# Explainable AI for Transparent Decision-Making in Autonomous Vehicle Systems

*By Dr. Yuliya Shylenok*

*Associate Professor of Applied Mathematics and Informatics, Belarusian State University of Informatics and Radioelectronics (BSUIR)*

## 1. Introduction

One of the ultimate goals of AVs is to achieve high levels of traffic safety. According to a report of the WHO, the Halving Global Road Traffic Deaths and Injuries report (2015), worldwide every year about 1.35 million deaths and 20 to 50 million injuries occur in traffic. AI systems in AVs are not only supposed to reduce most of such accidents caused by human error but also generally to increase the reliability of traffic substantially. AVs are supposed to collect a huge number of test kilometers for proving their absolute reliability statistically before they are allowed to participate in traffic. If the majority of (test) vehicles are rather 'green horns' on public roads at first, and a small fraction of 'super professionals' exist at best, accidents will happen all too frequently because test kilometers in traffic are countable in principle only. However, examining every single faulty decision arising from a test kilometer or infinitesimal drive in this test hall is prohibitively expensive, while AI-driven ASs especially enjoy producing statistical learning models. On top of this, due to strongly integrated hardware, software, and safety gear, faulty AV drives are indeed worthy of detailed consideration; one of the lasting issues is sensor and actuator reliability. Safety is the overriding principle in AV design – AVs are classified as ASCAS (SAE Level 5), meaning that ASs are the sole drivers and nannies. Hence, the decisions of the AS at any point in time have to be reasonable, and subsequently, the consequences have to be predictable. There is an increasing demand from the general public, legislators, and trade associations for transparency, reliability, and safety in all AI at the current time, but especially for ASs with their varied applications such as autonomous driving.

Safety-critical Autonomous Systems (AS) such as Autonomous Vehicles (AVs) require transparent decision-making since their core goal is to transport not only the vehicle but also

**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.

passengers and even cargo safely from an origin to a destination without major human intervention or assistance [1]. On the other hand, in machine learning and AI, there is a trade-off between predictive power and transparency [2]. While black-box algorithms in ML offer impressive performance, their decision-making process is essentially uninterpretable, i.e., not amenable to human understanding [3]. In this paper, we discuss the importance of explain-ability and interpretability in AV design and optimization. Public acceptance strongly depends on how accurate, reliable, and transparent the activities of the AS are. It is essential for public trust in autonomous driving that the decisions of AVs are intelligible and transparent in real-time. In this context, transparency means the understanding of the learnable AS decision-making process and intelligibility means the explainability of decisions in a user and/or other system-interpretable form. AVs have to act within the boundaries of the regulations. Thus, abiding by the sensor, processor, and actuator chain, the ASs must grant, just like a human driver would in a similar situation, that their current and future courses of action are permissible within the legal framework.

…

## 1.1. Background and Significance

As a consequence, AVs rely on decision-making algorithms that should be carefully designed and tested, and their output explicitly needs to be not only validated, calibrated, and certifiable, but also traceable and comprehensible to human operators and competent authorities [1]. Assistability can contribute to the development of interpretable models in a systematic, insightful, and fast way by utilizing symbolic reasoning, exploiting the available knowledge in the form of explicability rules, and allowing the users to control conformity to expectations for the learned models. Expert systems built with a knowledge base and an inferencer are principled but tend to either lack explanation capacity or generalized knowledge. Combining expert systems with machine learners—inductive expert systems—can integrate empirically acquired knowledge during learning. Such systems can learn from both symbolic explanations and feature relevance (e.g., SHAP) to further enhance the expert system.

Autonomous vehicles (AVs) have the potential to revolutionize transportation in terms of safety, efficiency, and accessibility. As part of the automotive industry's movement toward AV development, many companies use and rely on AI models for perception and control [4].

**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.

However, since these models are black boxes, it is difficult to understand their underlying decisions. This is especially problematic in real-world scenarios, as it is difficult to create test data that covers scenarios as diverse as real-world traffic. For further reasons, even deep learning algorithms are interpretable only if the learned weights are referring to "transparent" inputs describing the input space. Furthermore, some aspects of the decision space—even if human interpretable—need not be of interest to humans such as telling whether each sensor measurement is used to determine lane-keeping (LK). LK is typically the hardwired response to lateral deviations and modeling this input action mapping in the depth of the involved artificial neural networks costs valuable resources. To be able to focus on what aspects AI decisions are based on is even of more interest when it comes to adversarial attacks on AI-based perception systems [5]. Also for safety assessment plus functionality and behavior analysis of AI, a concise insight into relevant input/output connections and fitted artifacts in terms of weights, locations, feature priors, data candidates, data-trained weights, and values is crucial.

### 1.2. Research Objectives

The line graphs in Figure 2 show the metrics' internal scores $z_1 \dots z_m$ Charcoal in the studies covered in Table 1. Each metric captures part of interpretability, while the m.charcoal number of useful features among these needs to be determined explicitly. Sometimes, the properties acquired under $z_1 \dots z_m$ help us determine that some AI models are globally explainable to at least a certain extent (Section 3). A common issue identified from these studies is that the trustworthiness of the whole autonomous vehicle (AV) is often tested separately from that of an individual model. Therefore, we are unable to analyse how the interpretability/transparent-ness of the overall system is influenced by the component error inferences. This often leads to the scenar-IO fio that ignore-error statistics could be the culprit behind most of the seemingly contradiction results in the reliability, causality and trust literature.

[4] [6] To address the credit assignment problem, XAI for AV needs to focus on how and why an AI model may commit an error. This is due to the fact that in many high-level decision-making scenarios (e.g., lane changing), neither the training data nor the univariate input features by themselves provide enough information to explain the performance of the AI (Table 2). In addition to such knowledge-based explanations, dynamic mechanisms may also

**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.

help. Since dynamic real-time data can be readily collected, their explanation and optimization might lead to a transient solution of Explainable Robotics.

## 2. Fundamentals of AI in Autonomous Vehicles

When developing complex autonomous vehicle systems, especially those applied to society, the requirement for explainability is often raised by stakeholders. Humans are explanatory creatures, and it can be unacceptable to trust that the system, albeit complex, makes the right decision without providing explanations to the decision-maker. This is particularly critical in the field of intelligent vehicles. For example, if the system does not provide explanations regarding why it takes a specific action,: then it is very hard for the users to accept the system. Moreover, in case of a serious event, relevant legal considerations need to be considered. Cognitive issues (or just a lack of information) may lead to the user not making the best of a fault, or operating unintentionally in a fault-tolerant mode that minimizes safety [7]. Properly explaining system behavior can help to solve these problems, improving the safety, robustness, and acceptability of intelligent agents. The demand for explainable AI in autonomous vehicle systems leads to many unanswered questions. For example: How much visualized information is enough? What is the best way to visualize complex neural network decision-making information? And, most importantly, how do we make sure that the visualization information is true and reliable?

Since the dawn of motor manufacturing, mankind has always been looking for safer, more intuitive ways to move. More intelligent autonomous vehicle systems have become one of the most important applications of AI and information technology advancements. Autonomous vehicle technology has become an important focus of research in the automotive industry [8]. Autonomous vehicles rely on the perception, control, and planning systems to achieve their capabilities and functions. AI technologies are widely used in this field to process information, drive behavior generation, and optimize control targets. With the application of deep learning, AI has shown strong capabilities in autonomous driving areas such as environment perception, behavior generation, and control system optimization which pushes autonomous vehicles closer to normal traffic conditions [3]. An AI-driven advanced driver assistance system is developed that can handle complex situations such as autonomous lane changes merging into traffic with sufficient safety margins. The method can be executed with

**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.

significant off-line planning time to ensure that safety remains paramount even in the occurrence of other road users making unpredictable or unsafe maneuvers.

## 2.1. Machine Learning and Deep Learning Basics

In this era, AI models also need to justify their outputs to be considered trustworthy and eventually replace human decision-making. For example, deep learning models frequently show a supremely accurate performance in complex tasks. But an important capability is built into deep learning models: transparency or reasonability or explainability. The key properties that need to be satisfied by trustworthy AI models are fairness, safety, and feature significance for credible decision-making. A predictor is said to be fair if the predictor makes the same estimate on the true label given any sensitive attribute (like race, gender, etc.) value [1]. In brief, any trustworthy AI study aims to gain human acceptance and trust by making accurate decisions honest, transparent, interpretable, and general."

"Machine learning (ML) is a field of computer science that uses statistical techniques to give computer systems the ability to "learn" with data, i.e., to progressively improve performance on a specific task without being explicitly programmed [5]. ML uses data-driven training where system processes patterns from the data and learns to make decisions instead of being explicitly programmed. Deep learning as a subfield of ML particularly shines in the data-rich regime. Deep learning algorithms automatically discover interactions from different features and learn highly nonlinear functions. For the same task, deep learning models usually require more data, computation, and training time compared to classical ML (statistical) models. But in many cases, they often show better performance than traditional ML models. An increasing number of AI systems are being developed, some of which have already been deployed in real-world operations such as in finance, medicine, and manufacturing.

## 2.2. Types of AI Models in Autonomous Vehicles

An innumerable amount of features characterizes the variability of the AVs devising operation and representation. Therefore, several state representations need to be used for representation of the system [6]. The total number of all taken actions will result in the method's full architecture. Features inform the man about the problem at hand and help to further investigate the action taken, when an explainer is augmented to the related decision-making module. Suppose, some other similar architectures are omitted, then the action space

**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.

of the AI module would have to be exchanged and all the subsequent dependant modules would require an overhaul. Crucially, the trained explainer network still has to observe the distribution of some features, to capture dynamics of mapping the features of a given state to the correct action.

An XAI-model, by definition, aims to show a human what it's done and how it is computed [5]. Basically, if a model uses being on the left or on the right as a decision, it is better to show a human right and left instead of a binary value. The choice is a function of the vehicle behaviour: speed, acceleration, orientation of the sensors, etc. The AI-modules in AVs currently use different kinds of models and some need to be explainable as well to gain transparent decision-making [9]. First, choosing the path and the action by looking for man-made marks to choose correct lane lines, prefer right over left-side to adhere to road signs and sometimes also considering the track borders. Second, state models serve the sole goal of dealing with a situation from the sensors to the decisions. While AVs are gaining prominence, explainability has taken the hot seat as a big issue when training them and understanding their decisions in various scenarios. It is safe to say that the proposed explainable AI for AVs is directly beneficial to state-of-the-art vehicle autonomy.

## 3. Explainable AI Techniques

Ref: 98a23868-d9de-439d-8867-78a3d2c03d1e Needless to say, from robotics to computational creativity, scientific theories to financial market predictions, all require the principles of explainable AI (XAI) to unfold their decision-making outputs. Often knowledge representation and reasoning (KRR) based XAI methods connect to a machine theoretics module which both controls and updates its internal axioms. Focusing on these tasks, the active inference principle in part offers regulatory significance, in part underpins an action selection procedure in its pursuit of optimality, and in part model prediction errors appear to direct plausible causes of them ultimately leading to transparency. The focus of the paper is to provide a framework for a transduction of the theory of motivated cognition, humans principles of inferential understanding and decision-making, to the design of AI systems that are naturally explainable with validity at both levels. A number of biologically-inspired techniques are thereby naturally related to active inference, providing a principled framework for flexible AI systems that can be less prone to adversarial attacks and have naturally improved explainability [10].

**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.

Ref: 32b6eeb5-a44c-4ef8-a721-32664723a739 As deep neural networks have gradually established their dominance across various application fields, we need to ensure that they are not treated as a "black box". Developing explainable AI (XAI) has therefore emerged as a necessary goal. For instance, when using AI for image/voice recognition, it is very necessary to have a strong background knowledge and have a deep understanding of the features that these models are focusing on for a particular task. By making use of the XAI techniques, decisions of AI systems can be made more transparent to people who may not be specialists in this domain. It also finds applications in shared control systems involving collaboration between intelligent agents and experts; AI models' explanations allow human experts to predict what the AI system will behave in particular situations and to decide when the AI system's decision is best intervened for optimum outcomes. To list another instance, lawyers and forensic analysts sometimes require interpretable explanations for representing evidence in the courtroom. AI systems, as part of forensic expert systems, can well classify whether an image has undergone tampering or not and, in turn, the XAI techniques can provide the rationale for image deployments [11].

### 3.1. Interpretable Machine Learning

One method to achieve this is to adopt a model that can explain itself. In this approach, first, a transparent predictive model (e.g., decision trees) is trained using a wide range of populations including biopsy images. Based on the predictions of this transparent model, demographic (clinically meaningful) models have been proposed to get interpretable models. The demographic model prediction only depends on the relevant features and the trained model has a good performance and at the same time has a feature explanation mechanism. Increments accused only primary myeloma were excluded from the analysis. EIM's approach is to overlay an interpretable decision rule that explains a black-box model's output based on high-level domain knowledge,where is the probability of the extracted attribute in the set A when the object was assigned to class CI, and the number of times the attribute is contained in the relevant class subset in the set A, can be relevantly obtained by effectively splitting the data set.

In the field of computer vision, it has been difficult to obtain interpretability on an image-by-image basis and to understand classifications. Some recent studies convert black-box predictions to white-box classifiers by learning to generate images that can be explained by a

**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.

white-box classifier. In the medical imaging field, predictive models are usually trained to predict image-level labels, making it difficult to obtain interpretable predictions for each pixel. For example, a chest X-ray classification model for pneumonia can detect lesions but it is difficult to judge what feature is used for Leisons detection. Therefore, in addition to biopsy images, images with some kind of interpretability are constructed from black-box models.

[12] Interpretable machine learning / interpretable models (IMs) emphasize feature interpretability over performance. When a decision needs to be explained to a user, effective explainability is more necessary than high performance. Nevertheless, most machine-learning models can be explained only on the basis of complex mathematical mechanisms.

### 3.2. Local Explanations vs. Global Explanations

To achieve transparency in BF-driven AFI and to give useful explanations about model decisions, we may use local, as well as global explanation methods [13]. Local explanations describe the prediction of an individual instance, which will help stakeholders to understand why they were personally affected. They typically focus on simple linear models in the vicinity of the instance to be explained. In other words, we zoom into a small region in the neighborhood of a single test case [14]. Global Explanation: These methods offer insights on the model as a whole. They typically aim to highlight important features, relationships or dependencies within the model. They do not only give you information about one specific prediction but are often used to understand which are the most important driving factors of the decisions your model makes. These are useful if you aim to verify and validate a global configuration of the system, or show the dominance of certain features within the input domain. Despite these differences, these two categories are used in the very same application areas, and we are going to focus on both of these topics due to their important characteristics and how much influence they have on the AV system, if they were to be implemented. [15].

### 4. Challenges in Implementing Explainable AI in Autonomous Vehicles

This transparent and interpretable real-time decision-making from AVs are hindered by AI safety, i.e., interpretability issues and ethical concerns, as per ethical standards, and concerns among the general public. For regulating and legal decisions based on adversarial cases, the AI systems should be safe, more transparent, interpretable, and explainable. AI methods, which are complying with the above-mentioned considerations for ensuring responsibility

**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.

with respect to ethical standards and laws and regulations, are mainly considered in the present paper. Thus, we aim at Understanding the safety aspects, which are to be agreed upon in order to make the intelligent, self-driving systems transparent in their operational and reactive decisions and ensure public trust in that way. Planning for regulation of those AI systems is also being included to ensure their operational and reactive decisions to be safety compliant and transparent. Understanding the global mapping of driving behavior of unintentional/unmanned vehicles and human drivers as well, so the vehicles would be justified to take that kind of actions as per the regulations for AVs [7].

Advanced driver assistance systems (ADAS) have shaped the latest state-of-the-art autonomous vehicles (AVs) for ensuring safer and comfortable mobility for passengers and other road users. These ADAS and AVs rely on the latest advancements in technologies such as artificial intelligence (AI), machine learning (ML), sensors, and data processing alongside real-time decision-making algorithms. Deploying deep learning techniques for handling vast, complex datasets has allowed these systems to exhibit improving performances during navigation, object detection, image and vision analyses, and automated decision-making functions [3]. However, during safety-critical decision-making scenarios, the AVs need to be transparent and explainable for public trust, regulatory compliances, and the safety of passengers and other road users. This requirement for real-time decision-making from ML systems is considered as the next AI safety frontier, to make the vehicles' AI/ML systems explainable, safe, and regulatory compliant.

### 4.1. Complexity of AI Models

Decision-making in dynamic safety-critical environments like traffic still poses substantial balancing acts with decisions often needing to be computed between conflicting goals such as short-term response and long-term future implications, efficiency and safety, convention and emerging pragmatics, as well as valid assumptions for reactive control behavior and exploratory actions (with yet unknown implications) for learning in long-term online decision-making (planning). Making planned decisions certain is computationally relatively straightforward, whereas demonstrating the decision-making process is often not a properly represented goal in efficient computational implementation [3]. This fact is reflected in the centrality of simplicity vs. inefficiency/complexity/trustworthiness trade-off everywhere. Autonomous vehicles—by being agents taking decisions—required balanced design 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.

transparent ML models with their computational models for planning (being hidden) and the corresponding reliable decision making should become a significant part of AI research for the (near) future.

Despite the promising prospect of deploying machine learning (ML) and artificial intelligence (AI) for autonomous vehicle decision-making, their increasingly high complexity leads to difficulty for stakeholders to understand or trust the systems sufficiently to deploy them [6]. Therefore, the usual perception of vehicles as complex and opaque "black-box" decision makers causes challenges in various respects including trust and safety. To address the aforementioned concerns and transparency and safety for transparent and safe (and hence ethical and security compliant) autonomous-vehicle decision-making, a paradigm shift is needed in AV research: the exploitability of ML/AI models used in autonomous vehicles (AVs) should be taken as crucial as the achieveability of high-fidelity models [16]. Explainable AI (XAI) is a multi-disciplinary paradigm area that addresses the design, use and development of AI models that can be effectively understood by humans. The design and development of AI models for transparent AI (XAI)-based AV decision making is an essential and crucial research topic, that remains basically unexplored.

### 4.2. Trade-off Between Performance and Interpretability

The discussion above assumes people can be informed enough about the physical properties of their vehicles to feel confident about their choices, but one development that is quite the opposite is owned-by-insurance services, in which dissimilar treatments can be the norm. When autonomous vehicles become normal, the self-insurability concepts of it may co-develop with explanations and justifications around the drivers of the (multinational) cars, and feedbacks from the potential risks that we plan for. The trade-off between transparency of how decisions were reached versus what insurance analysis should include in such decision-making should be central. An important one to us is can a mass-market vehicle in safe times feel confident enough to receive the prediction of intersection management? Retrospective dependence between explanations of how decisions were formulated do not provide us confidence in dealings with the parents; they should have explanations that build on forward-looking scenarios [17].

Autonomous vehicles are a key part of the next-generation transport system, expected to improve public transport cost- and energy consumption-share while allowing passengers to

increase productivity and safety, normally allocated to seeing, hearing, and sensing while driving. Since it's almost exclusively proprietary advances—mainly licensed from universities or other industries—there remains little transparency in addressing critical issues such as the trade-off between interpretability and performance. This situation brings with it a growing demand for finding methods that provide explanations of how a decision was achieved. Explainable AI (XAI) is the solution for reliable decision-making, making decision-making possible where justifications are reasonable. Given the explanation of decision-making, people can gauge autonomous vehicle actions in the same manner they gauge other vehicle and pedestrian actions. So, there will be demand for going beyond the development of AV models with just optimal physical performance in terms of interpretability. The integrated optimization not only of both performance and interpretability, but also of both models and sensor functionalities should become a research focus [18].

## 5. Applications of Explainable AI in Autonomous Vehicles

In the context of L4+ and L5 autonomous driving, Explainable AI clearly exhibits benefits with respect to trustability and transparent decision-making, ideal for infrequent but critical events, sensor object recognition in critical traffic situations, or system failures [19]. Arguably, additional benefits in terms of a safety metric through improved connection of perception and prediction in the web of cause-effect reasoning cannot be achieved, if the AI cannot explain its actions and predictions. In detail, X–AI can provide more relevant learning insights, particularly when common stimuli for mentoring are scarce under real operational conditions. Therefore, future supporting systems are also discussed in the following. In terms of Adaptive Learning capabilities, ethical or moral decision-making, and also natural language interaction, AI/3D graphics can, for instance, boost the guidance of transport professionals or even optimise the uptake of the respective transfer solutions (drone including operation as described). In terms of those exemplified fields, any automotive stakeholder being actively involved, i.e., individuals, vehicle or drone missions, are hence not only to be addressed; also, the interaction of the Adaptive Learning system to be supported as well as the supportive structures Jankres et al. can increase the user experience in Context-aware User Interface Systems.

Recent results of Explainable AI (XAI) demonstrate significant improvement on algorithm-encoded explanations, enhancing user understanding and trust levels – a cornerstone for the

acceptance and DAO compliance by the vehicle manufacturer or maintenance provider. Typically, users are sensitive about AI decision-making, in particular with regard to the execution of highly autonomous tasks. The authors have to provide the user with not only the decision made, also why and how an algorithm arrived to that decision. Moreover, the obtained insights can foster AI configuration learning without injecting the respective human reasoning skills. Also, the capabilities of drivers, passengers, operators et cetera can be taken into account (e.g., applications only causing trust-relevant explanations).

### 5.1. Safety-Critical Decision-Making

In the first level of decision-making, future-oriented planning decisions relate to not single AV components, but require full-stack AV modelling interdependencies. This is shown to be even more critical when the vehicle executes the interaction decisions in the future. The other two decisions, often considered Fail-Operational (FO), are closer to real-time planning decisions and can be viewed as peace-time interaction executions as per International Law of Armed Conflict, also operationalized as an additional set of both adversarial and forensical scenarios for model robustness testing within and outside their typical operational design domain [15]. Within Autonomous Vehicle systems, this work pertains to the road traffic conflict, extending that to adversarial interactions, and including additional considerations such as teleoperations.

Autonomous Vehicles (AVs) are integrating complex and high-performing Machine Learning (ML) algorithms into critical decision-making modules. This inherently introduces black-box behaviours across various AV systems, demanding improvements towards an increased level of transparency, interpretability, and explainability [2]. These qualities are essential for the Trust-and-Verify requirements of safety-critical, high-stakes AV decision-making and are required to hold bad actors accountable in the event of an accident. The AV must explain its decisions, typically under a metaphorical auditor's lens. Here, we present the three levels of AV decision-making from the 2017 Federal Automated Vehicle Policy (FAVP) and explain how AV decision-making can be translated across these levels: the passive safety decisions, such as trajectory tracks selected and acted upon years before completion, Fail Operational (Lane Assist) decision-making regarding the short-term path updates (Quarter-First), and the full-stack vehicle mission management decisions that are made over the far horizon (Trip-Whole) [16].

**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.

## 5.2. Regulatory Compliance

The call for legally informed AI and robotics technologies articulates the need for the development and implementation of explainability, an essential element necessary to understand why and how the various downstream responses occur in autonomous systems. The NHTSA, echoing the above technical definition of explainable AI, noted its importance as a way to ensure transparency, understandability, and security in autonomous vehicle policy. Hence, the regulatory compliance of a system that relies on the explainability of AI for ensuring transparency is at risk if AI development does not center around the design of trustworthy explainable models that facilitate a clear and transparent reasoning process [16].

The rapid evolution of AI technological capabilities is incumbent with a concomitant and urgent need for adequate regulation and policy to ensure safety and transparency regarding their operational decision-making. Within the last several years, regulatory frameworks have primarily focused on the actual operation of an autonomous vehicle, imposing stringent safety compliance or testing requirements. These regulatory and legal standards all address the output of the autonomous vehicle's system, but none explicitly address the decision-making process that enables that output [1].

## 6. Case Studies and Experiments

Explainable AI (XAI) advocates for rendering machine learning models transparent and interpretable while staying accurate. A white-box strategy allows the AI to be inherently interpretable. For safety-related applications, clear policy representation and intuitions about the main network behavior make deep learning models more human-centric and understood. Adding to that, interpretability can enhance user trust for the system and make the prediction output more acceptable. Hence, one of the main ambitions of this paper was to make an in-depth investigation of driving-oriented interpretability for various types of deep learning architectures and subsequently perform a comprehensive uncertainty analysis [9].

autonomous vehicles, interpretable decision-making model, Advanced Driver-Assistance Systems (ADAS), computer vision, uncertainty analysis Advancing self-driving cars promises significant improvements in road safety, driving comfort, and time savings. Autonomous driving development focuses heavily on the capabilities of neural-networkbased perception, planning, decision-making, and control modules. Within these algorithms, deep learning is

**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.

most prevalent, leveraging great expressiveness and superior performance. However, deep learning black boxes often raise interpretability, trust, and compliance concerns that are particularly pronounced in safety-critical driving scenarios [7]. To adopt self-driving technology in real-world complex environments, the ability to understand, debug, and efficiently analyze driving policy remains a crucial problem.

### 6.1. Real-world Use Cases

Indeed, the societal deployment of fully autonomous vehicles is multifaceted. Many technical challenges are still present—for example, the transition period during which human-driven transportation shares the streets with machines, the reassurance provided by regulations, policy, and the legal system. But one of the biggest challenges is the possibility of failures caused by the different feature interactions. In some domains, such as the medical field, human error can be considered as some kind of an explanation or summary of failures or unexpected side effects such as discovering a rare disease of the patient. In other words, in some domains, we develop methods to provide some predictions about different possibilities of the failures or how to recognize some risks from observation. In other cases, such as in autonomous vehicles, the forecasting, preventing, and recovery from errors side is not the only way to provide safe, efficient, and reliable systems, but a transparent reason for visiting off-nominal conditions is a mandatory requirement [20].

[1]Autonomous vehicles (AVs) and robots are transitioning from being opaque systems that take decisions in the bowels of their artificial intelligence software to becoming invitees in our physical environment and cognitive space, acting as teammates. The transparency of the decision-making process, or explainability, is fundamental to such transparency. Several attempts in the literature have aimed to provide better user–agent interaction in the form of intention, next moves, and trust [16]. In this article, we provide a brief survey of our research in Explainable AI (XAI) and the transparency of decision-making in autonomous robots and vehicles so as to improve our understanding of what is happening as the field deploys AVs.

### 6.2. Experimental Evaluations

In recent years, the development and deployment of autonomous vehicle systems have been a hot topic thanks to the major improvements in both hardware and compute power, the development of new deep learning-based methods, and the development of novel datasets.

**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.

Autonomous vehicle systems should have the ability to understand and follow the environment together with the transportation rules to reach a selected destination without causing any accident, and when they have controlled manually, the vehicle should have the capability to follow the human driver amendments without any risk. The promise of self driving automotive will revolutionize individual and public ground transportation networks around the world, but autonomous vehicles comes with systematic challenges and interesting social questions. Cutting across driving safety, efficient driving, fuel economy, traffic, the reliance on automotive in daily life, and public trust the transparency of behaviors for intelligent automated systems has a crucial impact moreover ensuring equivalent safety for all passengers are also ethical implications. Future autonomous vehicles operators legal, insurance entities discuss the further dissemination of this technology which may undo traditional set of legal liability in the event of an accident.

[21] [14] We continue our abovementioned discussion and experimentally evaluate a hypothetical self-driving car system with the help of the provided explanation. More specifically, three main contributions are introduced: (1) different classes of driving scenarios are defined and analyzed; (2) a new well-posed problem regarding a hypothetical self-driving system is defined and modeled to represent the problem according to machine learning techniques; and (3) different classes of driving scenarios in the context of which a hypothetical self-driving car system is evaluated. To do this, we first experimentally evaluate the average performance of the hypothetical self-driving car system when the correct maneuver (e.g., turn left, proceed straightforward, or slow down) is selected without the help of any explanation provided by the AI system. It can be observed that the average portion of the correct maneuver predictions of the hypothetical self-driving car system is lower than the expectations. Next, averaged performances of the provided explanation technique (Proxy) and a state-of-the-art alternative explanation technique (Lime) for the hypothetical self-driving car system are compared and it is observed that the Proxy explanation technique generally provides the most accurate and useful explanations in the high-level scenario modeling of driving scenarios. In contrast, we found out that it is not easy to experimentally define the best explanation paradigm through simulation environments.

## 7. Ethical and Legal Implications

**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.

xAI is an extended, multi-disciplinary domain, requiring efforts and active involvement by different disciplines. xAI systems need to provide understandable reasoning that can be assessed by stakeholders involved and affected by the applications of these systems. xAI's transparency, fairness, explanation, and opinion generation systems will largely act on stakeholders in the case of non-compliance with legal conditions and ethical norms as set by different legal jurisdictions [22]. In this development process, fundamental changes will also have to be implemented in the relevant traditional liability systems. In addition to providing technical measures with various models and algorithms to explain decisions taken by machines, the juridical-legal measure should maintain the axiomaticity of the operation of AI-based systems. Given the risks to be managed, it is in this direction that extensive developments aimed at understanding and ensuring that for decisions taken by AI, responsibility will be handled not only in the phase of accidents but with continuous transparent operation and in a manner compatible with a sociotechnical mindset are required. Unfortunately, developments in the established legal systems are not yet visible in these terms. At present, the only solution is that AI's legal responsibility be managed by the general principles and solutions which the protection of rights and freedoms ensured where there is a shortcoming is considered as the basis of protection [23].

Explainable AI (xAI) systems are likely to modify the decision regimes in relation to decision-making systems in dealing with complex problem domains [24]. Particularly in autonomous decision-making systems, societal attention is captured by the opaque nature of this technology. The acceptance, trust, understandability, and usability of such systems in society is significantly influenced by the question of explainability of AI-based outputs. Ethical and legal discussions revolve around how much of the decision is opaque and to what extent this opacity can be tolerated by society in public domains and by affected stakeholders. A layered regulation model including a set of ethical assumptions and thorough legal guidelines is expected to substantially contribute to creating trust and acceptance in society and to avoiding techno-moral dilemmas. A proper regulation of xAI is, hence, fundamental for explaining decision-making results of AI-based systems and for fostering trust and confidence in such systems in society.

### 7.1. Bias and Fairness

**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.

One chapter of the book Levchugin [25] deals with the Birman-Schwinger principle (BSP), introduced in quantum mechanics in 1960 to describe the interaction between stable sizeable objects. It has later been adopted, and applied to different domains of physics. In 2018, the author of this chapter extended the BSP from quantum mechanics and electrodynamics to the linearized Einstein's general relativity. The resulting system of linear non-local differential equations describes the propagation and the scattering of waves by cylindrical objects (wormholes) immersed in Minkowski space-time. If the cylinder is placed around a black hole, and vibrates, then the equations describe the nonlinear interaction between the slowly varying living matter and electric gravitational waves. This is the first rigorous model implementing the idea of surrounding the black hole by thin horror shell, closest to the black hole woman.

Moreover, XAI cannot be implemented in an off-the-shelf manner. Instead, XAI solutions should consider multiple requirements, including the portability of the AI model, in-application and small-footprint solution, the capacity to update the model and restore operability in case of faults or updates [1]. An alternate strategy is to have a trusted AI model that is "human informed" with individual human intentions and utilities. Care should be taken when deploying such an approach because the personalized objective might be unfair towards other agents. Essentially, XAI can shift the moral blame, leading to additional issues. Given the variety and complexity of the possibilities, the standardization of the XAI models may be the most sustainable approach.

### 7.2. Privacy Concerns

Users should not only interact with these AVs but also actively participate in the decision-making process, in the form of trading throwbacks and instructing on new actions. At the same time, privacy principles should also be taken into account so that the users' privacy is not breached, either by external malicious entities or other co-travelers. Thus, users must understand the AI-driven processes employed in an AV's decision-making and must feel wholly in control of, and equal to, the AV-user relationship. For example, the interface of an intelligent off-road vehicle (IOV) may provide a visual feedback module that demonstrates how the environmental knowledge module of each IOV in a congregation arrived at their individual decisions. Hence, this module can be expected to highlight the intention of each IOV. According to the beneficiary perspective, this feeds into the public accountability 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.

AIOV decision-making [24]. Editorial note: The reference information is incomplete and the reference needs to be listed in full in the referenced article section.

Autonomous vehicles and intelligent sensors, along with their support infrastructure networks, have brought about a new wave of concerns and anxieties related to privacy, including those raised in related studies [16]. These concerns are far-reaching and pertaining to various aspects, such as user-level privacy, security, surveillance, and social acceptability. To this end, privacy is one of the determinants of technology adoption in any industry, including transportation. The adoption of a personal autonomous vehicle means the individual becomes a part of the smart city, and always remains under observation. Hence, privacy-related acceptance is important in addition to the various other dimensions of technology acceptance.

## 8. Future Directions and Research Opportunities

Self-explaining AI is important because 1) AI applications need to be trusted by users, 2) there's a growing need for personalization, 3) AI decisions need to be transparent to be accountable, 4) AI oversight needs to remain understandable by developers, and 5) AI should be able to explain fallback plans to humans, in case handover is needed (and vice-versa). We provide a brief overview of different techniques for building explainable AI, encompassing different levels including the model's architecture, the model's opacity (i.e., built-in or post-hoc explanations), and the goodness of fit between the model's explanation and human cognition [15]. Despite significant headway, trade-offs still remain among model complexity, model performance, and model explainability, making XAI a lively and essential area of research going forward. Therefore, making sure AI explanations are perfectly acceptable is a crucial cognitive task.

Explainable AI (XAI) is an increasingly popular research field aiming to increase the transparency and interpretability of AI models so that their decisions are understood by humans. This is crucial in system critical areas like medical diagnosis, fraud detection, and legal decision-support systems, among many others. In this paper, we focus on the use of XAI techniques in autonomous vehicles in order to explain AI decisions in closed-loop scenarios rather than to interpret feature importance in order to increase the performance of traditional learning models [26]. Specifically, we discuss XAI approaches focusing on i) rule-extraction from convolutional neural network classifiers, ii) sensor desensitization in LiDAR pre-

**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.

processing, iii) Abductive Reasoning to obtain pseudo-GPS velocity, and iv) grammatical inference to create formal model within a complete autonomous vehicle system. We hope these case studies demonstrate some of the various existing and ongoing directions in explainable learning for autonomous vehicles that can guide future works and new ideas [2].

The output text has to refer to f4e0c730-a7ba-48ee-a5e1-089cfed90dc7 and f49abc52-87ed-456d-b164-3049643fe0b9.

## 8.1. Advancements in Explainable AI Technologies

Although the "Transparency" might be considered a sub-aspect of the idea of explanation, at many places, these two concepts were considered separately since two should be used as the end to not level off new Technological change. The importance of explicating the different between both of them, also in the same area of job opportunities, has been mentioned earlier can, further, evidently impact law making and its legal code for; full framing along with the Accidents, Technology and Legal law's basis and intellectual property & PYAs Advantages legal right in order to respected the other. Except for viewing it with legal eyes; however, the concept of "Transparency" plays the meaning here for the business field as, a concept that includes consideration of transparency as a broader concept in the business, law resources, comprising: Transparency has been defined by UNESCO as an important aspect of business civil life. And it is very identical to good governance, too [17]. A complete relationship between these two strategies and Technology Law's transparency on one side relates to the algorithm while pursuing this programmed code, which transparency mode of operation and management of technology field in resources, analyzed as well the Transparency and its desirable laws steps in abovementioned arrangement were the pointed indicators of both research core report tomatinaflrmation ratification, perspicuity, following the narrative in all of the contents.

Explainable AI (XAI) pertains to AI systems' ability to consistently provide high-quality explanations, and at the same time, these explanations will be understandable to humans, depict the system's internal processes accurately, and communicate knowledge boundaries [15]. Inbrief, explanations provided by the AI will be meaningful for creator and end-users. Some of the most known advancements in XAI are LIME (Ribeiro et al., 2016), SHAP, Deep SHAP (Lundberg & Lee, 2017) and the Attention Mechanism. LIME is a technique for transparently uncovering what an AI model uses as a basis for predictions on a certain input

example. The benefits of this method are that it can be used without modifying the model and can work with any uninterpretable machine learning model, linear and non-linear. As it's been mentioned as the third type of mechanism, the doctrinal rule first appeared in Industrial Revolution's era concerning human behavior previously (French Revolutionary Wars (France): Civil Code of 1804 regulating the relations between people public, employees and a finished work, stand as far from it as it is acceptable), thereafter the input column goes with the company's context and the output result goes with the job market's observations. In the whole there has been the resource, that has been, in general, viewed as the most involved one, neglect those flow in which the elucidation by "Transparency" and "Explanatory" can always be opt for and accordingly play necessary roles, nonetheless.

### 8.2. Integration with Human Factors

Three key levels, in increasing order of resistance against outside intervention, are identified: (1) Manual driving: full human control. (2) Supervision: Human monitoring of the environment and verification of the system's decisions. (3) Shared autonomy: Involves a more continuous or inherent interaction between human and machine, with potentially entangled motion. As shared autonomy moves from a supervision level to a manual level, the car gradually 'hides' the potential safety opportunities and constraints from the human. For that reason, thought must be given to how the level of system autonomy (and so transparency) can change with the operational context and user experience. Furthermore, careful consideration of the transparency of each of these elements will suggest how much information should be exchanged when control is handed off between driver and system at the interface. For this purpose, such feedback mechanisms are particularly important in low and shared autonomy, and proposed for the future development lacking pattern recognition capabilities. For (high) shared autonomy in low visibility conditions, interaction may not be limited to safe navigation. In extending out of bounds motion, continual interaction may be an important contender in harmonic physical human robot interactions as a factor to be considered in a transparent system in these environments.

Explainability increases the trust that users, drivers, and other road users place in autonomous vehicles and their technologies [20]. The ideal situation is that systems either make decisions explicit, or that they can provide explanations when they (or their decisions) are interrogated [27]. This aspect is important when designing systems for shared autonomy and in creating

**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 that handle fault tolerance with more resilience. In this chapter, judicious choices of reactions to the environment (both feedback responses and operational distances) keep the human in the loop, and present situations that require human awareness. Thus, this capacity for interacting with and handing over to humans naturally lead to transparency; how human awareness and input can intervene in an autonomous system at various levels [1].

## 9. Conclusion

At the moment, different stakeholders, including victims or users, courts or regulators, might have different purposes in understanding the ML-based system's decision [28]. As the result, transparency adversarial scenarios are arising according to which the explanation is considered by one of the stakeholders for his interests but nobody has incentives to provide cognitive capabilities to this decision explanation framework which should be done by design. Indeed, explain the UVAV decision simply with very few raw feature can be challenging and potentially ambiguous. It is hard to provide simple yet accurate explanations while dealing with the trade-offs between true negatives and true positives. This is necessary to acknowledge that these conflicting goals should be properly balanced according to the AV scenario where XAI systems are employed. In fact, false positives, as well as false negatives, need to be minimized.

"Explainable AI was long believed to be roughly divided into two main areas, model explanation and decision explanation which means that we care to explain decisions taken by the model, and the irrelevant aspects of the input in which the model prediction relies on [15]. However, the definition of XAI is far from being final and the status of the art is messy and incomplete. At AV level, the main scope we are interested in is to know the causes of different valuations of input features observed at prediction time, i.e. the importance in training of input data points. This is because decision transparency is crucial for blame management in case of undesired results generated by XAI systems, i.e. in case that XAI models are the cause of an AV crash. Of course the system configuration and the specific scenario are also relevant in this sense and it is extremely important to guarantee full clarity on these aspects too [26]. This is to have clear from the beginning what is and what is not managed by the XAI system, which is a crucial aspect about explainability in AV systems. In this scenario, the main goal of this work is proposing a new metric to assess XAI system understanding in AV context."

**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.

**Reference:**

1.  Tatineni, Sumanth, and Anjali Rodwal. "Leveraging AI for Seamless Integration of DevOps and MLOps: Techniques for Automated Testing, Continuous Delivery, and Model Governance". Journal of Machine Learning in Pharmaceutical Research, vol. 2, no. 2, Sept. 2022, pp. 9-41, https://pharmapub.org/index.php/jmlpr/article/view/17.

2.  Prabhod, Kummaragunta Joel. "Advanced Machine Learning Techniques for Predictive Maintenance in Industrial IoT: Integrating Generative AI and Deep Learning for Real-Time Monitoring." Journal of AI-Assisted Scientific Discovery 1.1 (2021): 1-29.

3.  Tatineni, Sumanth, and Venkat Raviteja Boppana. "AI-Powered DevOps and MLOps Frameworks: Enhancing Collaboration, Automation, and Scalability in Machine Learning Pipelines." *Journal of Artificial Intelligence Research and Applications* 1.2 (2021): 58-88.

**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.