www.ijernd.com

Explainable Deep Learning Models for Autonomous Vehicle Decision-Making

Sai Kiran | Sanketika Institute Of Technology and Management

**Abstract** 

Autonomous vehicles (AVs) rely heavily on deep learning models for perception, planning, and control.

While these models achieve high accuracy, they operate as black boxes, making their decision-making

process opaque. Lack of interpretability in safety-critical environments limits trust, hinders

debugging, and complicates regulatory approval. This paper proposes an explainable deep learning

framework for AV decision-making by integrating state-of-the-art perception and control models

with post-hoc explainability techniques, specifically SHAP (Shapley Additive Explanations) and

LIME (Local Interpretable Model-Agnostic Explanations). Using simulated driving scenarios in

the CARLA simulator and real-world datasets such as KITTI and nuScenes, the framework generates

actionable explanations for vehicle actions, including lane changes, braking, and steering. Our results

demonstrate that XAI techniques can highlight critical features influencing decisions, uncover model

biases, and assist developers in improving AV reliability. The proposed framework enhances safety,

transparency, and accountability, providing a practical path toward regulatory-compliant

autonomous systems.

**Index Terms** 

Explainable AI, Autonomous Vehicles, Deep Learning, SHAP, LIME, CARLA Simulator, Safety-

Critical AI.

Introduction

Autonomous vehicles are at the forefront of

intelligent transportation, offering the promise

of enhanced safety, reduced traffic accidents,

and improved mobility. Modern AVs rely on

deep neural networks (DNNs) for tasks such

**Perception:** Object detection, semantic segmentation, and lane identification.

Planning: Path generation, obstacle

avoidance, and trajectory prediction.

as:

• Control: Steering, throttle, and braking commands.

Despite impressive performance, **DNNs are opaque**. For instance, a network may decide to **swerve suddenly or brake abruptly**, but
engineers cannot determine whether the
decision was triggered by a **genuine obstacle**, **sensor noise**, or **spurious correlations** in
training data. This lack of interpretability raises **critical safety concerns** and impedes trust
among regulators, engineers, and the public.

**Explainable AI (XAI)** offers techniques to generate **human-understandable insights** into black-box models. Two widely used post-hoc methods are:

- SHAP (Shapley Additive
   Explanations): Assigns feature
   importance scores using game-theoretic principles to quantify
   contributions of each input feature.
- LIME (Local Interpretable Model-Agnostic Explanations): Creates local surrogate models by perturbing inputs to approximate model behavior near a specific decision.

Applying these techniques to AVs can provide actionable insights into why a vehicle performs a certain maneuver, helping to identify biases, incorrect reasoning, or model weaknesses.

## **Challenges in AV Explainability:**

- High-dimensional input from multimodal sensors (cameras, LiDAR, radar).
- Temporal dependencies in sequential decisions.
- Need for real-time or near-real-time explanations in simulation or operation.

## **Contributions of this work:**

- Integration of SHAP and LIME with AV perception and control pipelines.
- Analysis of simulated CARLA scenarios and real-world datasets
   (KITTI, nuScenes).
- Quantitative and qualitative evaluation
   of model interpretability and
   decision transparency.

 Insights for safety verification and trust-building in autonomous driving systems.

#### **Literature Review**

### A. Deep Learning in Autonomous Vehicles

Deep learning has transformed AVs by enabling end-to-end perception and control:

- CNNs: Used for image-based object detection, segmentation, and lane recognition.
- RNNs/LSTMs: Capture temporal patterns for sequential decisionmaking.
- Transformers: Emerging for multimodal temporal reasoning, integrating camera and LiDAR data.

AV systems such as **Tesla Autopilot**, **Waymo**, and **OpenPilot** demonstrate real-world applicability but remain **black-box systems**, which limits interpretability.

## B. Explainable AI (XAI) Techniques

XAI provides tools to make models transparent:

- LIME: Perturbs input features (e.g., pixels, sensor readings) and fits a local linear model to approximate decision boundaries. Useful for understanding single-instance predictions.
- SHAP: Assigns Shapley values to input features based on their contribution to the model's output.
   Captures global and local importance and is mathematically grounded in cooperative game theory.

Applications of XAI span finance, healthcare, and autonomous systems, with focus on trustworthiness, debugging, and regulatory compliance. In AVs, XAI can highlight which objects, lanes, or environmental cues influenced a vehicle's action.

#### C. Explainable AV Research

Prior work in explainable AVs includes:

- Saliency maps: Highlight regions of camera images influencing control decisions.
- Attention mechanisms: Visualize which features the model "attends to" during prediction.

Feature importance analysis:
 Quantifies the influence of specific inputs (e.g., LiDAR points, lane markings).

Challenges remain in real-time applicability and integration with multi-sensor inputs. Our work addresses these by combining SHAP and LIME with both perception and control models, enabling comprehensive explainability across the AV pipeline.

#### **Problem Statement**

Autonomous vehicles face **critical challenges** due to the opacity of deep learning models:

- Lack of Interpretability: Engineers cannot verify whether actions are based on valid features.
- 2. Safety and Compliance Risks:
  Without explainable reasoning, AVs
  cannot guarantee safety in edge cases
  (e.g., unexpected pedestrians, poor
  weather).

## **Research Question:**

"How can post-hoc explainability techniques (SHAP and LIME) be effectively integrated with AV deep learning pipelines to provide interpretable, actionable insights for real-time decision-making without compromising performance?"

The goal is to provide **both qualitative and quantitative explanations** of vehicle actions

to improve **trust**, **debugging**, **and regulatory compliance**.

## Methodology

### A. System Overview

The framework consists of four primary components:

- Data Acquisition: Multi-modal sensor data (camera, LiDAR, radar) from CARLA, KITTI, and nuScenes.
- 2. Model Training:

CNN/LSTM/Transformer models for **perception and control**, predicting object locations, lane positions, and control commands.

- 3. Explainability Module: Apply SHAP and LIME to generate feature importance and saliency maps for decisions such as braking, lane changing, and obstacle avoidance.
- 4. Visualization & Analysis: Visualize explanations using heatmaps, feature rankings, and temporal sequences to interpret AV actions.

#### **B.** Datasets and Simulation

 CARLA Simulator: Generates urban, highway, and adversecondition scenarios. Provides full control over environmental conditions and ground-truth labels.

- KITTI Dataset: Real-world driving data with camera images and LiDAR for object detection and depth estimation.
- nuScenes Dataset: Multi-modal dataset including 360° perception, dynamic objects, and vehicle states, supporting complex driving scenarios.

**Data Preprocessing:** Align multi-modal data temporally, normalize sensor inputs, and convert LiDAR point clouds into voxel grids for CNN input.

## **C. Model Implementation**

- Perception: CNN-based networks for object detection and semantic segmentation.
- Control: LSTM/Transformer models predict steering, throttle, and braking.
- Loss Functions: Cross-entropy for object detection; mean squared error (MSE) for control commands.
- Training: Conducted on GPUenabled TensorFlow/PyTorch environments, with hyperparameter tuning (learning rate, batch size, sequence length).

## D. Explainable AI Integration

 LIME: Perturbs input frames/features to approximate local linear models, providing explanations for individual vehicle decisions.

- SHAP: Calculates Shapley values for multi-modal inputs, ranking feature contributions globally and locally.
- Explanations highlight **critical features**, such as pedestrian positions,
  lane markings, and traffic lights,
  affecting vehicle actions.

#### **E. Evaluation Metrics**

- Quantitative: MAE (steering/throttle/brake), IoU (object detection), trajectory deviation.
- Qualitative: Expert evaluation of interpretability, clarity, and usefulness of SHAP/LIME visualizations.

#### **Results and Discussion**

#### A. Quantitative Results

Model	Steeri ng MAE	Throttl	Ü	IoU (Segmentat ion)
CNN+LS TM	0.12	0.08	0.07	0.82
Transfor mer	0.09	0.06	0.05	0.85

 Transformer-based models perform slightly better, producing more stable

**predictions** that benefit interpretability.

Reduced MAE indicates more
 accurate vehicle actions, which enhances trustworthiness of explanations.

## **B.** Explainability Analysis

- SHAP Heatmaps: Highlight which features (lane lines, pedestrians, obstacles) contributed to a decision.
- LIME Analysis: Illustrates how perturbations in the input affect predicted actions, revealing temporal dependencies.
- Identified model biases, e.g., overreliance on lane markings over dynamic objects.
- Experts found explanations intuitive, useful for debugging and scenario analysis.

## C. Real-Time Performance

- LIME: ~180 ms per frame
- SHAP: ~220 ms per frame

 Both methods are feasible for offline evaluation; potential exists for realtime integration with optimization.

#### Conclusion

This research presents a comprehensive framework for explainable deep learning in autonomous vehicles:

- Demonstrates integration of SHAP and LIME for perception and control decisions.
- Evaluated on simulated CARLA scenarios and real-world datasets (KITTI, nuScenes).
- Provides both quantitative and qualitative insights into vehicle behavior.
- Enhances trust, safety, and debuggability of autonomous systems.

Explainable AI is critical for regulatory compliance, public trust, and safe deployment of AVs.

#### **Future Work**

- Real-Time XAI Deployment:
   Optimize SHAP/LIME for live autonomous driving.
- Multi-Modal Integration: Combine camera, LiDAR, and radar explanations.
- User-Centric Explanations: Translate feature importance into natural language explanations for operators.
- Adversarial Robustness: Evaluate explanations under sensor noise and attacks.
- Adaptive Learning: Continuous improvement with online feedback and explanations.

#### References

 Doshi-Velez, F., Kim, B., "Towards a Rigorous Science of Interpretable Machine Learning," arXiv preprint, 2017.

- Lundberg, S. M., Lee, S.-I., "A Unified Approach to Interpreting Model Predictions," NeurIPS, 2017.
- 3. Ribeiro, M. T., Singh, S., Guestrin, C., "Why Should I Trust You?"

  Explaining the Predictions of Any Classifier," *KDD*, 2016.
- Dosovitskiy, A., et al., "CARLA: An Open Urban Driving Simulator," *1st* Annual Conference on Robot Learning, 2017.
- Caesar, H., et al., "nuScenes: A multimodal dataset for autonomous driving," CVPR, 2020.
- Geiger, A., Lenz, P., Urtasun, R., "Are we ready for autonomous driving? The KITTI vision benchmark suite," CVPR, 2012