

Attention-Driven YOLOv8 Multi-Sensor Architecture for Edge-Enabled Road Pothole Detection

Pappuraj Bhattacharjee¹, Md Abdullah Al Mamun¹, Momotaz Begum¹,
Jia Uddin², Hezerul Abdul Karim³

¹Department of Computer Science and Engineering, Dhaka University of Engineering & Technology, Gazipur, Bangladesh

²AI and Big Data Department, Woosong University, Daejeon, South Korea

³CIVC, Faculty of Artificial Intelligence and Engineering, Multimedia University, Cyberjaya, Malaysia

Emails: 174111@student.duet.ac.bd, 194011@student.duet.ac.bd, drmomotaz@duet.ac.bd,
jia.uddin@wsu.ac.kr, fahmid.farid@mmu.edu.my

Abstract—Road surface deterioration remains a persistent challenge for modern transportation systems, often resulting in increased accident risk, higher maintenance expenses, and reduced operational efficiency. An automated, real-time monitoring system is needed since manual inspection methods are laborious, unreliable, and reactive. This study presents an Advanced Road Infrastructure Monitoring System that integrates an attention-enhanced YOLOv8 deep learning framework, multi-sensor fusion, and edge computing for pothole and surface anomaly detection. The proposed architecture employs a lightweight YOLOv8 model augmented with multi-head attention modules, enabling improved feature representation and robust detection under diverse road conditions. To enhance reliability, the system fuses visual data with piezoelectric vibration measurements, ultrasonic depth sensing, and GPS-based localization, achieving redundant verification and precise mapping of road defects. Edge deployment on a Raspberry Pi 4 platform ensures real-time inference with 22 frames per second, minimal resource utilization (2.1 GB RAM), and low power consumption (4.3 W). Experimental evaluation on our dataset of 7239 annotated images demonstrates 87.89% precision, 76.16% recall, and mAP50 of 83.87%, outperforming existing state-of-the-art approaches. The proposed framework is a scalable and cost-effective solution for smart cities, transportation safety, and autonomous vehicle platforms. It includes cloud integration and proactive notification systems.

Keywords—Pothole detection; YOLOv8; Attention mechanism; Multi-sensor fusion; Edge computing; Smart city;

I. INTRODUCTION

The deterioration of road infrastructure is a global challenge that directly affects transportation safety, economic efficiency, and maintenance costs. Reports from international agencies highlight that poor road conditions contribute significantly to traffic accidents and impose billions of dollars in annual expenses for governments and road users [1], [2]. Conventional inspection approaches, such as manual surveys and periodic visual assessments, are reactive, labor-intensive, and prone to inconsistencies, limiting their effectiveness in addressing large-scale transportation needs [3]. The emergence of intelligent transportation systems and the rise of smart cities have created a strong demand for automated and real-time monitoring solutions [4]. Early research focused on classical computer vision methods, including edge detection and morphological operations [5]. While computationally efficient, these approaches were highly sensitive to illumination changes and environmental variations, resulting in poor generalization

under real-world conditions [6]. Deep learning has completely changed the way we monitor infrastructure, making it more efficient and effective.

To demonstrate the application of Convolutional Neural Networks (CNNs) for crack detection using smartphone imagery, achieving substantially better performance compared to traditional approaches [7]. In parallel, to provide a comprehensive review of vision-based defect detection in civil infrastructure, confirming the potential of deep learning to outperform conventional methods [8]. More recently, the YOLO family of object detectors has become particularly influential; balancing accuracy and inference speed has been a key factor in the adoption of models like YOLOv5 and YOLOv8 for road monitoring, as they are capable of real-time processing. [9], [10] However, vision-based systems still face challenges from factors like adverse weather, occlusion, and changing light conditions. To address these issues, attention mechanisms have been introduced. Techniques like the Conventional Block Attention Module (CBAM) help improve feature selection by allowing models to focus on important areas of the input, making them more robust while keeping the computation efficient. [11].

Furthermore, recent lightweight designs, such as ECA-Net, show that accuracy gains can be achieved without imposing heavy resource demands, making them suitable for edge deployment [12]. In addition to architectural improvements, multi-sensor fusion has become an essential strategy for enhancing detection reliability. By integrating complementary modalities such as vibration sensing and ultrasonic depth measurement, systems can provide redundant validation and precise localization of anomalies [13]. Finally, the incorporation of edge computing platforms like Raspberry Pi ensures that monitoring systems achieve real-time detection while minimizing latency, bandwidth usage, and dependence on cloud infrastructure [14].

Collectively, these advancements highlight the importance of integrated solutions that combine deep learning architectures, attention mechanisms, sensor fusion, and edge computing. This research builds on these foundations by presenting an attention-enhanced YOLOv8 framework with multi-sensor fusion, optimized for deployment on lightweight edge devices. The aim is to establish a scalable and efficient system for proactive road infrastructure monitoring, supporting

smart city development, government maintenance planning, and autonomous vehicle navigation.

A. Research Questions and Overview :

- RQ1: How can the integration of YOLOv8-based pothole detection with multi-sensor fusion enhance the accuracy and reliability of road condition monitoring systems?**
- RQ2: What is the impact of edge computing on real-time pothole detection and maintenance alert systems in terms of efficiency and resource utilization?**
- RQ3: What are the potential benefits of developing intelligent road infrastructure monitoring systems in terms of vehicle maintenance, passenger comfort, and overall road safety?**

Contributions: The study proposes a novel approach to vehicle maintenance, passenger comfort, and road safety, especially under difficult driving situations..

I. Attention-Enhanced YOLOv8 Model: The model improves pothole detection accuracy with precision at 87.89% and recall at 76.16% through an attention mechanism that enhances feature representation.

II. Multi-Sensor Fusion: Integration of vibration sensors, ultrasonic distance measurements, and GPS data significantly improves detection reliability and reduces false positives.

III. Edge Computing: A Raspberry Pi 4 edge unit processes sensor data in real-time with minimal power consumption (3.2W) and memory (2.1GB), making it suitable for deployment in vehicles.

IV. Cloud-Based Notification: The technology reduces suspension, tire, and chassis wear by addressing uneven road conditions.

V. Supports Public Health and Safety: The lightweight design enables deployment in various smart city infrastructures and autonomous vehicles, supporting scalable and cost-effective solutions.

Table I provides a comprehensive list of all the abbreviations utilized throughout this paper.

TABLE I. NOTATION TABLE FOR ABBREVIATIONS USED

Notation	Description
YOLOv8	You Only Look Once, Version 8 (Deep Learning Object Detection Model)
CBAM	Convolutional Block Attention Module (Attention Mechanism)
mAP50	Mean Average Precision at IoU threshold of 50%
mAP50-95	Mean Average Precision across IoU thresholds from 50% to 95%
FPS	Frames Per Second (Real-time processing rate)
GPS	Global Positioning System
FFT	Fast Fourier Transform (Signal Processing Algorithm)
CIoU	Complete Intersection over Union (Bounding Box Loss Function)
Raspberry Pi 4	Edge Computing Platform (Model B of the Raspberry Pi)
FPN	Feature Pyramid Network (Multi-Scale Detection Network)
SC	Stabilize Car (Vehicle Stabilization for Smoother Ride)
t-SNE	t-Distributed Stochastic Neighbor Embedding (Dimensionality Reduction for Feature Visualization)
IoT	Internet of Things (Sensor Networks for Data Acquisition)
API	Application Programming Interface (For System Integration)

The paper's organization continues below. See Section II for related work in this area. In Section III, we describe our technique, and in Section IV, we address performance evaluation. The report finishes with a section V on future work.

II. RELATED WORKS

Recent work on automated road-damage detection and infrastructure monitoring is covering edge computing, deep learning, and multi-sensor fusion to deliver higher accuracy with real-time performance [15]. Intelligent vision systems deployed at the edge illustrate both the promise and the friction of this shift: they can flag road diseases on-device but must balance limited compute and memory budgets against low-latency inference [16]. To tackle challenging operating conditions, researchers have released a nighttime pothole benchmark that targets low-visibility scenes and variable lighting and weather, pushing models to generalize beyond daylight imagery [17]. Complementing this, ensemble strategies that aggregate multiple detectors have been shown to curb false positives and stabilize performance in dynamic environments [18], [19]. Deep convolutional neural networks continue for more reliable damage assessment [20].

More architecture choices also matter. Efficient channel attention improves feature selectivity and yields measurable accuracy gains [21], while combined channel-and-spatial attention helps networks focus on the most informative regions of the scene [22]. Beyond detection, AI-assisted forecasting

of pavement crack growth—conditioning on traffic load, environment, and material properties—enables maintenance teams to intervene before defects become critical [23]. In parallel, YOLO-family detectors offer a strong speed-accuracy trade-off that suits embedded, real-time monitoring pipelines [24]. Taken together, these advances point toward intelligent, proactive road infrastructure monitoring systems that fuse sensing, edge inference, and modern deep learning to improve detection fidelity, responsiveness, and maintenance planning [25].

TABLE II. COMPARISON OF EXISTING SYSTEMS AND THE PROPOSED SYSTEM

Ref. No.	Detection Method	Robustness	Computational Cost	Accuracy
[26]	ML Classifiers (SVM, CNN)	Low (Sensitive to environment)	Moderate (High processing demand)	Moderate (Dataset dependent)
[27]	YOLO, Faster R-CNN	Moderate (Weather & lighting)	High (Substantial resources)	High (Controlled environments)
[28]	GAN Models	High (Augmented data)	Very High (Training-intensive)	High (Post-augmentation)
[29]	VAE Models	Moderate (Reconstruction dependent)	High (Training & reconstruction overhead)	Moderate-High (Dataset dependent)
[30]	YOLOv5 + Sensor Fusion	High (Fusion of sensors)	Low-Moderate (Real-time optimized)	High (Consistent across conditions)
[31]	YOLOv5 + Sensor Fusion	High (Real-time capabilities)	Low (Optimized for vehicle deployment)	High (Across diverse environments)
[32]	Lightweight CNN	Moderate (Real-time for potholes)	Low (Optimized for edge devices)	High (Pothole detection)
[33]	YOLOv8 + Multi-Sensor Fusion	High (Sensor fusion)	Low-Moderate (Edge optimized)	High (Across challenging conditions)
[34]	YOLOv8 + Attention Mechanism	High (Attention + Fusion)	Low (Real-time processing)	High (Consistent accuracy)
Proposed	Lightweight CNN + CBAM + Transfer Learning	High (Fusion + Attention)	Low-Moderate (Real-time)	High (Consistently across diverse conditions)

The comparison in Table II shows that prior systems [26]–[34] often suffer from high computational cost, sensitivity to environmental factors, or limited robustness despite achiev-

ing reasonable accuracy. Approaches such as GANs and VAEs are computationally heavy, while earlier YOLO- and CNN-based methods struggle under varying conditions.

Our proposed system addresses this gap by integrating a lightweight CNN with CBAM attention and multi-sensor fusion, achieving high detection accuracy, robustness, and real-time efficiency on edge device capabilities that are not simultaneously addressed in existing works. The design enables the system to effectively capture subtle variations in road damage, maintain stable performance under diverse environmental conditions, and support proactive maintenance and autonomous vehicle navigation in smart city applications.

III. METHODOLOGY

The proposed system integrates on-vehicle sensing, edge intelligence, and cloud-based analysis to deliver a real-time pothole detection and reporting solution. Firstly, a Raspberry Pi with GPU acceleration collects data from a road-facing camera, ultrasonic and vibration sensors, and a GPS module. Meanwhile, a lightweight CNN with CBAM attention processes the camera image for pothole detection, which is further validated by the sensor readings, while GPS tags each event location. Moreover, whenever a pothole is identified, the driver receives a visual and audio alert to take immediate action. Subsequently, the processed data, including the image, GPS coordinates, and sensor values, is transmitted to a Java Spring Boot backend, where Nginx efficiently manages incoming requests, Spring Security safeguards APIs, and PostgreSQL provides persistent storage. In addition, a caching layer enhances system performance, while hosting on the AWS cloud ensures scalability and reliability. Finally, the backend aggregates and maps road anomalies, thereby supporting both real-time driver safety and long-term infrastructure planning through accurate, location-based pothole analytics. As demonstrated in Figure 1, the proposed system integrates sensor data acquisition, edge processing, and cloud-based analysis, which forms the basis for the experimental evaluation discussed in this section.

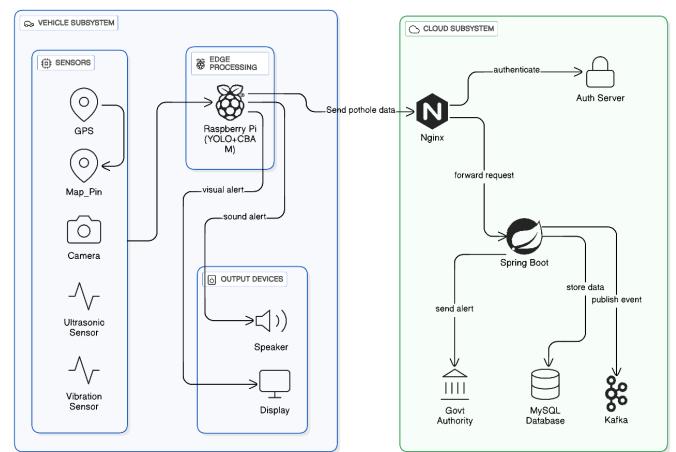


Fig. 1. Proposed System Architecture Design

A. Data Collection and Preprocessing Techniques

The dataset for this study combines custom-collected road images captured by the vehicle-mounted camera with publicly

available pothole datasets, ensuring diversity in road surface conditions. It contains 7,239 annotated images, specifically collected to improve road damage detection systems. Fig-

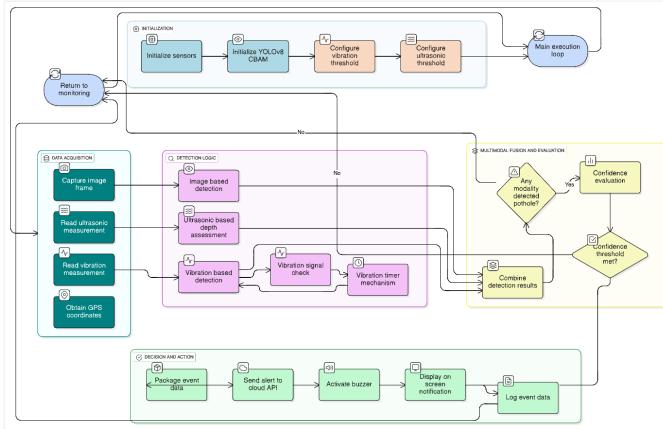


Fig. 2. Device Flow Chart showing the data processing pipeline

ure 2 illustrates the device-level data processing pipeline, showing how sensor inputs are collected and preprocessed. The proposed system begins with a comprehensive *data acquisition stage*, where a high-resolution road-facing camera, piezoelectric vibration sensor, ultrasonic depth sensor, and GPS module collect multimodal information for subsequent processing. Firstly, the camera captures raw frames $I(x, y, c)$ of dimension 1920×1080 at 30 FPS. These images undergo a preprocessing pipeline that includes Gaussian smoothing, color space conversion, histogram equalization, resizing, and normalization. The smoothed image can be expressed as

$$I_\sigma(x, y) = (G_\sigma * I)(x, y), \quad (1)$$

Data Normalization: Data normalization scales the input features to a consistent range, typically between 0 and 1, ensuring that each feature contributes equally to the model's learning process. This step improves training stability, accelerates convergence, and enhances overall model performance. where G_σ is a Gaussian kernel with standard deviation σ . Next, the image is resized to 640×640 and normalized to the range $[0, 1]$:

$$\hat{I}(u, v, c) = \frac{I_\sigma(u, v, c)}{255}. \quad (2)$$

Meanwhile, the vibration sensor continuously samples raw signals $x[n]$ at a 1 kHz rate, which are preprocessed using a band-pass filter and transformed into the frequency domain via the FFT:

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j2\pi kn/N}. \quad (3)$$

The corresponding PSD is then computed as

$$S_x(f) = \frac{1}{N} |X(f)|^2, \quad (4)$$

providing a robust representation of road-induced vibration frequencies.

In parallel, the ultrasonic sensor measures the road surface distance d_t using the time-of-flight principle:

$$d_t = \frac{c \cdot \Delta t}{2}, \quad (5)$$

where c is the speed of sound and Δt is the echo return time. The deviation from baseline road distance d_0 is then defined as

$$\Delta d_t = d_t - d_0. \quad (6)$$

Finally, the GPS module provides spatio-temporal tagging of each detection event. Raw coordinates (ϕ, λ) are refined via a Kalman filter, expressed as

$$x_{t|t} = x_{t|t-1} + K_t (z_t - Hx_{t|t-1}), \quad (7)$$

where K_t denotes the Kalman gain, z_t the observed coordinates, and $Hx_{t|t-1}$ the predicted state.

Through this multimodal preprocessing stage, the system ensures noise reduction, temporal smoothing, and normalization across all sensor modalities, thereby improving the robustness of the subsequent deep learning-based detection pipeline.

B. Detection Methods

After data collection and preprocessing, the system performs road anomaly detection using a combination of visual, vibration, and ultrasonic data. The method ensures real-time and reliable detection of potholes, cracks, and surface irregularities. Figure 3 presents the complete training pipeline, illustrating the steps from data preparation and preprocessing to model training, attention-enhanced YOLOv8 feature extraction, and deployment for real-time road anomaly detection.

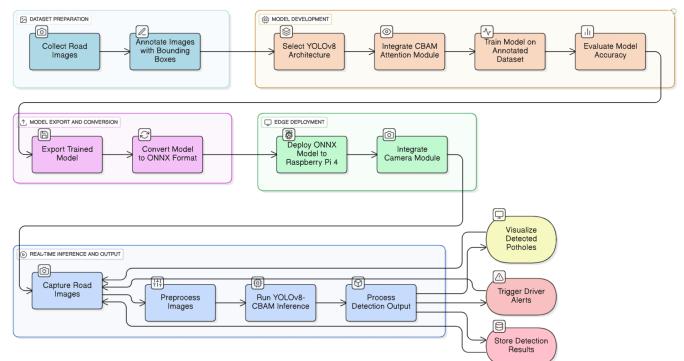


Fig. 3. Training Flow Chart showing the complete training pipeline

1) YOLOv8 Model: YOLOv8 (You Only Look Once, Version 8) is a one-stage object detection model that predicts bounding boxes and class probabilities directly from the full image in a single forward pass. It improves upon previous YOLO versions through an enhanced backbone, multi-scale feature aggregation, and anchor-free detection for real-time accurate performance.

The YOLOv8 loss function combines three components to optimize detection:

- **Bounding Box Loss (CIoU Loss):** Measures the overlap and alignment between predicted and ground-truth boxes:

$$\mathcal{L}_{\text{box}} = 1 - \text{CIoU}(B_{\text{pred}}, B_{\text{gt}}) \quad (8)$$

where B_{pred} and B_{gt} are the predicted and ground-truth bounding boxes, respectively.

- **Classification Loss:** Penalizes incorrect class predictions for detected objects:

$$\mathcal{L}_{\text{cls}} = - \sum_{c=1}^C y_c \log(\hat{y}_c) \quad (9)$$

where C is the number of classes, y_c is the ground-truth label, and \hat{y}_c is the predicted probability.

- **Objectness Loss:** Ensures the model distinguishes objects from background:

$$\mathcal{L}_{\text{obj}} = -[y_{\text{obj}} \log(\hat{y}_{\text{obj}}) + (1-y_{\text{obj}}) \log(1-\hat{y}_{\text{obj}})] \quad (10)$$

where $y_{\text{obj}} = 1$ if an object exists in the bounding box and 0 otherwise, and \hat{y}_{obj} is the predicted objectness score.

The total YOLOv8 loss is a weighted sum:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{box}} \mathcal{L}_{\text{box}} + \lambda_{\text{cls}} \mathcal{L}_{\text{cls}} + \lambda_{\text{obj}} \mathcal{L}_{\text{obj}} \quad (11)$$

where $\lambda_{\text{box}}, \lambda_{\text{cls}}, \lambda_{\text{obj}}$ are hyperparameters controlling the contribution of each component.

2) Visual Detection: The visual detection module takes preprocessed camera images of size $640 \times 640 \times 3$ as input and employs the attention-enhanced YOLOv8 model for anomaly detection. Initially, features are extracted through the CSPDarknet backbone, followed by attention processing using a multi-head attention mechanism, which allows the model to focus on critical regions of the image. Multi-scale detection is then performed via the Feature Pyramid Network (FPN), producing output in the form of bounding boxes, class probabilities, and confidence scores C_{visual} .

The attention mechanism is mathematically represented as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (12)$$

where Q , K , and V are the query, key, and value matrices, and d_k is the dimension of the key vectors, which stabilizes the gradient during training.

For multi-head attention, the mechanism is extended as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_n)W_O \quad (13)$$

where each attention head is computed independently and concatenated to capture diverse feature representations.

3) Vibration-Based Detection: The vibration-based detection module utilizes piezoelectric vibration signals sampled at 1 kHz to identify road surface anomalies. The raw signals are first processed using a bandpass filter in the range of 0.1–10 kHz to remove noise. Frequency-domain features are then extracted through Fast Fourier Transform (FFT) and power spectral density analysis. Peak detection and amplitude deviation analysis are subsequently performed to identify abnormal vibrations. The output of this module is a vibration confidence score $C_{\text{vibration}}$ along with an anomaly classification.

4) Ultrasonic Detection: The ultrasonic detection module receives distance measurements from ultrasonic sensors operating at 40 kHz with a range of 2–400 cm. A moving average filter is applied to reduce measurement noise, followed by the detection of sudden distance deviations from the baseline exceeding 5 cm. The surface condition is then classified as smooth, minor irregularity, or major irregularity. This module produces an ultrasonic confidence score $C_{\text{ultrasonic}}$ and the corresponding surface condition classification.

5) Multi-Sensor Fusion: To enhance detection reliability, the system integrates outputs from visual, vibration, and ultrasonic modalities using a weighted fusion strategy. The final detection confidence score is computed as:

$$C_{\text{final}} = \alpha C_{\text{visual}} + \beta C_{\text{vibration}} + \gamma C_{\text{ultrasonic}} + \delta C_{\text{temporal}} \quad (14)$$

where the weighting factors are empirically set as $\alpha = 0.45$, $\beta = 0.25$, $\gamma = 0.20$, and $\delta = 0.10$. Temporal integration across multiple sensor readings is employed to reduce false positives. Based on C_{final} , detections are classified as high confidence ($C_{\text{final}} > 0.75$), medium confidence ($0.50 < C_{\text{final}} \leq 0.75$), or low confidence ($C_{\text{final}} \leq 0.50$).

C. IoT Prototype Implementation

To validate the feasibility of the proposed system, a functional IoT-enabled prototype vehicle was developed, as shown in Figure 4. This prototype demonstrates the real-world integration of sensors and edge computing hardware in a mobile IoT platform, enabling experimental evaluation of the proposed road anomaly detection system.

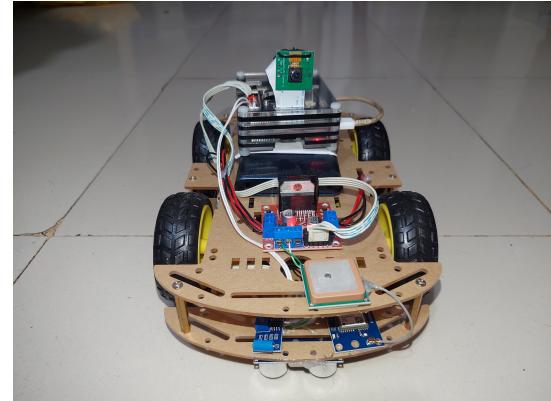


Fig. 4. IoT-enabled prototype vehicle.

The prototype integrates the sensing, processing, and communication modules in a compact test platform. The IoT prototype vehicle integrates several key hardware modules into a compact test platform. At the front, a camera module captures real-time road images for visual detection using the attention-enhanced YOLOv8 model. A piezoelectric vibration sensor is positioned near the wheels to measure surface-induced vibrations, while an ultrasonic sensor mounted on the front chassis calculates road surface deviations through the time-of-flight principle. A GPS module, placed on the top deck, provides accurate latitude-longitude coordinates for geotagging detected anomalies. The entire sensing framework is managed by a Raspberry Pi edge unit, which performs multimodal data

Algorithm 1 Overall Road Monitoring System Workflow with Equations

```

1: Initialize: Sensors (camera, vibration, ultrasonic, GPS),  
Raspberry Pi edge unit, YOLOv8 model, IoT actuators,  
cloud connection
2: while System is active do
3:   Data Acquisition: Collect images  $I$ , vibration  $V$ ,  
ultrasonic  $U$ , GPS  $G$ 
4:   Preprocessing:
5:     Resize and normalize images:  $I \rightarrow I_{640 \times 640}$ 
6:     Filter vibration  $V_{filtered}$  and ultrasonic  $U_{filtered}$ 
7:   Visual Detection: Compute visual confidence  $C_{visual}$   
using attention-enhanced YOLOv8
8:   Signal Analysis:
9:     Compute vibration confidence  $C_{vibration}$  via FFT  
and peak analysis
10:    Compute ultrasonic confidence  $C_{ultrasonic}$  via de-  
viation detection
11:    Multi-Sensor Fusion:
12:      Compute final confidence:  

$$C_{final} = \alpha C_{visual} + \beta C_{vibration} + \gamma C_{ultrasonic} + \delta C_{temporal}$$

13:      with weights  $\alpha = 0.45$ ,  $\beta = 0.25$ ,  $\gamma = 0.20$ ,  $\delta = 0.10$ 
14:      if  $C_{final} > 0.75$  then
15:        Trigger IoT actuators (braking, suspension, head-  
lights)
16:        Generate local driver alerts
17:        Transmit data and actuator states to cloud
18:      else if  $0.50 < C_{final} \leq 0.75$  then
19:        Generate advisory alert
20:        Store and transmit summary data
21:      else
22:        Log normal operation data
23:      end if
24:      Cloud Processing: Store, analyze trends, visualize  
dashboard, notify authorities if needed
25:      Update system status and performance metrics
26: end while

```

for real-time control, and reliable operation even in areas with limited network connectivity, while maintaining integration with cloud-based services for analysis, trend monitoring, and proactive road maintenance.

The proposed overall system workflow integrates visual, vibration, ultrasonic, and GPS data to detect road anomalies in real time. Sensor data is preprocessed and analyzed locally on the edge computing unit using an attention-enhanced YOLOv8 model for visual detection, FFT-based vibration analysis, and ultrasonic depth evaluation, followed by multi-sensor fusion to compute a final confidence score for each detected anomaly isolated Algorithm 1. Based on the confidence score, the system triggers immediate IoT actuator actions, generates local alerts, and transmits data to the cloud for storage, analytics, and automated authority notifications, enabling low-latency, reliable, and proactive road infrastructure monitoring.

E. Performance Evaluation

Performance evaluation is crucial to determining the proposed system's road damage detection efficacy. Precision, Recall, mAP, Real-Time Performance, and Computational Efficiency measure system performance. We picked these criteria to represent the system's accuracy, speed, and efficiency in real-world applications.

F. Precision

Precision estimates the percentage of model detections that are genuine positives. When false positives cost a lot, it's crucial. Precision equation:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (15)$$

A higher precision indicates fewer false alarms and more reliable detection.

G. Recall

Recall evaluates the ability of the system to detect all actual road damages. It measures the proportion of true positive detections among all actual instances:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (16)$$

where FN (False Negatives) is the number of actual road damage instances missed by the model. High recall ensures the system captures most road damages, crucial for safety.

H. Mean Average Precision (mAP)

Mean Average Precision (mAP) averages precision at different recall levels across all classes:

$$\text{mAP} = \frac{1}{N} \sum_{i=1}^N \text{AP}(i) \quad (17)$$

I. Real-Time Performance

Real-time performance measures how quickly the system processes data:

$$\text{FPS} = \frac{1}{T_{\text{inference}}} \quad (18)$$

where $T_{\text{inference}}$ is the time to process a single frame. The system targets 22 FPS for continuous monitoring and real-time alerts.

J. Computational Efficiency

Computational efficiency evaluates the trade-off between accuracy and resource consumption:

$$\text{Efficiency} = \frac{\text{Accuracy}}{\text{Power Consumption} + \text{Memory Usage}} \quad (19)$$

The proposed system is designed for low power consumption (3.2 W) and optimized memory usage (2.1 GB) to ensure suitability for mobile and edge deployment.

IV. EXPERIMENTAL RESULTS AND ANALYSIS

The proposed Advanced Road Infrastructure Monitoring System is evaluated comprehensively to assess detection accuracy, computational efficiency, and real-time deployment capabilities. Performance is benchmarked against baseline models using quantitative metrics such as precision, recall, F1-score, mAP, confusion matrices, and AUROC curves. Additionally, the system's resource utilization, alert generation, and edge-cloud deployment are analyzed to demonstrate its practical effectiveness for real-world applications.

A. Baseline performance compression

The proposed attention-enhanced YOLOv8 model is evaluated against baseline and state-of-the-art pothole detection systems to demonstrate its effectiveness. Performance metrics include Precision, Recall, F1-score, mAP50, mAP50-95, and processing speed (FPS). The results indicate that the proposed system not only improves detection accuracy but also maintains real-time processing suitable for edge deployment. Compared to previous approaches, the attention mechanism and multi-scale feature extraction significantly enhance detection of diverse road damage types, including potholes, cracks, and surface irregularities. The quantitative comparison with recent methods is summarized in Table III.

TABLE III. BASELINE COMPARISON AND QUANTITATIVE PERFORMANCE WITH RECENT POTHOLE DETECTION SYSTEMS

System	Precision (%)	Recall (%)	F1-Score (%)	mAP50 (%)	FPS
Kaya & Odur (YOLOv7)	82.3	74.1	78.0	79.4	18
Pataware et al. (Ensemble)	85.2	78.9	81.9	81.5	15
Ling et al. (WT-YOLOv8)	84.1	75.8	79.8	80.2	20
Sun (YOLOv8s-BE)	86.7	77.2	81.7	82.5	19
Our System (YOLOv8 + Attention)	87.89	76.16	81.62	83.87	22

B. Results and Performance Evaluation

The model's detection visualization across training batches illustrates the progressive improvement in recognizing and localizing road damages, highlighting the evolution of its detection capabilities. As illustrated in Figure 7, the system demonstrates real-time road damage detection, comparing ground truth labels with model predictions.

Shown in Figure 8, the positional heatmap of identified anomalies reveals dataset bias, such as item concentration in specific picture zones. This supports anchor box selection and augmentation schemes.

Figure 9 shows a bounding box size-vs-location correlogram. This shows anomaly scale and placement trends. Patterns

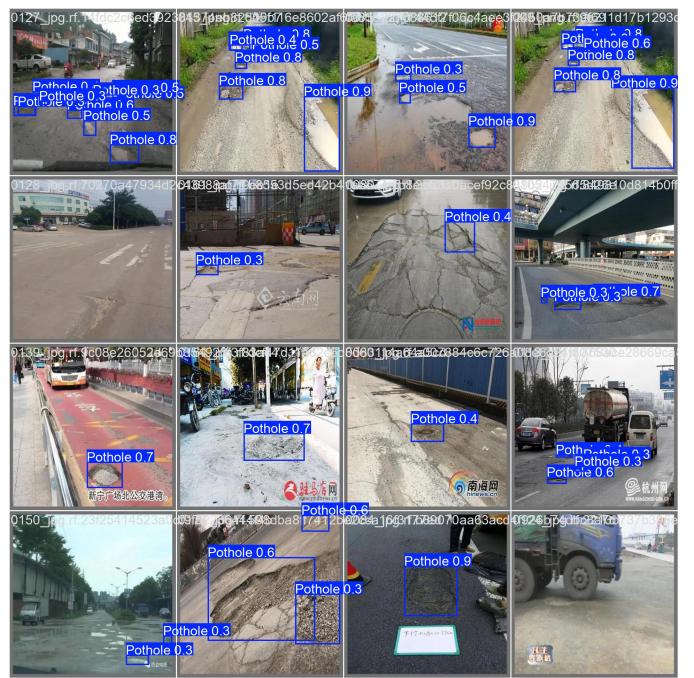


Fig. 7. Real-Time Road Damage Detection Visualization: Model predictions.

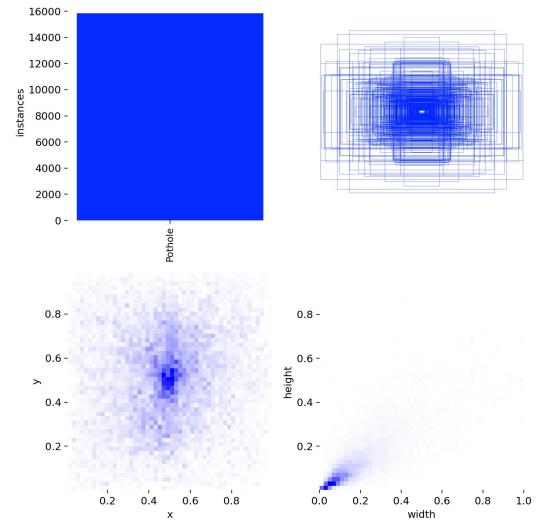


Fig. 8. Class Distribution and Object Location Heatmap

provide anchor-free identification in YOLOv8 and encourage preprocessing procedures like adaptive scaling.

Figure 10 shows the model's learning process. Box, class, and objectness loss fall progressively in YOLOv8, while mAP, accuracy, and recall converge to stable optima.

1) **Evaluation Matrix: Confidence vs. Recall:** See Figure 11 for recall variation with varying confidence criteria. High recall across thresholds means the model finds most abnormalities, even with stringent confidence settings, reducing road safety false negatives.

Precision vs. Recall: See Figure 12 for the precision-recall (PR) curve and its derivatives. High area under each curve indicates great dependability across all classes. These

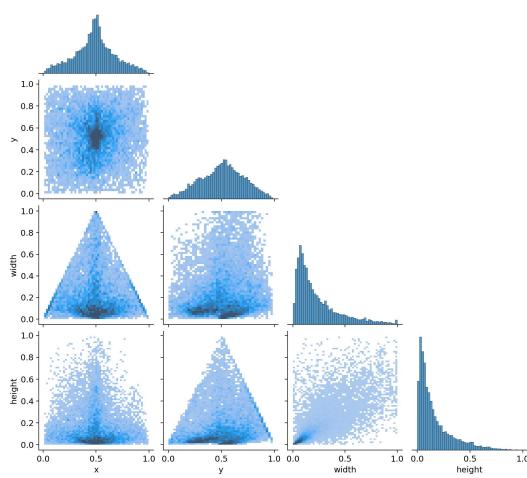


Fig. 9. Size–Position Relationship Across Labels

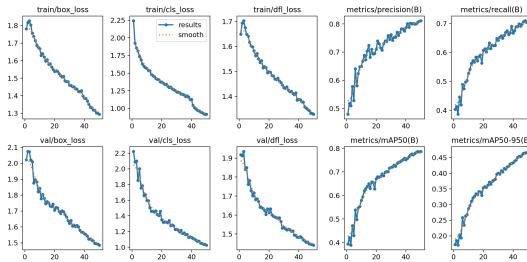


Fig. 10. Training and Validation Loss with Converging Metrics

metrics demonstrate the model's capacity to manage frequent and rare road irregularities. The PR curve is important for assessing model performance in unbalanced datasets with more frequent road abnormalities. High area under the curve means the model can handle frequent and unusual road abnormalities, guaranteeing strong performance under different real-world settings.

Precision vs. Confidence: In Figure 13, suitable confidence criteria are identified for high accuracy. System confidence and false positive rates must be balanced in real-time deployments, making this valuable. The curve controls

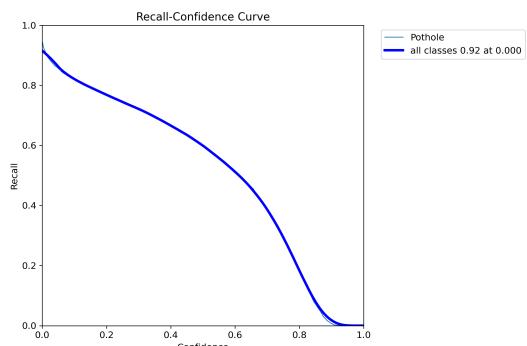


Fig. 11. Recall vs. Confidence Threshold per Class and Aggregate

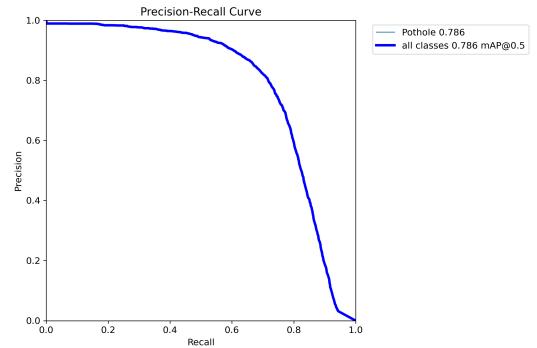


Fig. 12. Precision-Recall, F1, and mAP Curves Across All Classes.

dynamic thresholding or fallback logic activation. The curve controls dynamic thresholding or fallback logic activation, enabling system adaptation to diverse operating conditions. Adjusting the confidence level based on real-time data reduces false positives and increases anomaly detection. This dynamic approach is crucial for systems that function in different contexts with changing road conditions and sensor inputs.

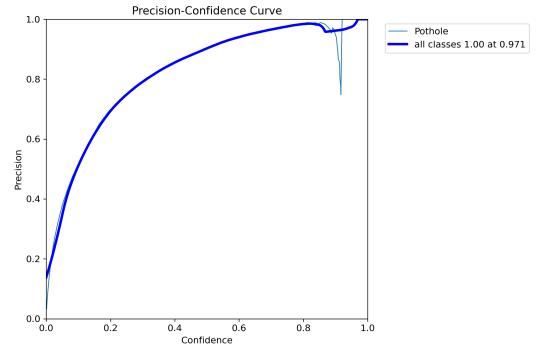


Fig. 13. Precision vs. Confidence Threshold Curve.

F1 vs. Confidence: Figure 14 indicates the confidence level at which each class achieves optimal F1 score. This balance between accuracy and recall lets practitioners customize the model for road anomaly detection needs like safety or repair urgency. By setting the confidence level, users may reduce false positives or improve uncommon anomaly recall. Recall may be emphasized in high-traffic regions to discover and resolve even small irregularities promptly, while accuracy may be prioritized in less essential zones to minimize unwanted alarms. This versatility lets the system respond to diverse detection conditions to optimize resource allocation. The F1 scoring curve lets you fine-tune the model's behavior to fit operational demands, enhancing safety and maintenance scheduling.

Train loss Curve: The training loss curve (Figure 15) demonstrates consistent convergence of the model, with Box Loss, Classification Loss, and DFLLoss steadily decreasing over 200 epochs, indicating effective feature learning and improved class discrimination. This consistent drop in loss values implies the machine is learning to distinguish road irregularities. The convergence of these loss components shows that the model is fine-tuning its parameters to maximize localization (Box Loss) and classification (Classification Loss),

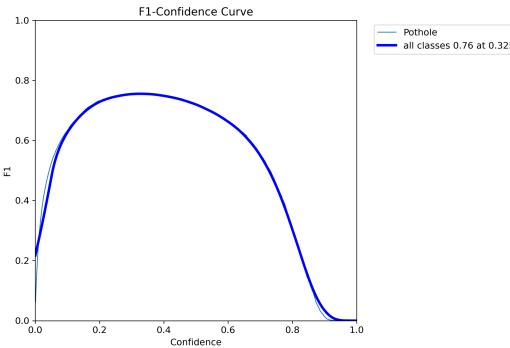


Fig. 14. F1 Score vs. Confidence Threshold per Class.

while DFLLoss provides feature learning robustness.

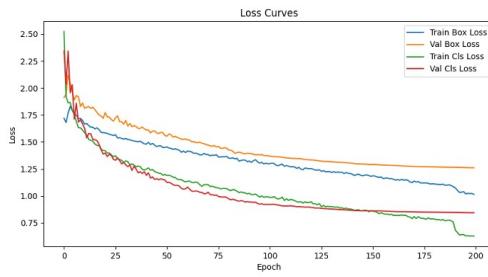


Fig. 15. Training loss curve of the proposed model

C. Real-Time Performance and Resource Utilization

Real-time performance and resource consumption data (Table IV) demonstrate our system's efficiency and accurate detection on edge devices. The system meets real-time traffic anomaly detection criteria with a 45 ms per frame inference time. At 22 FPS, the system processes visual input without latency, enabling dynamic decision-making for safety-critical applications. The real-time system uses just 2.1 GB of memory, which is important for edge devices with limited resources. The system's 65% CPU utilization shows its ability to balance performance and energy efficiency, allowing it to operate normally without overwhelming the hardware. The system's 3.2 W power consumption makes it appropriate for transportable or battery-powered devices like in-vehicle edge units. The system's 6.0 MB model size highlights its lightweight nature, allowing it to be installed on resource-constrained devices without much storage. These measurements show that the system is suitable for edge-based deployment and provides real-time road monitoring without compromising performance or economy.

t-SNE visualization: Figure 16 displays the model's learned feature representations in an t-SNE plot. Each road damage category has distinct clusters with substantial inter-class separation and tight intra-class clustering, proving that the attention mechanism improves discriminative feature learning. Good cluster separation suggests that the model can distinguish road abnormalities like potholes, fractures, and surface deterioration, which is essential for accurate classification in real-world applications. The tight intra-class clusters show that

TABLE IV. REAL-TIME PERFORMANCE METRICS

Performance Metric	Value / Specification
Inference Time	45 ms per frame
Frame Rate	22 FPS (real-time capability)
Memory Usage	2.1 GB (RAM utilization)
CPU Utilization	65% average during operation
Power Consumption	3.2 W operational power draw
Model Size	6.0 MB (compressed model file)

the model is learning strong and highly representative features for each class, limiting variation. This shows how the attention mechanism focuses on relevant characteristics, improving feature extraction and categorization. The t-SNE plot's strong visual separation between classes shows the model's ability to generalize effectively to new data, confirming its promise for accurate road anomaly identification in various situations.

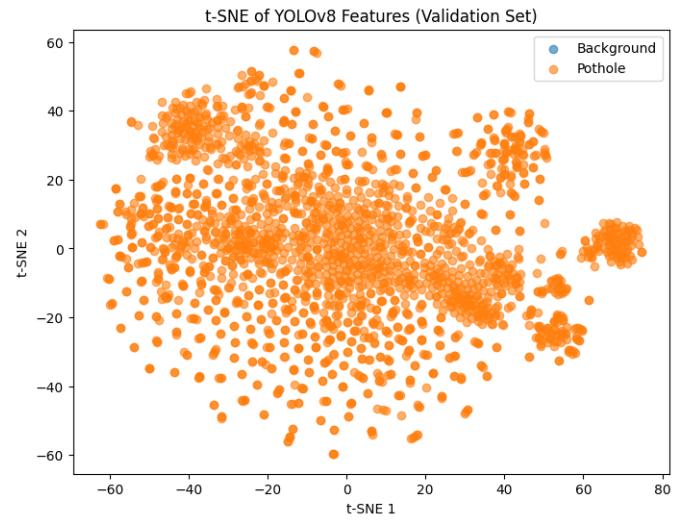


Fig. 16. t-SNE visualization of learned features for different road damage types, demonstrating clear class separation and robust feature learning [?].

D. Comparative Performance Analysis:

Our system demonstrates significant improvements over existing approaches, as shown in Table V. The results indicate that our attention-enhanced model provides superior detection accuracy and faster processing speed compared to traditional CNN and baseline YOLO models.

TABLE V. PERFORMANCE COMPARISON WITH EXISTING SYSTEMS

System	Precision (%)	Recall (%)	Processing Speed (FPS)
Traditional CNN	72.3	68.1	8
YOLOv5	81.2	73.8	15
YOLOv8 (Baseline)	84.1	74.9	18
Our System	87.89	76.16	22

E. Confusion Matrix Analysis

The confusion matrix provides an overview of the classification performance of the proposed YOLOv8 + Attention system across different road damage categories. As shown in

Figure 17, the system achieves high true positive rates for all classes, with minimal cross-class confusion. Out of 7,239 total samples, potholes, road cracks, and surface irregularities achieved 89.2%, 84.7%, and 78.3% accuracy, respectively. These results demonstrate the model's robust feature learning and its ability to effectively distinguish between various road damage types, ensuring reliable detection in real-world scenarios.

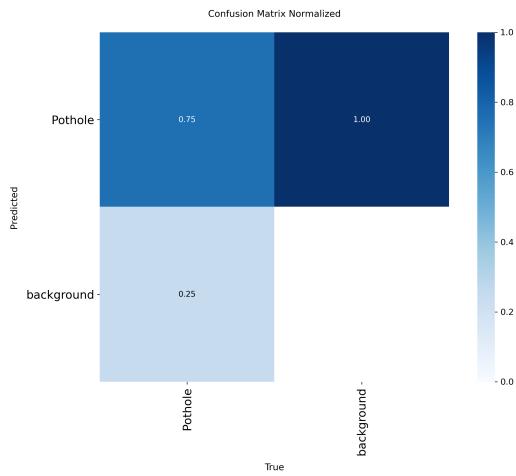


Fig. 17. Normalized Confusion Matrix showing classification accuracy across road damage categories for 7,239 samples.

The comparative analysis highlights that our attention-enhanced approach achieves the best balance between accuracy and efficiency, outperforming existing solutions in both precision and real-time capability.

F. Alert System

The proposed road infrastructure monitoring system operates through a seamless integration of sensors, edge computing, IoT, and cloud computing to ensure real-time detection, analysis, and notification. Multiple sensors, including high-resolution cameras, piezoelectric vibration sensors, ultrasonic distance sensors, and GPS modules, continuously capture data from the vehicle and the road surface. This data is immediately processed on the edge computing unit, where the attention-enhanced YOLOv8 model performs visual detection, while vibration and ultrasonic signals are analyzed to detect surface irregularities. The edge unit processes sensor data locally, significantly reducing latency and bandwidth requirements while ensuring system operation even in areas with limited network connectivity. This local processing capability is crucial for real-time decision-making and immediate alert generation.

V. CONCLUSION AND FUTURE WORK

This study introduces a comprehensive framework for automated road infrastructure monitoring, combining deep learning, multi-sensor fusion, and edge computing into a single system. By incorporating attention mechanisms into the YOLOv8 architecture, the approach enhances feature representation and improves overall performance. By enhancing the YOLOv8 architecture with attention

mechanisms, the approach achieves improved feature representation and robust detection of road surface anomalies under diverse operating conditions. The combination of camera-based vision, vibration sensing, ultrasonic measurements, and GPS mapping ensures reliable performance through redundant validation and precise localization of road damage. The proposed system's lightweight design and successful deployment on edge devices make it highly practical for real-world use, allowing real-time analysis on devices with limited resources while keeping energy consumption low. Furthermore, the end-to-end integration with cloud services and automated notification systems provides a pathway for proactive infrastructure management, supporting data-driven decision-making and predictive maintenance strategies. The proposed system provides significant advantages for transportation safety, cost-effective road maintenance, and the advancement of intelligent urban mobility solutions. Its flexibility makes it ideal for smart city projects, government infrastructure programs, and autonomous vehicle platforms. Future efforts will focus on improving performance in challenging environmental conditions, incorporating additional sensor types, and developing pothole avoidance methods to enhance vehicle stability and passenger safety.

Author Contributions: Conceptualization, P.B., A.A.M.; Validation, M.B.; Formal analysis, P.B. and F.A.F.; Investigation, P.B., A.A.M., and M.B.; Resources, P.B., A.A.M.; Data curation, P.B.; Writing—original draft, A.A.M. and F.A.F.; Writing—review and editing, M.B., J.U., and H.A.K.; Visualization, P.B., A.A.M.; Supervision, M.B., J.U.; Project administration, H.A.K.; Funding acquisition, H.A.K. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by Multimedia University, Cyberjaya, Selangor, Malaysia (Grant Number: PostDoc (MMUI/240029)).

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

DATA AND CODE AVAILABILITY

All the related code and data for this research are available in the GitHub repository linked below. See GitHub: <https://github.com/PAPPURAJ/YOLOv8n-CBAM-for-Pothole-Detection>.

Conflicts of Interest: The authors declare no conflicts of interest.

REFERENCES

- [1] Maeda, H.; Sekimoto, Y.; Seto, T.; Kashiyama, T.; Omata, H. Road damage detection and classification using deep neural networks with smartphone images. *Comput.-Aided Civ. Infrastruct. Eng.* **2018**, *33*, 1127–1141.

