.机器学习驱动汽车相关任务是一个高耦合、长期的研究领域,依然有很多未知之处。RNNs show very promising improvement over traditional Machine Learning models due to their ability to understand temporal sequential information [Y. LeCun et al., 2015]. RNNs are built with the ability to maintain sequences of past information. Even though vanilla LSTM has its limitations such as RNNs are usually not better than conventional model of statistical approach on time series problem and Another limitation of vanilla LSTM is its restriction of memory storage. Different kind of modifications have also been proposed to deal with these limitations, such as modifications to RNNs which have shown promised improvements.

Standard neural networks are confined to a single layer with feedforward connections, while RNNs are depicted as chains which is pretty handy when trying to capture temporal information. Registered advancements in 3D Convolutional Networks has started incubating train set like sequential scan patterns because one could just use a 3D ConvNet and this alone. However this always comes to the cost of tiny resolution limitations, and power consumption, processing time stretch and most crucial of all the raw data integrity. Similar to attitudes in general image recognition areas, incorporating hybrid designs made idoneous senese was the obvious next step in vehicular data where the data nature is inherently sequential, which comes through the mercury of RNNs [33]. Wide normal LSTMs and Gated Recurrent Units (GRUs) preprocessed and processed significant attention towards vehicular data, nearly

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

equally, until Alexandrapoulos et al. (Alexandropoulos et al., 2016) brought Convolutional LSTMs in they proposed Run-length encoded Spatio-temporal Convolutional Networks (ReST), aiming better prediction accuracy and reduced memory requirements. 2DCNN, 3DCNN, LSTMs, GRUs and ReST are the most evaluated RNN variants in the recent literature w.r.t. vehicular data processing. ReSTs, for instance, have been reported to provide superior performance in the transportation literature when compared to former RNNs.

[34] Recurrent neural networks (RNNs) [Y. LeCun et al., 2015] are a group of artificial neural networks where connections between units form a directed cycle. This creates an internal state of the network which allows it to exhibit dynamic temporal behavior. Fundamental principle of RNNs lies in that each activation is a combination of the input and the last internal state of the network. Therefore assuming that the network is fully recurrent, in every moment of time a node's activation combines only information from the past. However there are a variety of architectures that work around this limitation such as Gated Recurrent Units (GRU) and Long Short-Term Memory (LSTM), this idea is of utmost importance when discussing RNNs because huge proportion of vehicular data is sequential in its nature.

## 4.3. Generative Adversarial Networks (GANs)

Gating Recurrent GAN (GRU-GAN) combines the sequence estimating capability of Gated Recurrent Unit (GRU) with the ability of GAN to generate real and unpredictable data. 3D-CGAN is a simulational technique which can generate synthetic samples. A naive simulation could be impractical in computer graphics simulations, but when the process is based on GANs, we can sure to get a sample which builds upon a distributed property, effectively from scratch. Relativistic average GAN models enable conditional high-resolution face image generation. These models can help determine a person's expression from a morphed image, and can be associated with task orientated social model scenarios. Kitani et al, study the problem of collective trajectory prediction by viewing human motion as a multi-agent system. GANs are used to backpropagate the critic's output back through the generator, which tries to improve upon the data and the critic cannot distinguish.

[35] [36] In the recent years, various works have been proposed in literature based on generative networks and in particular on Generative Adversarial Networks (GANs). GANs employ two networks: a generative network G and a discriminative network D. These two networks are trained simultaneously in a zero-sum game, aiming to fool one another.

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

Learning in this manner results in the generative network generating samples that resemble a real dataset. GANs are not only capable of duplicating a given dataset, but they can also mimic a distribution of data once equipped with a dataset. Thus, GANs learned to generate images with certain styles. The simplicity and effectiveness of GANs in image generation has led to the exponential increase of its applications. GANs have been widely used in video prediction and generation of data for several applications in the computer vision domain.

## 5. Sensor Fusion and Multi-Modal Perception

The authors in [37] have introduced a comparison of camera-LiDAR, fusion-LiDAR, and image-camera sensing techniques. It is widely believed that sweet spot for sensor fusion corresponds to mid-level fusion, which represents a flexible sweet spot between the expressiveness of high-level fusion and the computational/practical difficulty of low-level fusion. Therefore, we provide a detailed summary of this line of research that involves different degrees and levels of fusion for different tasks and modalities in autonomous driving research. Finally, we summarize popular strategies for lightweight multi-modal perception for faster inference times. In this section, we introduce different fusion strategies, and describe popular multi-modal tasks in autonomous driving.

Although individual sensors, such as cameras, LiDARs, and radars, can provide rich but incomplete scene information, multiple sensors can complement each other and result in robust scene perception. The fusion process in perception is called sensor (and/or modal) fusion, which can take place at different levels of the information flow. For instance, low-level fusion typically corresponds to feature-level image processing, where data from different sensors are combined resulting in the generation of multiple feature representations. After that, other multi-modal approaches refine the individual modality predictions with other modalities' features such as aligning the predictions of LiDAR detection with the learned features extracted from the camera. These activity estimates are usually used to perform a high-level fusion process, where objects from different modalities are determined by using the learned representations. Such modal predictions can then be directly utilised for driving control or more sophisticated applications such as decision making, routing, etc .

[38] [39]

### 5.1. Fusion of Lidar, Radar, and Camera Data

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

Therefore, integrating information from different sensors becomes a crucial and prior task in perception systems. It is necessary to set up a fusion system where the individual sensor data could be merged. It should also be learning capable since models for perception might be different for different sensors, and developing an unsupervised ways of reducing biases and computationally heavy registration are also crucial challenges. There are many works in the literature for sensor fusion in autonomous driving. These working focus on different types of sensors and algorithmic methodologies. In this study, we discuss the fusion of Lidar, radar, and camera sensors as they are the golden triangle in autonomous driving. We focus on the inherent difficulties in defining fusion methodologies of this trio of sensors since they are among the most widely used sensors in autonomous driving in the literature. In this section, we mainly focus on the fusion of these three primary sensors for autonomous driving. The rest of the objects detector models are beyond the border of this article. After getting the fusion data, it is commonly detected from a multi-level network by identifying the objects and their position.

The integration of data from different vehicle sensors can lead to a more comprehensive understanding of the vehicle's surroundings, providing safe, efficient decision-making for autonomous driving systems. This is because different types of sensors possess complementary characteristics on sensory modality. For instance, Lidar sensors can detect the depth information of objects with high precision. On the other hand, radar sensors are better at extending perception range, are true multi-modal sensors, work reliably in adverse weather conditions, and have the ability of detecting speed without being confused with the relative motion of other objects such as tire rotation and pedestrian movement. Cameras serve as the primary sensor to perceive lane markings, traffic signs, and traffic signals. Furthermore, humans and animals can be well dealt by cameras. It is necessary for autonomous driving to leverage these sensors together to exploit their diverse characteristics. Furthermore, it can improve robustness and reliability of the environment perception systems which is essential for the safety-critical application of autonomous driving [40].

## 5.2. Challenges and Solutions in Sensor Fusion

The only visual-object dataset from the weather complexity perspective is the Carla dataset, covering seven weather and lighting conditions. In contrast, the ApolloScape dataset is limited to conventional and synthesized-six other weather and lighting conditions. Therefore,

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

we are close to a lack of an extensive research project detailed dataset-weather complexity approach to understanding AV perception in different weather and lighting conditions [41]. In addition, the monocular cameras capture the information as a 2D image, while the LiDAR and radar perceive the environment as point cloud information. Furthermore, the technologies for sensor calibration and multimodal alignment are mature, and the camera-LiDAR sensor fusion remains a hot research topic. These technologies require advanced sensor fusion methods for the generation of trustworthy fusion results, such as learning-based and geometric-based methods.

When autonomous vehicles operate in complex urban and highway environments, they face numerous challenges that can only be addressed using multi-sensor fusion technology. Insufficient perception is among the main technical challenges in the development of autonomous vehicles. A standalone monocular camera cannot provide as accurate a perception as both a high-precision LiDAR and a radar in various complex road environment conditions [37]. The fusion of 2D and 3D point clouds from the LiDAR and 3D radar is meaningful for these tasks and can greatly improve the detection performance of autonomous vehicles. Specifically, the improved sensor fusion improves the detection range and the perception accuracy, creating multiple perspectives to understand the environment from more viewpoints and providing an effective visual complement to mitigate the problem of sparsity [42].

## 6. Evaluation Metrics for Environment Perception Systems

On the other hand, detection and tracking object are carried out by joint ML and DL methods in, where user can sent a video with his specific queried object that is detectable by a CNN model, to obtain the video in which all frames contains at least one queried object. In the field of robotics, distance estimation is another essential task, especially when the depth sensors are not available. The principal objective of this work is to evaluate the performances of the DL model for vehicle tracking, target detection, and distance estimation, considering road and external conditions that might influence visual perception. Among all the factors, road conditions are making accurate and reliable obstacle perception and localization rather challenging due to deformation, occlusion, and lighting sensibility on moving obstacles, especially on the image processing-based methods. An intelligent system for decision-making

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

not only needs to collect information through different sensors, but also fuses the information efficiently.

Artificial intelligence (AI) has gained increasing attention as a key-enabling method for intelligent vehicles. Most AI methods, for example, vision-based deep learning (DL) techniques, require potentially heavy computations, where real-time or near real-time processing is not guaranteed. While the AI has shown abilities to outperform traditional approaches, especially in the field of environment perception in which no model or rules are available, traditional methods could still be included to achieve near-real-time processing, by employing an ensemble mechanism. In this work, two novel evaluation metrics have been adopted to evaluate the environment perception systems achieved by machine learning (ML) and AI methods: deep learning (DL) sensor fusion algorithms for autonomous vehicle in, and the detection and localization problem. Recognition is performed frame by frame, and when the given frame does not contain a queried object, detector proposes the top few most likely locations in which the queried object may be located within the video, and run the tracker at each of these candidate locations.

## 6.1. Accuracy, Precision, and Recall

For collision-critical safety evaluation, standard precision scores and recall scores are insufficiently tailored to detection safety goals and do not incorporate the considerable differences in collision safety between different, often interacting, objects. For example, lots of objects hover or float closely above the ground. That means that a collision with, e.g., an otherwise harmless pedestrian can be less severe than a collision with the autonomous vehicle (AV) traveling on the road. Therefore, there is utility in having more sensitive, more safety-relevant metrics than standard ones that do not incorporate the different interaction dynamics. A next level of development of these measures would be to quantitatively compare results for both an AV performing the perception action and a human annotator who makes the same decisions. This is an important metric to see whether the performance of these perception systems is at a level where errors are likely to be noticeable by someone. The vast majority of high- and low-level perception error evaluator measures today ignore the evident fact that these errors are safety errors of some kind that should be perceived by the autonomous vehicle systems. A higher degree of flexibility would open up and all systems that have a specifically tailored function as a key output could be measured in a safety-specific

way. Much more work is needed and these safety errors should be taken seriously in further development of if these safety errors need to be taken seriously in further development of perception systems [43].

The accuracy, precision, and recall of an object detection system are three key measures that are often used by researchers and developers to evaluate how well an algorithm can recognize large datasets, like ImageNet or other public benchmarks. Although these basic measures are widely used, there is an increasing need to develop evaluation measures specifically tailored to the requirements for autonomous vehicles (AVs). These are environments with significantly more complex traffic interactions and where the end effects of decisions made based on the sensors' outputs are particularly critical. Within a safety-critical application such as a self-driving car, correct detections are more important than false ones, and risk-related information, such as collision safety, is arguably more informative [44].

## 6.2. F1 Score and Confusion Matrix

The survey also identifies future research trends and gaps that need to be addressed to shape the future of AVs. While the quantity of work reviewed is satisfactory, the review of the literature [45] on automated environmental perception and analysis provides a unique overview of the contribution of machine learning in this technology. Section 6.2 trains a CNN optimized with the ADAMW optimizer on NVIDIA triple 30 Ti GPUs and configures it to measure the performance of effective models such as precision, recall, accuracy and F1 score. The F1 score provides a balanced distance for the recall and precision, which allows for a simultaneous consideration of both metrics to measure the effectiveness of binary classification models, making it a more robust model for evaluating highly imbalanced datasets. The present research advocates that the proposed assessment could be conventionally used in future research as a standard outcome for the benchmark recognition of autonomic models [46].

Automated learning for autonomous vehicles is at the forefront of technology-based solutions and enables autonomous vehicles to mimic and enhance human experiences [43]. In previous surveys, this has been covered from diverse dimensions including environmental modeling, navigation, and route planning, human-like driving behavior modeling, social autonomy, and traffic-aware motion planning. This survey builds on the work aforementioned and provides an in-depth exploration of the survey of Machine Learning techniques and methods used for

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

the perception and analysis of the environment. In light of recent efforts to understand environmental perception, mapping, localization and high-level route planning modules take advantage of environmental modeling data which is obtained usually through Telematics data, SLAM-based mapping, and computer vision.

## 7. Applications of Environment Perception in Autonomous Vehicles

- In AVs, every system part functions in cohesion, and it is envisaged that any part of these systems will have to interoperate with neighboring units to provide a benefit to the host-system in the future. Environmental perception methods will also be required to open the door between their machinery layers to found viable operation strategies. The fusion of multiple diverse sensors will allow the AVs to understand their surroundings more effectively. Data from the involved sensing sources with different natures will be employed together to make better decisions about the future scene dynamics. In another balcony of AI, work has focused on decision-making processes within one ecosystem. On this agenda, it is essential that EF on the scene and SF from the sensors of the vehicle can behave in synergy with keeping a constant exchange of decision-making processes between them. A diagrammatic view of that point is devised in this compartment, and also a few of the extractable models of decision-making process mechanisms will be introduced afterwards. Configuration range of sensor technology and design of available fields of view of involved drivers should be the assessment main estimations of the sensors. Besides, using AI-enabled SF methods, AVs can collaboratively communicate with each other for road performance inference and mutual warning for potential off-nominal conditions.

- Safety has always been one of the main objectives in AVs, but cannot be achieved without secure and accurate environment perception capabilities [38]. All decisions about motion and interaction are established on the acquired sensory data, and different numbers of available and possibly redundant sets of sensors could be mounted for AVs created for different operational spheres. However, vehicles can interact with advanced approaches while utilizing the given hardware set in different sceneries such as open off-road sections, semiurban areas and dense urban streets. A useful environmental perception suite should include limpid software and energetic sensor fusion methods; these techniques should be able to receive raw data from available sensors and should convey important information to higher-tier systems after multi-level categorization – raw to object and scene level. In this review, each regard of

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

object-oriented digital scene representation will be highlighted. Furthermore; present-day environmental perception approaches that most of the automotive industry are relying upon for AV prototyping will be evaluated [2], and following an extended NET will be addressed [47].

### 7.1. Object Detection and Tracking

This article aims to provide a systematic review of the object detection and tracking algorithms for different types of static environments by deep learning with both labelled and synthetic datasets using camera sensors [10]. This can enhance the holistic understanding of the sensor-based environment perception for different levels of autonomy in the use case of an autonomous vehicle. We also review traffic datasets for perception, especially the tracking datasets for the development and testing of perception systems. The article also provides a systematic review of the techniques to fuse camera, LiDAR and GPS based input data to perform interactive, multi-scaled object detection and tracking using a fusion of classical and deep learning models. Finally, the article provides a list of open issues and research challenges that need to be addressed to develop holistic environmental perception systems for future autonomous platforms.

[22] [48] The task of environment perception for autonomous vehicles is critical and challenging, and it consists of object detection, instance segmentation, and multiple-object tracking. Real-time object detection and tracking are crucial for ensuring safe driving of autonomous vehicles by ensuring accurate object identification, state estimation, and predicting future behaviors. Environment objects can be represented by both static and dynamic objects (e.g., pedestrians, vehicles, and bicycles). For example, two essential tasks in environment perception are to recognize moving obstacles and predict the future trajectories for perception and motion planning to avoid collisions. Therefore, object detection and tracking are two fundamental components in the entire framework of autonomous vehicle systems.

### 7.2. Lane Detection and Path Planning

Based on lane detection, given the real-time lane information and juncture turning directions, path planning makes the vehicle reach the desired terminal location. There are two different categories for lane detection in traditional methods [49]. (1) Parametricized lane detection

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

models: These models use six defining parameters for the left/right lanes, and the equation should contain distance, width, the lane function direction, and the lane continuity and curvature variation. - Continuous lane model parameters: The usage of higher-order polynomials and the direct representation as parametric space. - Elongated lane models: An infinite-lane-long- region parametric space to handle long-lane model boundaries compared to the image region. -Piecewise models: v is augmented with the lane lines for better detection handling very curved situations tension side. (2) Structured detection approaches: This category contains the following methods. - Semantic Segmentation Based: The method that detects the lanes based on the semantic segmentation maps by appropriately post-processing the networks' output. - CropleNet: USA visionary method cards an unusually significant to assist various autonomous concerns, for example, lane detection, ranging, technique cubes the image to aggregate the least containing the lane line. - Patch-based one-stage method: This replaces the staged appearance of Robust Lane Recognition and Tracking with deep learning for the select of Fast R-CNN.

The next section explains the lane detection and path planning techniques. Lane detection is a fundamental component and is essential in path planning and trajectory selection, and modeling of lane detection principally depends on the input sensor [50]. Grayscale, color, depth, range, infrared, and night vision sensors have all been popular and utilized. Traditionally, road segmentation techniques divide the process into the following major parts: pre-processing, feature extraction, and post-processing [51]. The pre-processing step removes noise and irrelevant input from the sensor. The main task of the feature extraction stage is to highlight the regions of interest and eliminate the undesired information. The final post-processing step is to fill in the gaps, smooth the obtained road mask, or to remove irrelevant segments. Lane detection and classification in the bird's-eye view is performed in stereo calibrated images. The approach involves the following major steps: (1) image acquisition; (2) rectification of stereo images; (3) feature extraction and stereo matching; and (4) post-processing of the stereo data. A strategy for lane recognition and trajectory selection strategy using a dense 3D map of the environment has also been proposed using stereo vision.

## 8. Challenges and Future Directions

Security: Ensuring that systems are free from unwanted interference is another major issue. An unauthorized person cannot gain access to the car, either from its digital infrastructure or

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

physically. Disaster Situations: Systems require further development for critical and disaster situations. In case of an unavoidable accident, the system should abandon the protocol and try to minimize the accident's effects [ [13]].

Real-Time Requirements: In real-time adaptive systems, the requirements of long-term prediction, optimization, consecutive actions, and hardware requirements should be combined into one purpose. Safety: Various sensors have different localization error margins. This requires the application of sensor fusion methods for these sensors.

Environment: It is difficult to determine the arrangement of lanes and other vehicles. These dynamic environments require accurate and fast sensors and information fusion methods.

Legal Framework: Nearly all experiments in the area of autonomous cars are currently being carried out under experimental regulations; therefore, formulating the regulatory framework for autonomous vehicles, as well as insurance, remains the major challenges.

Future and real-time requirements demand more intelligent systems []. The future problems and challenges in autonomous vehicles include the following. Design and Development: Designing cars that can operate in a legal and safe manner on roads remains a c hallenge.

## 8.1. Robustness and Generalization in Perception Systems

As we mentioned in the introductory section, many challenges and opportunities drive neural network research for autonomous vehicles. We concentrate on problems related to robustness and generalization. The tools coming from the neural network framework in point cloud processing are opening new ways and stimulating research connected with robustness and generalization issues. In this review, we provided a list of the biggest problems present to give the reader an idea of the technology maturity and help understand what kind of solution is needed in the field. That provided a context for the listed future challenges connected with robust and generalized solutions in such modules as 3D object detection, tracking and fusion, room semantic extraction and dynamic scene graph tracking and the review. Considering fast development of data collection models, we foresee that we need clever solutions talking into account considerable robotic deployment, which in fact is not represented in the reviewed literature. Instruments to significantly reduce still prevailing carbon footprint are also desired. And the last problem is powerful and informative evaluation hypotheses development of

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

generalization capability of the approaches developed with robots in the loop for all the tasks present in the review [52].

Perception is the foundation of autonomous vehicle systems, and its associated modules critically influence the performance of other driving systems. With increasingly added sensors, significant amounts of high-dimensional data are collected, which increases the workload for the processing modules. These systems raise several concerns such as, what happens when individual units of the modules are faulty or work with deficiencies due to sensor malfunctions (e.g., failure)? What happens if a certain sensor is not available? How will systems behave in new and/or different conditions from the training phase [19]? Researchers have different opinions on the "ifs and buts" of the situation that will occur more than expected in real life. However, what is voiced by everyone is the necessity of a comprehensive perception system that can robustly and generably perceive the environment we are in [53].

### 8.2. Ethical and Legal Implications of Autonomous Vehicles

Although some of these questions are ethical in nature, they also have legal implications, and a regulatory framework covering these issues needs to developed urgently. The 2019 European Parliament's resolution on artificial intelligence and the 2018 German ethics commission's report for autonomous driving urge that ethical questions not be left to market players alone, but that Europe should establish a comprehensive and binding legal framework [54]. Public health and safety, liability and broader questions of regulatory frameworks are on top of the lists of ethical and legal questions. Moreover, an autonomous vehicle's operating system raises important privacy questions: The system consists of a multitude of interconnected systems communicating not only with each other but also transmitting and saving data. The potential for large scale influence, damage or theft through manipulation of such autonomous vehicles appears to be huge.

The introduction of machine learning and artificial intelligence to autonomous vehicle perception systems and the road raises a multitude of ethical and legal questions. Machine learning algorithms can often not simply be designed to maximize utility or minimize harm because their operationalization in decision-making systems is fraught with ethical concerns [55]. This will prove a challenge for the development of machine learning algorithms for autonomous vehicles, which can analyze the surrounding environment, and require to

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

continuously perceiv changes thereof, and predict future changes as well in order to react in time [19].

## 9. Conclusion

To conclude, in this paper, we have reviewed articles covering traditional sensor based perception and processing techniques, modern LIDAR, RADAR, Camera based novel algorithms for a self-driving car and modeling it with 3D animations and testing it for accuracy for usage in various environments (outdoor of technical institute and fuel refueling station for underground environment and take off and landing for aerial robot). We have also studied the navigation planning for self-driving cars by mapping the vision based input 3D model, and optimizing it by using GA on a mobile testing platform (aerobot) in software PIONEER, and 3D planning for GPS constrained aerial robots. Although, there are other sensors available which were not included in the process, like sound sensors, imu, ultrasonic sensors, In-Frayed sensors like found in Boston's Spot robots of the 2000s but through this research in the visual field, it is concluded that it is possible to have a self-driving car/an autonomous robot with the vision based simple or advanced techniques. However we still need to keep improving the design, like by finding other real time processing techniques (FPGA, Ghost, Neuromorphic chips) rather than in the computers for LIDARs and monocular or stereo cameras to reduce the time. Mobile Robots and technology are in 4IR era and there are still many opportunities in this field.

End-to-end and modular planning systems for autonomous vehicles have been majorly fuelled by deep learning technology. They have been widely studied and reported for micro to macro environmental perception [42]. In particular, the state-of-the-art deep learning models for scene segmentation and object detection in autonomous vehicles are PixelLink, Mask R-CNN, and Advanced Driving Assistance. Some typical deep learning-based monocular depth estimation models are CSS, Features Transform Models, R-MVSNet, Wang 3D, and Pyramid Depth Estimation. Visual SLAM which directly generates 3D information from 2D video streams is another widely used tool for autonomy, mainly for prediction and semantic mapping. Among the fronts of environment perception and analysis for autonomous vehicles, deep learning has been persistently popular. From all the discussed topics in this survey paper on the deep learning-based environment perception and analysis for autonomous vehicles, it is evident that scene segmentation, object detection, object tracking,

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

monocular visual SLAM, and end-to-end behavior prediction have become very popular, but they have been mostly studied independently. Cross-technique studies looking into the correlations between them would be helpful for this community. Further, recognizing other types of road users such as two-wheelers and pedestrians/further studying the road scene harmonious with the availability of traffic lights, zebra crossings, and other sign boards will improve the scope of autonomous driving operation [56]. Additionally, there are challenging tasks at an individual and international level for the research community, like planning in extreme scenarios, safe and trustable perception, and autonomous vehicle testing and benchmarking protocols. In conclusions, this review paper on the deep learning-based environment perception and analysis in autonomous vehicles has been presented as a baseline for future reference.

**Reference:**

1. Tatineni, S., and A. Katari. "Advanced AI-Driven Techniques for Integrating DevOps and MLOps: Enhancing Continuous Integration, Deployment, and Monitoring in Machine Learning Projects". *Journal of Science & Technology*, vol. 2, no. 2, July 2021, pp. 68-98, https://thesciencebrigade.com/jst/article/view/243.

2. Shahane, Vishal. "Optimizing Cloud Resource Allocation: A Comparative Analysis of AI-Driven Techniques." *Advances in Deep Learning Techniques* 3.2 (2023): 23-49.

3. Abouelyazid, Mahmoud. "Comparative Evaluation of SORT, DeepSORT, and ByteTrack for Multiple Object Tracking in Highway Videos." International Journal of Sustainable Infrastructure for Cities and Societies 8.11 (2023): 42-52.

4. Prabhod, Kummaragunta Joel. "Advanced Techniques in Reinforcement Learning and Deep Learning for Autonomous Vehicle Navigation: Integrating Large Language Models for Real-Time Decision Making." *Journal of AI-Assisted Scientific Discovery* 3.1 (2023): 1-20.

5. Tatineni, Sumanth, and Sandeep Chinamanagonda. "Leveraging Artificial Intelligence for Predictive Analytics in DevOps: Enhancing Continuous Integration and Continuous Deployment Pipelines for Optimal Performance". Journal of Artificial Intelligence Research and Applications, vol. 1, no. 1, Feb. 2021, pp. 103-38, https://aimlstudies.co.uk/index.php/jaira/article/view/104.

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.

**African Journal of Artificial Intelligence and Sustainable Development**
**Volume 3 Issue 2**
**Semi Annual Edition | July - Dec, 2023**
This work is licensed under CC BY-NC-SA 4.0.