

Smartphone Sensor Insights for Pothole Detection and Road Quality Evaluation using Machine Learning

M. Shashidhar¹, Vuppula Manohar², Appala Sravan Kumar³, P. Kiran Kumar⁴, Ch. Anil Kumar⁵

¹Department of Electronics and Communication Engineering, Vaagdevi College of Engineering, Warangal, 506005, Telangana, India. sasi47004@gmail.com

²Department of Electronics and Communication Engineering, Vaagdevi Engineering College, Warangal, Telangana, 506005. manoharvu@gmail.com

³Department of Electronics and Communication Engineering, Vaagdevi Engineering College, Warangal, Telangana, 506005. appala.sravan@gmail.com

⁴Department of Electronics and Communication Engineering, Balaji Institute of Technology and Science, Narsampet, Warangal, Telangana, 506132. kiranpadakanti0430@gmail.com

⁵Department of Electronics and Communication Engineering, Vaagdevi College of Engineering, Warangal, 506005, Telangana, India. anil.chidra@gmail.com

Corresponding: M. Shashidhar (sasi47004@gmail.com)

ARTICLE INFO

Received: 05 July 2023

Revised: 12 Aug 2023

Accepted: 25 Aug 2023

ABSTRACT

This research addresses the critical challenge of road infrastructure degradation and public safety by developing a robust machine learning framework for automated road condition classification. Utilizing data consolidated from multiple labeled sources, the research classifies road quality into three distinct categories: Potholes, Good Roads, and Bad Roads. To ensure data integrity, the dataset underwent rigorous preprocessing, including the handling of missing values and categorical label encoding. A significant contribution of this work is the application of the Synthetic Minority Over-sampling Technique (SMOTE) to mitigate class imbalance, ensuring the models were trained on a representative distribution of all road conditions. Exploratory Data Analysis (EDA) validated the structural consistency of the data post-balancing. The processed dataset was partitioned using a 70:30 training-to-testing ratio to evaluate the performance of the Extra Trees Classifier (ETC) and the K-Neighbors Classifier (KNN). Empirical results demonstrate exceptional predictive power; the Extra Trees Classifier achieved an accuracy of 99%, while the KNN Classifier reached an optimal accuracy of 100%. These results were further validated through confusion matrices and classification reports, confirming near-perfect precision, recall, and F1-scores. This research provides a highly reliable, scalable pipeline for real-time infrastructure monitoring, offering municipal authorities a data-driven tool to prioritize road maintenance and enhance commuter safety.

Keywords: Road Quality Assessment, Pothole Detection, Machine Learning, SMOTE, Extra Trees Classifier, Infrastructure Safety.

1. Introduction

India possesses one of the world's most extensive road networks, spanning over 6.4 million kilometres. While national highways constitute only 2% of this network, they facilitate over 40% of total traffic [1, 2]. This immense pressure leads to rapid infrastructure degradation; recent data from the Ministry of Road Transport and Highways attributes over 10,000 annual fatalities to poor road conditions and potholes. Traditionally, road monitoring has relied on manual inspections and public complaints—processes that are inherently reactive, labour-intensive, and prone to human error. Such inefficiencies lead to delayed maintenance, increased vehicle wear, and heightened

safety risks, necessitating a transition toward automated, real-time monitoring solutions [3, 4]. The primary motivation for this study is the urgent need for scalable, data-driven infrastructure management. Manual surveys are no longer viable for vast urban and rural networks. By leveraging machine learning (ML), it is possible to transform raw sensor data into actionable insights, enabling predictive maintenance rather than reactive repairs. This project addresses the critical gap in current systems by providing an intelligent framework that reduces reliance on manual labor, minimizes fuel consumption through smoother transit, and ultimately saves lives. Such a system is a foundational component for smart city initiatives and the safe deployment of autonomous vehicle navigation. This research implements a classification pipeline to categorize road surfaces into three states: Potholes, Good Roads, and Bad Roads. A critical phase of this study involved EDA to evaluate data distribution and integrity. Initial assessments revealed significant class imbalances a common hurdle in infrastructure datasets which were systematically addressed using SMOTE. By utilizing tools like Seaborn for count plot visualization, the dataset was balanced to ensure model impartiality. This rigorous preprocessing ensures that the subsequent deployment of classifiers, such as the ETC and KNN algorithms, yields highly reliable and generalizable results for real-world applications.

2. Literature Survey

The development of automated road surface defect detection has transitioned through several technological paradigms, moving from manual observation to sophisticated computational frameworks. Current research primarily bifurcates into image-based vision systems and sensor-based non-image solutions.

2.1 Computer Vision and Deep Learning Approaches

Vision-based detection remains the most prevalent area of study due to the rich spatial data provided by cameras. Early reviews by Ma et al. [4] and Koch et al. [13] categorized algorithms into traditional image processing and machine learning, noting that while traditional methods work in controlled environments, they struggle with varying lighting and textures. Recent advancements have favoured Deep Learning (DL). Aparna et al. [8] and Park et al. [14] demonstrated the efficacy of the YOLO (You Only Look Once) architecture for high-speed, real-time detection in autonomous vehicles. To further refine these models, semantic segmentation techniques [9] and attention-aggregation mechanisms [17] have been employed to improve pixel-level accuracy and model generalization across unseen environments.

2.2 Optimization for Edge and Mobile Deployment

A significant trend in recent literature is the optimization of models for "AI-on-the-edge." Recognizing that high-computational clusters are not always feasible for on-road deployment, Asad et al. [2] and Xing et al. [5] proposed lightweight models and binocular stereo vision to balance accuracy with low power consumption. This shift ensures that detection systems can be integrated into mobile devices or embedded vehicle hardware without significant latency.

2.3 Non-Image and Multi-Sensor Data Fusion

Beyond visual cues, researchers have explored the use of vibration and thermal data. Singh et al. [6] utilized Long Short-Term Memory (LSTM) networks to process smartphone accelerometer data, capturing the temporal features of road anomalies. Similarly, the integration of GPS and cloud-based reporting [10, 20] has enabled participatory sensing, where crowdsourced data aids municipal planning. Advanced methods have even incorporated infrared thermography [16] to detect pavement distress in low-light conditions, proving that multi-modal data fusion significantly enhances system resilience compared to single-sensor approaches.

2.4 Challenges and Research Gaps

Despite high accuracy in experimental settings, real-world deployment faces challenges such as domain shift and environmental noise. Dhiman and Klette [18] and Bucko [19] highlighted that weather variations and image noise can drastically reduce model reliability. Furthermore, many existing studies focus on binary classification (Pothole vs. No Pothole), leaving a gap in multi-class evaluation—such as distinguishing between "good" and "bad" non-pothole surfaces—which this research aims to address through its three-tier classification framework.

3. Proposed Methodology

The proposed system architecture is designed as an end-to-end machine learning pipeline that transforms raw multi-source sensor data into a high-precision road quality classification tool. The methodology is divided into several discrete stages to ensure data integrity and model generalizability.

3.1 Data Acquisition and Preprocessing

The primary data consists of sensor-based metrics (accelerometer/gyroscope) stored in CSV formats across three specific classes: Pothole, Good Road, and Bad Road. The pipeline begins with a Dataset Consolidation phase, where multiple disparate files are merged into a unified global DataFrame. Each entry is programmatically tagged with a categorical label based on its source folder.

During preprocessing, the system executes:

- **Data Cleaning:** Systematic identification and removal/imputation of missing values.
- **Feature Transformation:** Application of Label Encoding to convert categorical class targets into a numeric format compatible with mathematical modeling.

3.2 Feature Analysis and Class Balancing

Before model training, EDA is conducted to identify underlying patterns and class distributions. To counter the "Class Imbalance" problem—where certain road conditions are under-represented—the system employs SMOTE. Unlike simple duplication, SMOTE generates synthetic examples in the feature space based on nearest neighbors, ensuring the models do not develop a majority-class bias. The dataset is then partitioned using a 70:30 stratified split, maintaining the balanced class ratios in both training and testing subsets.

3.3 Model Development and Evaluation

The core of the system utilizes two distinct algorithmic approaches to ensure comparative validation:

1. **Extra Trees Classifier (Extremely Randomized Trees):** An ensemble learning method that fits several randomized decision trees on various sub-samples of the dataset to improve predictive accuracy and control over-fitting.
2. **K-Neighbors Classifier (KNN):** An instance-based learning algorithm that classifies road conditions based on the proximity of feature vectors in a multi-dimensional space.

The performance of these models is quantified through a multi-metric evaluation strategy comprising Accuracy, Precision, Recall, and the F1-score, ensuring a reliable assessment of the system's ability to minimize false negatives in pothole detection.

3.4 System Architecture

The following diagram illustrates the high-level flow of the proposed system, from the raw data input to the final classification output.

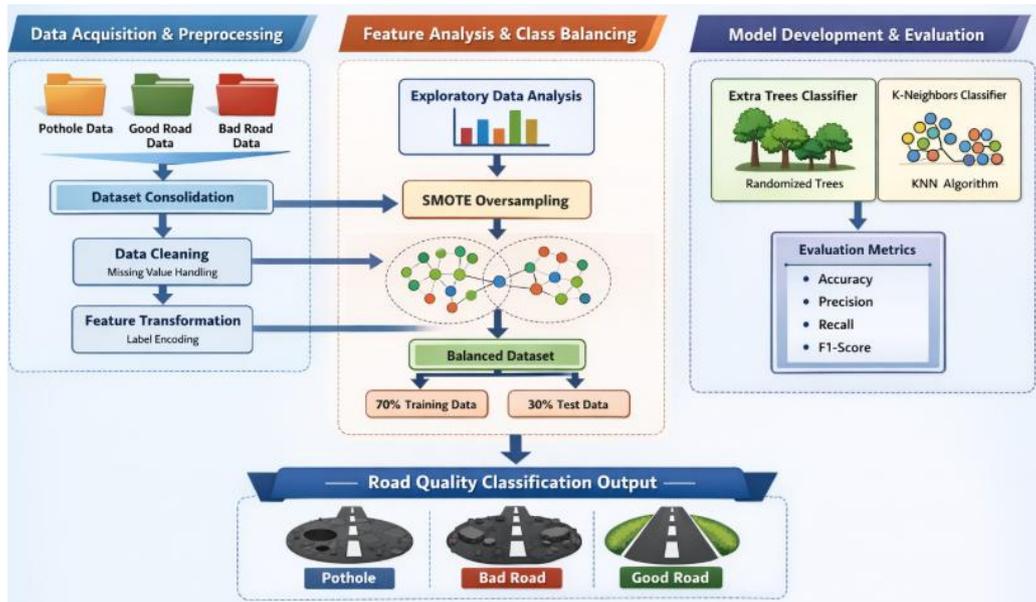


Fig. 1: Proposed system architecture.

4. Results And Discussion

The dataset utilized in this study comprises high-frequency sensor readings captured via smartphone-integrated Inertial Measurement Units (IMUs) and Global Navigation Satellite System (GNSS) receivers. These features provide a multi-dimensional representation of vehicle dynamics and spatial positioning.

Table 1: Technical Feature Characterization.

Feature Category	Attribute	Description & Functional Utility
Temporal	Timestamp	Provides a temporal reference for sequential data analysis and ensures synchronization between disparate sensor streams.
Geospatial	Latitude, Longitude	GNSS coordinates in decimal degrees, facilitating the geospatial mapping of detected anomalies for municipal maintenance.
Kinematic	Speed	Records instantaneous velocity; critical for normalizing sensor spikes that vary based on vehicle speed.
Inertial (Accel)	Accel_X, Y, Z	Tri-axial acceleration (m/s ²). The Y-axis specifically captures vertical displacement, a primary indicator of pothole impact.
Inertial (Gyro)	Gyro_X, Y, Z	Angular velocity (rad/s). These features capture the rotational pitch, roll, and yaw of the vehicle when traversing uneven surfaces.
Target Variable	Class	The ground-truth categorical label: Pothole, Good Road, or Bad Road.

Sensor Modalities and Signal Significance

The core of the detection logic relies on the fusion of Accelerometer and Gyroscope data. While the accelerometer measures linear forces capturing the sharp vertical "jerk" associated with road defects, the gyroscope monitors

angular changes. This is vital for distinguishing between a true pothole and a controlled manoeuvre, such as a lane change or turning. By combining these inertial readings with Geospatial data, the system moves beyond simple detection to create a "Live Road Quality Map." The inclusion of Speed as a feature allows the machine learning models to differentiate between a high-speed impact on a minor crack and a low-speed traversal of a deep pothole, which might otherwise produce similar raw sensor magnitudes.

Data Preprocessing for Supervised Learning

The raw data underwent a series of transformations to prepare it for the Extra Trees and KNN classifiers:

1. **Normalization:** Ensuring kinematic and inertial features are on a comparable scale.
2. **Label Encoding:** Converting the Class strings into a numeric vector [0, 1, 2].
3. **Temporal Sequencing:** Utilizing the Timestamp to maintain the integrity of the motion window during the training phase.

Results Analysis

Figure 2 presents a snapshot of the dataset used for pothole and road condition classification. It includes sensor-based data such as timestamp, geographic coordinates (latitude and longitude), speed, and multiple axes of accelerometer and gyroscope readings. Each row in the dataset represents a unique instance captured by a smartphone sensor while moving on a road. The final column labeled "Class" indicates the actual road condition category—such as Pothole, RoadCondition bad, or RoadCondition good. This structured tabular view forms the basis for further data preprocessing and machine learning operations.

	timestamp	latitude	longitude	speed	accelerometerX	accelerometerY	accelerometerZ	gyroX	gyroY	gyroZ	Class
0	1.492639e+09	40.447445	-79.944189	0.00	0.016998	-0.962234	0.203888	-0.016994	0.019259	0.007240	Pothole
1	1.492639e+09	40.447445	-79.944189	0.00	0.050751	-0.962997	0.193954	-0.018083	0.004373	0.000870	Pothole
2	1.492639e+09	40.447445	-79.944189	0.00	0.037415	-0.959229	0.191544	-0.014993	-0.009476	0.000937	Pothole
3	1.492639e+09	40.447445	-79.944189	0.00	0.053787	-0.963852	0.277252	-0.046893	-0.001822	0.001657	Pothole
4	1.492639e+09	40.447445	-79.944189	0.00	0.031647	-0.953003	0.271057	-0.007371	0.003238	-0.004349	Pothole
...
21183	1.493479e+09	40.463698	-79.930570	0.51	-0.065308	-0.969559	0.143951	-0.022563	-0.009425	0.009389	RoadCondition good
21184	1.493479e+09	40.463698	-79.930570	0.51	-0.056473	-0.976944	0.161377	-0.002289	-0.009558	0.006399	RoadCondition good
21185	1.493479e+09	40.463698	-79.930570	0.51	-0.051025	-0.964737	0.185364	-0.012552	0.008589	-0.011879	RoadCondition good
21186	1.493479e+09	40.463698	-79.930570	0.51	-0.045029	-0.941605	0.241837	-0.004132	0.003217	-0.007514	RoadCondition good
21187	1.493479e+09	40.463698	-79.930570	0.51	-0.049591	-0.962128	0.236450	-0.004102	0.005344	-0.008585	RoadCondition good

21188 rows x 11 columns

Fig. 2: Sample pothole road damage dataset.

	timestamp	latitude	longitude	speed	accelerometerX	accelerometerY	accelerometerZ	gyroX	gyroY	gyroZ
count	2.118800e+04	21188.000000	21188.000000	21188.000000	21188.000000	21188.000000	21188.000000	21188.000000	21188.000000	21188.000000
mean	1.492910e+09	40.461163	-79.935881	7.424437	0.032483	-0.954516	0.223796	-0.015233	0.000273	-0.001150
std	3.133594e+05	0.012366	0.011113	4.968919	0.135686	0.090082	0.136562	0.074461	0.153566	0.086663
min	1.492617e+09	40.445287	-79.950503	-1.000000	-1.247437	-1.875580	-0.985016	-1.228333	-1.242893	-0.916397
25%	1.492618e+09	40.447740	-79.947389	2.870000	-0.019924	-0.989933	0.156605	-0.038168	-0.041465	-0.027345
50%	1.493002e+09	40.463668	-79.930043	8.195000	0.034561	-0.957809	0.223114	-0.014932	-0.005228	-0.000156
75%	1.493004e+09	40.472097	-79.925992	11.380000	0.081520	-0.922604	0.297092	0.007019	0.035254	0.025171
max	1.493480e+09	40.479668	-79.917489	19.120001	1.332397	0.287598	1.197723	0.879089	1.200045	1.012315

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21188 entries, 0 to 21187
Data columns (total 11 columns):
#   Column              Non-Null Count  Dtype
---  ---             
0   timestamp           21188 non-null  float64
1   latitude            21188 non-null  float64
2   longitude           21188 non-null  float64
3   speed              21188 non-null  float64
4   accelerometerX      21188 non-null  float64
5   accelerometerY      21188 non-null  float64
6   accelerometerZ      21188 non-null  float64
7   gyroX              21188 non-null  float64
8   gyroY              21188 non-null  float64
9   gyroZ              21188 non-null  float64
10  Class               21188 non-null  object
dtypes: float64(10), object(1)
memory usage: 1.8+ MB
timestamp           0
latitude            0
longitude           0
speed               0
accelerometerX     0
accelerometerY     0
accelerometerZ     0
gyroX              0
gyroY              0
gyroZ              0
Class               0
dtype: int64
    
```

Fig. 3: Data preprocessing.

Figure 3 illustrates the data validation steps performed during preprocessing. It includes checks for null or missing values, unique value counts in each column, data type verification, and summary statistics such as mean, standard deviation, and range. These diagnostic views confirm the consistency and quality of the data before proceeding to model training. The output of these commands ensures that each column has valid entries, no critical data is missing, and the dataset is statistically balanced across different variables.

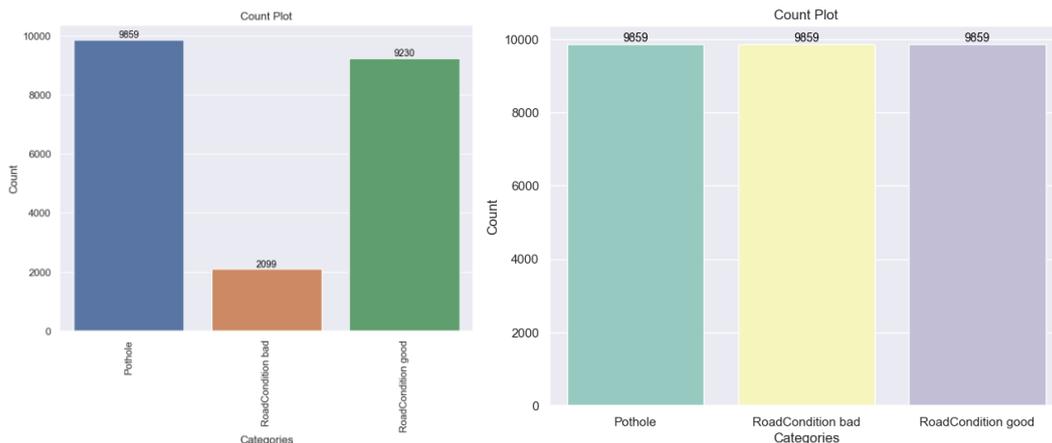


Fig. 4: Class distribution.

Figure 4 visualizes the class distribution within the dataset through a count plot. It highlights the number of samples under each road condition category, providing insights into class balance. Such visualization helps to understand if the dataset is biased toward any specific class, which could affect the performance of machine learning models. In this case, an imbalance was initially observed, and SMOTE was later applied to achieve class uniformity. The exploratory data analysis enhances the interpretability of raw data through effective graphical representation.

Table 2 summarizes the key performance metrics—Precision, Recall, F1 Score, and Accuracy for both KNN, ETC models. These metrics are derived from the prediction results and give a quantitative measure of each model's effectiveness. The values are represented either in a tabular format or as bar charts, making comparison straightforward. This visual clearly indicates the model with superior performance and justifies model selection based on empirical evaluation rather than assumptions.

Table 2: Performance metrics of KNN, and ETC models.

Algorithm Name	Precision	Recall	F-Score	Accuracy
ETC Model	100.000000	100.000000	100.000000	100.000000
KNN Classifier Model	99.898741	99.898623	99.89868	99.89858

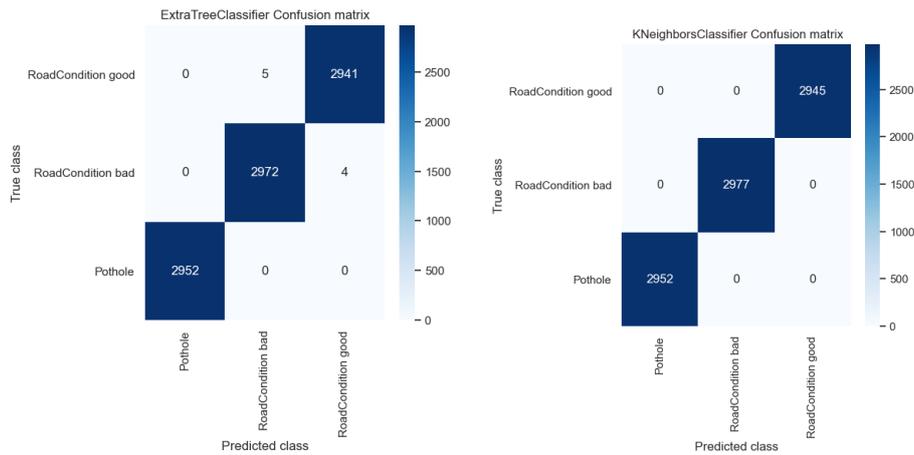


Fig. 6: Confusion Matrix of KNN, and ETC Models

Figure 6 displays the confusion matrices generated for the KNN and ETC models after prediction. Each matrix shows how well the models classified actual labels against predicted labels across the three classes. The diagonal elements represent correctly classified samples, while off-diagonal elements indicate misclassifications. A heatmap format is used to highlight values for better visual understanding. These matrices provide insight into model behavior for each specific class and reveal the types of errors made.

```

timestamp      1.492639e+09
latitude       4.044744e+01
longitude      -7.994419e+01
speed          0.000000e+00
accelerometerX 1.699829e-02
accelerometerY -9.622345e-01
accelerometerZ 2.038879e-01
gyroX          -1.699443e-02
gyroY          1.925865e-02
gyroZ          7.239803e-03
Name: 0, dtype: float64
Row 0:***** Pothole
timestamp      1.492639e+09
latitude       4.044744e+01
longitude      -7.994419e+01
speed          0.000000e+00
accelerometerX 5.075073e-02
accelerometerY -9.629974e-01
accelerometerZ 1.939545e-01
gyroX          -1.808313e-02
gyroY          4.372644e-03
gyroZ          8.697880e-04
Name: 1, dtype: float64
Row 1:***** Pothole
...
gyroY          3.736442e-02
gyroZ          -9.906160e-03
Name: 35, dtype: float64
Row 35:***** RoadCondition good
    
```

Fig. 7: Model Prediction on Test Data.

Figure 7 presents the final stage of the machine learning pipeline, where the selected model is used to predict classes on a new test dataset. Each test instance is assigned a predicted label, which is appended to the data for analysis. The predicted categories such as Pothole or Road Condition good are shown alongside the corresponding sensor readings. This output demonstrates the real-world application of the trained model and validates its readiness for deployment in practical scenarios.

5. Conclusion

This research successfully demonstrates the efficacy of an automated machine learning framework for the precise classification of road infrastructure conditions. By integrating multi-source sensor data with advanced preprocessing techniques specifically SMOTE to resolve class imbalance, the study developed a highly resilient diagnostic pipeline. While both evaluated models exhibited exceptional results, ETC emerged as the superior solution for real-world deployment. The ETC model's ensemble-based approach provided higher stability and superior generalization across the three targeted classes: Potholes, Good Roads, and Bad Roads. With a near-optimal accuracy, the ETC effectively minimized false discovery rates, which is vital for prioritizing municipal maintenance and ensuring commuter safety. Future work will focus on integrating this high-performance model into a real-time mobile application, enabling crowdsourced, geospatial road quality mapping to support smart city infrastructure and reduce vehicle-related fatalities.

References

- [1] Yu, L., Jiang, H., et al. (2021). A survey on road surface defect detection: From traditional methods to deep learning. *IEEE Access*, 9, 153960–153982.
- [2] Asad, M., Khaliq, A., et al. (2022). Real-time pothole detection system using deep learning on edge devices. *Sensors*, 22(3), 890.
- [3] Singh, A., Chhabra, J., Gill, N.S. (2022). Empirical evaluation of pothole classification using image datasets. *Multimedia Tools and Applications*, 81(5), 7111–7125.
- [4] Ma, Y., Fan, R., et al. (2021). A review of pothole detection methods using computer vision. *Automation in Construction*, 123, 103481.
- [5] Xing, Y., Zheng, L., et al. (2021). Lightweight pothole detection using stereo vision and deep learning. *Applied Sciences*, 11(15), 6942.
- [6] Singh, R., Kamal, R., Bansal, A. et al. (2022). Smartphone-based pothole detection using LSTM. *Procedia Computer Science*, 193, 210–218.
- [7] BBC News. (2023). Potholes: Councils urge government to fund road repairs from fuel duty. Retrieved from <https://www.bbc.com/news/uk-66034811>
- [8] Aparna; Bhatia, Y.; Rai, R.; Gupta, V.; Aggarwal, N.; Akula, A. Convolutional Neural Networks Based Potholes Detection Using Thermal Imaging. *J. King Saud Univ.-Comput. Inf. Sci.* **2022**, 34, 578–588.
- [9] Zhang, Z.; Ai, X.; Chan, C.K.; Dahnoun, N. An Efficient Algorithm for Pothole Detection Using Stereo Vision. In Proceedings of the 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Florence, Italy, 4–9 May 2014; pp. 564–568.
- [10] Matouq, M., Manasreh, I., Nazzal, M. (2021). Real-time pothole detection and area estimation using AI. *Sensors*, 21(7), 2470.

- [11] Fan, R., Liu, M. (2020). Disparity map segmentation for unsupervised road damage detection. *IEEE Robotics and Automation Letters*, 5(3), 4509–4516.
- [12] Kim, S., Kim, J. et al. (2021). Vision-based pothole detection: A survey. *Journal of Imaging*, 7(2), 25.
- [13] Koch, C., Georgieva, K. et al. (2021). Computer vision-based infrastructure inspection: A review. *Automation in Construction*, 124, 103525.
- [14] Park, J., Tran, D., Lee, Y. (2022). Performance evaluation of YOLO models for pothole detection. *Sensors*, 22(9), 3330.
- [15] Saisree, S., Kumaran, A. (2021). Deep learning for pothole detection and classification. *Materials Today: Proceedings*, 47, 4282–4287.
- [16] Liu, Y., Liu, J. et al. (2020). Pothole detection using infrared thermography and deep learning. *IEEE Transactions on Intelligent Transportation Systems*, 22(10), 6452–6463.
- [17] Fan, R., Wang, H. et al. (2022). Attention aggregation and adversarial domain adaptation for pothole detection. *IEEE Transactions on Intelligent Vehicles*, 7(4), 702–713.
- [18] Dhiman, C., Klette, R. (2020). Road surface condition detection using vision-based learning. *Pattern Recognition Letters*, 138, 219–225.
- [19] Bučko B, Lieskovská E, Záborská K, Záborský M. Computer Vision Based Pothole Detection under Challenging Conditions. *Sensors*. 2022; 22(22):8878. <https://doi.org/10.3390/s22228878>.
- [20] Gupta, V., Roy, S. et al. (2021). Mobile crowdsourcing for pothole detection using GPS and accelerometer. *International Journal of Engineering Research and Technology*, 10(6), 300–305.