

A Differentiable Gating Mechanism for DETR: Improving Attention Efficiency in Real-Time Road Anomaly Detection

Noor Misbah

Department of Computer Science and Engineering, JSS Science and Technology University, Mysore, Karnataka, India
misbah@jssstuniv.in

S. Srinath

Department of Computer Science and Engineering, JSS Science and Technology University, Mysore, Karnataka, India
srinath@jssstuniv.in (corresponding author)

R. Rakshitha

Department of Computer Science and Engineering, JSS Science and Technology University, Mysore, Karnataka, India
rakshitha.r@jssstuniv.in

Received: 3 November 2025 | Revised: 24 November 2025 | Accepted: 7 December 2025

Licensed under a CC-BY 4.0 license | Copyright (c) by the authors | DOI: <https://doi.org/10.48084/etasr.15977>

ABSTRACT

Accurate detection of road-surface anomalies such as potholes and bumps, along with safety-critical dynamic objects including vehicles and pedestrians, is essential for ensuring traffic safety and enabling reliable autonomous navigation. In this work, "anomalies" refer specifically to static road defects, whereas dynamic objects are treated as safety-relevant events that require immediate attention by intelligent systems. Conventional convolution-based detectors like Faster Region-based Convolutional Neural Network (Faster R-CNN), Single Shot MultiBox Detector (SSD), and You Only Look Once (YOLO) perform well on structured objects but struggle to capture long-range contextual dependencies, limiting performance in complex scenes. Transformer-based models such as the Detection Transformer (DETR) overcome these limitations through global self-attention but suffer from redundant attention activations and slow convergence. To address this, we introduce a Differentiable Gating Mechanism integrated into the encoder's self-attention layers of DETR, employing learnable sigmoid-based gates to selectively emphasize informative heads while suppressing redundant ones. Experiments on a custom COCO-annotated dataset of over 4,700 road images demonstrate that the proposed model improves mean Average Precision (mAP)_{@0.5} from 82.9% to 96.2%, increases mean Intersection over Union (mIoU) from 0.79 to 0.84, reduces trainable parameters by 56%, and achieves a 4.58× faster per image inference time (147.6 ms to 32.2 ms). These results confirm that adaptive gating enhances attention efficiency, accelerates convergence, and significantly improves detection accuracy for real-time road anomaly detection.

Keywords-DETR; transformer-based object detection; differentiable gating; attention mechanism; road anomaly detection; autonomous driving; deep learning

I. INTRODUCTION

Accurate detection of road-surface anomalies such as potholes, cracks, and bumps, along with safety-critical dynamic objects such as pedestrians and vehicles, is essential for safe and efficient transportation. While pedestrians and vehicles are not structural road anomalies, in autonomous navigation they function as dynamic anomalies, representing sudden or abnormal events in the driving environment that demand

immediate attention from intelligent systems. Early identification of these defects minimizes vehicle damage, supports preventive maintenance, and enhances passenger safety [1]. Deep learning-based vision systems have significantly improved road inspection accuracy [2], establishing computer vision as a core component of Intelligent Transportation Systems (ITS) and autonomous driving [3]. Early approaches relied on Convolutional Neural Network (CNN)-based detectors such as Faster Region-based

Convolutional Neural Network (Faster R-CNN) [4], Single Shot MultiBox Detector (SSD) [5], and You Only Look Once (YOLO) [6], which achieve strong accuracy and real-time performance. For road applications, variants including ML-YOLO [3], POT-YOLO achieving 96.8% mean Average Precision (mAP) and 52 Frames per Second (FPS) [7], and YOLOv5 reporting 94.7% mAP on diverse obstacle datasets [8] have shown good performance. However, CNN-based approaches depend on anchor boxes and Non-Maximum Suppression (NMS), and their limited receptive fields reduce effectiveness for small, irregular, or structurally complex anomalies [7-9].

Transformer-based detectors introduced global self-attention and fully end-to-end optimization. Detection Transformer (DETR) [10] removes anchors and NMS but suffers from slow convergence and redundant attention activations [11]. Improvements such as Deformable DETR (10× faster convergence) [12], Conditional DETR (up to 90% fewer training epochs) [13], DN-DETR (denoising-based acceleration) [14], and DINO (57.8% mAP on COCO) [15] alleviate training challenges but still treat all attention heads uniformly. Domain-specific adaptations like Pavement-DETR (91.3% mAP, 52 FPS) [15] and transformer-based autonomous driving models ($\approx 92\%$ accuracy on KITTI) [16] show strong potential but do not explicitly address head-level redundancy. Lightweight hybrids such as XFCOS [17] and multimodal systems like DamageQwen [18] further highlight the breadth of recent progress. Independent studies on attention redundancy reveal that many heads contribute little to performance. Authors in [19] showed that up to 80% of heads can be pruned without significant loss. Authors in [20] demonstrated that pruning 40% of BERT's heads can improve inference speed by 17%. Sparse Transformers [21] and Synthesizer [22] further support reducing redundant attention computation.

Despite these insights, transformer-based road anomaly detection still lacks adaptive, learnable mechanisms to regulate per-head attention activity, especially for heterogeneous datasets representing real-world roads. To address this gap, this work proposes a Differentiable Gating Mechanism for DETR

that adaptively modulates attention head activity in the encoder, suppressing redundancy and enhancing informative feature flow. The key contributions of proposed work are:

- A lightweight differentiable gating module for dynamic, per-head attention regulation in DETR.
- Seamless integration into the DETR encoder without modifying its end-to-end set prediction design or increasing computational cost.
- Extensive evaluation on a custom COCO-annotated road anomaly dataset, demonstrating improvements over baseline DETR.

II. METHODOLOGY

A. Dataset Description

The dataset used in this study was collected using an Android smartphone equipped with a 13 MP rear camera, capturing videos at a resolution of 1920×1080 pixels and a frame rate of 30 FPS. Recordings were carried out across diverse urban and suburban road environments, covering variations in surface quality, traffic density, and illumination conditions. Video sequences were processed to extract individual frames at 30 FPS, which were then manually annotated with bounding boxes for four target categories: potholes, speed bumps, vehicles, and pedestrians. The annotations were initially created in Pascal VOC format [23] and later converted into the COCO JSON format [24] to ensure compatibility with PyTorch-based object detection frameworks and widely used evaluation protocols. To enhance dataset diversity and model robustness, augmentation techniques such as brightness adjustment, darkening, rotation, and horizontal flipping were applied. The final dataset comprises over 4,700 annotated images. Each image may contain one or more target instances (bump, pothole, vehicle, pedestrian), representing both simple and complex real-world scenes. This multi-class, multi-label structure provides a realistic benchmark for evaluating object detection models in road environments. Representative sample images from the dataset are illustrated in Figure 1.



Fig. 1. Sample images from the private dataset with custom annotations.

B. Experimental Setup and Evaluation Metrics

All experiments were conducted using both the baseline DETR and the proposed Differentiable Gating DETR. Both models used the same ResNet-50 backbone, input resolution, and training hyperparameters for a fair comparison. Training was performed on Google Colab using an NVIDIA Tesla T4 GPU (15 GB VRAM), PyTorch 2.8.0, and CUDA 12.6. The dataset was split into 80% training, 10% validation, and 10% testing. Models were trained for 30 epochs using the AdamW optimizer with a learning rate of 1×10^{-4} , batch size of 8, and a cosine learning rate scheduler with weight decay of 1×10^{-4} . The proposed model integrates a lightweight differentiable gating module inside each encoder layer to regulate head activations with minimal parameter overhead.

To evaluate detection performance, Intersection over Union (IoU)-based metrics were employed to assess both localization and classification accuracy. The IoU between predicted and ground truth bounding boxes is computed as shown in (1):

$$\text{IoU} = \frac{B_{pred} \cap B_{gt}}{B_{pred} \cup B_{gt}} \quad (1)$$

where B_{pred} and B_{gt} denote predicted and ground truth bounding boxes, respectively. Predictions with $\text{IoU} \geq 0.5$ were considered correct detections. The mAP measures the average detection accuracy across all classes by computing the Average Precision (AP) for each class from the Precision-Recall curve at a fixed IoU threshold ($\text{IoU} \geq 0.5$), and then averaging these AP values. The mean Intersection over Union (mIoU) represents the average IoU between predicted and ground truth bounding boxes across all classes, providing an overall measure of localization accuracy.

In addition to IoU and mIoU, Precision, Recall, and the F1-score were computed to evaluate the classification reliability of the detected anomalies. These metrics quantify how effectively the model distinguishes true anomalies from background or non-relevant regions. Precision represents the proportion of correctly identified anomalies among all predicted anomalies, Recall indicates the proportion of correctly detected anomalies among all actual anomalies, and the F1-score provides a harmonic mean between Precision and Recall to balance detection accuracy and completeness. They are defined as follows:

$$\text{Precision} = \frac{TP}{(TP+FP)} \quad (2)$$

$$\text{Recall} = \frac{TP}{(TP+FN)} \quad (3)$$

$$\text{F1-score} = \frac{2 \times (\text{Precision} \times \text{Recall})}{(\text{Precision} + \text{Recall})} \quad (4)$$

where TP , FP , and FN denote true positives, false positives, and false negatives, respectively. These metrics collectively assess localization accuracy, classification reliability, and model robustness, enabling comprehensive comparison between the baseline DETR and the proposed Differentiable Gating DETR.

C. Revisiting DETR Model

The proposed work builds upon the DETR architecture by authors in [10], leveraging its fully end-to-end detection capability and fine-tuning it on a custom COCO-annotated road dataset. DETR removes anchor design and NMS by formulating object detection as a direct set prediction task using a transformer encoder-decoder architecture (Figure 2).

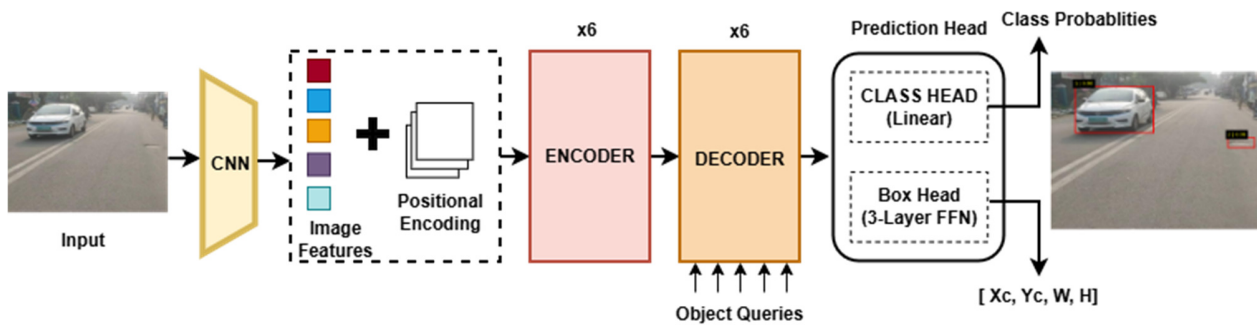


Fig. 2. Baseline DETR architecture with CNN feature extraction, positional encoding, transformer encoder-decoder, and prediction heads.

A ResNet-50 backbone extracts convolutional feature maps, which are flattened and enriched with sine-cosine positional encodings to preserve spatial information. The transformer encoder, composed of Multi-Head Self-Attention (MHSA) and feed-forward layers, aggregates global context to capture long-range dependencies. The decoder operates on a fixed set of learned object queries, refining them through cross-attention with encoder features. Each decoded query is fed into a classification head and a bounding-box regression head, producing normalized coordinates.

The baseline DETR model is implemented using the facebook/detr-resnet-50 checkpoint within the Hugging Face

AutoModelForObjectDetection framework. A ResNet-50 backbone extracts feature maps, which are flattened and enriched with sine-cosine positional encodings before entering the transformer encoder-decoder layers. The encoder employs MHSA to capture global contextual relationships, whereas the decoder refines a fixed set of learned object queries through cross-attention with encoder features. Each query is processed through classification and bounding-box regression heads following DETR's set prediction design. This architectural setup provides the foundation for integrating the proposed Differentiable Gating Mechanism, which modulates per-head attention contributions without altering DETR's end-to-end structure.

D. Enhanced Encoder with Differentiable Gating Mechanism

In the baseline DETR encoder, MHSA is used to learn global contextual relationships, where each attention head captures a different pattern or spatial dependency. However, all heads contribute equally during aggregation even though many may produce redundant or low-value representations. To regulate this imbalance, we introduce a Differentiable Gating Mechanism that assigns a learned importance weight to each head. A gate is a scalar in $[0, 1]$, obtained by applying a sigmoid function to a trainable logit (an unconstrained parameter). This allows the model to emphasize informative heads and suppress less useful ones in a fully differentiable manner.

1) Differentiable Gating Formulation

The mathematical formulation of the proposed Differentiable Gating Mechanism is presented in (5)–(8). Given query, key, and value matrices $Q, K, V \in \mathbb{R}^d$, the i -th head in a standard MHSA produces:

$$O_i = \text{Attention}(Q_i, K_i, V_i) \quad (5)$$

where $Q \in \mathbb{R}^{d_h}$ and $d_h = d/H$ for H heads. The importance of each head is controlled by a learnable gate, given by:

$$g_i = \sigma(\theta_i), \quad \theta_i \in \mathbb{R} \quad (6)$$

where θ_i is the trainable logit and $\sigma(\cdot)$ is the sigmoid function that maps it to $[0, 1]$. The gated attention output for the i -th head (O^i) is given by:

$$O^i = [g_1 O_1 \parallel g_2 O_2 \parallel \dots \parallel g_H O_H] \quad (7)$$

where " \parallel " denotes concatenation. This preserves full differentiability and allows the network to learn head-level importance during training.

E. Integration into DETR Encoder

The gating operator is inserted after each MHSA block and before the residual connection, as illustrated in Figure 3. The modification occurs immediately after the multi-head attention output and before the residual connection and layer normalization. Each encoder layer thus performs:

$$\text{Encoder_Layer}(x) = \text{LayerNorm}\left(x + \text{DG}(\text{MHSA}(x))\right) \quad (8)$$

where x represents the input to the encoder layer, and $\text{DG}(\cdot)$ denotes the differentiable gating operator.

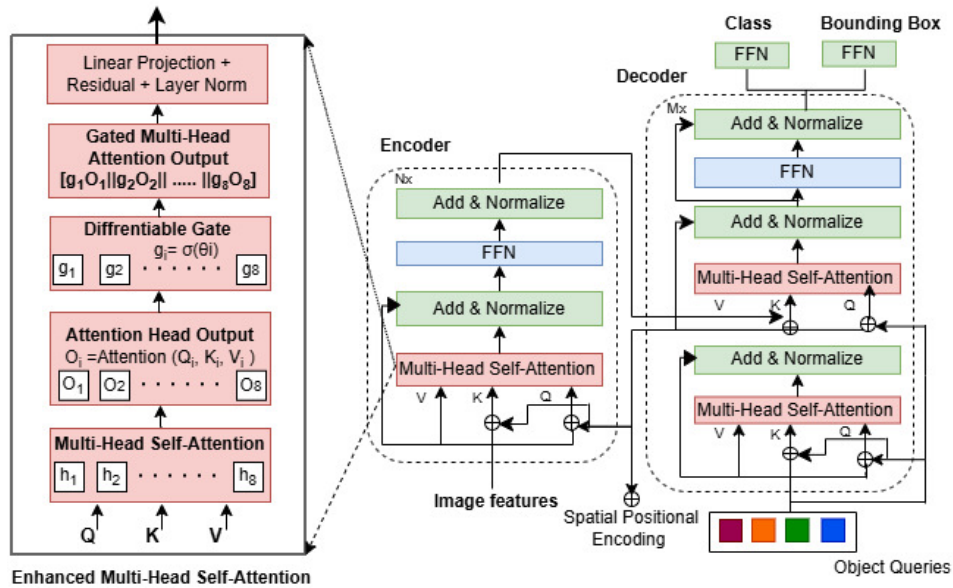


Fig. 3. Proposed Differentiable Gating Mechanism integrated into the DETR encoder.

This integration adds minimal computational overhead and does not modify DETR's architecture, matching loss, or training pipeline. A lightweight PyTorch wrapper attaches the gating module to all encoder layers of the pretrained facebook/detr-resnet-50 model, enabling end-to-end optimization under the same training settings as the baseline.

The enhanced encoder adapts the contribution of each attention head, reducing redundancy and suppressing noisy contextual relations. This improves feature selectivity, leads to more efficient attention computation, and accelerates inference. The mechanism remains fully differentiable and introduces minimal parameter overhead.

III. RESULTS AND DISCUSSION

A. Quantitative Results

This section presents the quantitative comparison between the baseline DETR and the proposed Differentiable Gating DETR. Table I reports the model complexity and computational efficiency, including model size, parameter count, and average training and inference times. Table II summarizes the detection performance on test and validation datasets using mAP, mIoU, Precision, Recall, and F1-score.

TABLE I. MODEL COMPLEXITY AND EFFICIENCY COMPARISON BETWEEN THE BASELINE DETR AND THE PROPOSED DIFFERENTIABLE GATING DETR

Model	DETR	Proposed DETR
Model size (MB)	158	159.01
Total parameters	41,502,666	41,524,816
Trainable parameters	41,280,266	18,069,904
Avg. training time per epoch (min)	14.21	7.68
Avg. inference time per image (ms)	147.61	32.2
Total epochs	30	30
Convergence epoch	16	18

TABLE II. QUANTITATIVE PERFORMANCE COMPARISON ON TEST AND VALIDATION DATASETS

Dataset	Test	Test	Valid	Valid
Model	DETR	Proposed model	DETR	Proposed model
mAP@0.50	0.829	0.9622	0.8332	0.9688
mAP@0.75	0.6535	0.784	0.6678	0.7943
mIoU	0.7961	0.8449	0.8145	0.8597
Precision	0.7205	0.8109	0.759	0.9
Recall	0.9378	0.9001	0.9545	0.842
F1-score	0.8149	0.8531	0.8453	0.8699

The proposed model offers substantial computational advantages. Training time per epoch decreases by 45.9%, inference speed per image improves by a factor of 4.58x, and the number of trainable parameters reduces from 41.28M to 18.07M. These improvements reflect the lightweight nature of the gating mechanism and make the model better suited for real-time road anomaly detection.

The proposed model shows consistent accuracy improvements over the baseline. mAP@0.50 increases from 82.90% to 96.22% on the test set, and mIoU rises from 0.7961 to 0.8449, indicating more accurate spatial localization. Precision, Recall, and F1-score also improve, demonstrating more confident and balanced detections across all classes. These accuracy improvements are consistent with prior findings that many attention heads contribute little to performance, and that regulating or pruning redundant heads leads to more efficient and expressive transformer representations [19, 20].

B. Training Behavior

Figure 4 illustrates the training loss curves for both models. The baseline DETR shows an early plateau and converges around the 16th epoch due to limited gradient flow and redundant attention paths, which restrict further loss reduction. In contrast, the proposed Differentiable Gating DETR continues to improve until approximately the 18th epoch, achieving a lower and more stable final loss. This slightly later convergence is not a drawback but a reflection of more effective learning, enabled by the gating mechanism that suppresses redundant attention and stabilizes gradients. Moreover, despite converging a few epochs later, the proposed model trains 45.9% faster per epoch, resulting in a more efficient overall training process.

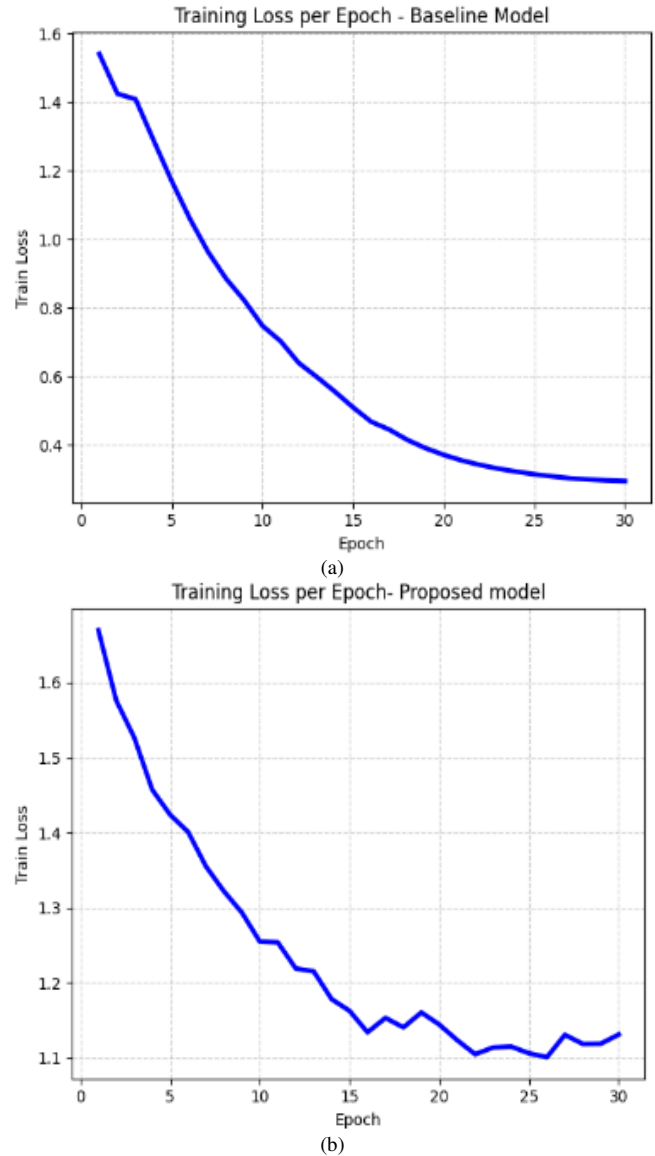


Fig. 4. Training loss progression over 30 epochs for: (a) the baseline DETR model, (b) the proposed Differentiable Gating DETR model.

C. Qualitative Results

Figure 5 presents qualitative comparisons between the baseline DETR model and the proposed Differentiable Gating DETR model on challenging road scenes containing vehicles, pedestrians, potholes, and speed bumps. For visualization clarity, in the baseline DETR outputs, red boxes denote ground-truth annotations and cyan boxes represent predicted detections, whereas in the proposed DETR outputs, green boxes indicate ground truth and red boxes represent predictions.

In the first two rows, both models detect the visible objects, but the proposed model consistently produces higher confidence scores, reflecting stronger feature discrimination and improved contextual understanding. In more challenging scenarios affected by occlusion, low lighting, or small object

size, the performance gap becomes more evident. The baseline DETR often produces weak detections (row 4), misses objects entirely, or fails to detect speed bumps in nighttime scenes (row 3), reflecting its difficulty in handling low-contrast and structurally subtle anomalies. In contrast, the proposed Differentiable Gating DETR provides higher-confidence detections with fewer false negatives, demonstrating improved robustness in challenging scenes and its ability to detect small or low-contrast anomalies.

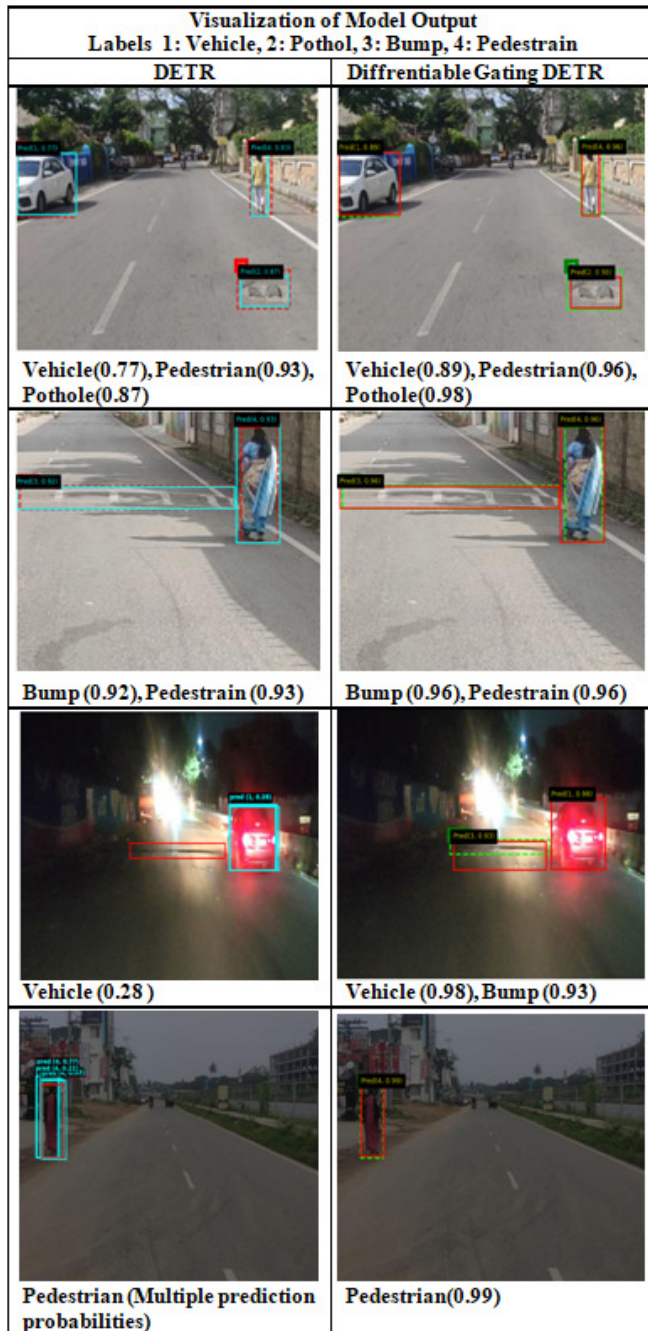


Fig. 5. Qualitative comparison between baseline DETR (left) and the proposed Differentiable Gating DETR (right) on challenging test images.

D. Comparison with Existing Studies

To contextualize the performance of the proposed model within current research, Table III compares the Differentiable Gating DETR with recent transformer- and CNN-based road-anomaly and road-scene detection methods. Because each study employs different datasets, sensing conditions, and evaluation protocols, not all metrics, such as FPS or parameter count, are consistently reported in the original papers. Therefore, the comparison focuses on the metrics that are available and places emphasis on relative performance trends, particularly in terms of detection accuracy, parameter efficiency, and suitability for real-time deployment. This approach ensures a fair, literature-aligned comparison without introducing speculative or unreported values that could bias the analysis.

TABLE III. COMPARISON OF THE PROPOSED MODEL WITH EXISTING STATE-OF-THE-ART METHODS

Method	Dataset	mAP / Accuracy	FPS	Parameters	Remarks
POT-YOLO [8]	Custom road	96.8%	52	Not reported	High speed, anchor-based
Pavement-DETR [15]	Pavement defects	91.3%	52	Not reported	DETR-based, domain specific
YOLOv5 (obstacle detection) [9]	Mixed road	94.7%	Not reported	Not reported	Strong CNN baseline
Deformable DETR [11]	COCO	Not reported	Not reported	Not reported	Faster convergence
Proposed Differentiable Gating DETR	Custom anomaly	96.2%	31	18.07M	Highest precision, fewer params

Compared to existing transformer- and CNN-based approaches, the proposed Differentiable Gating DETR offers competitive or superior detection accuracy while using significantly fewer trainable parameters. Anchor-based models such as POT-YOLO achieve high speed but depend on predefined anchors and NMS, whereas transformer-based methods like Pavement-DETR and Deformable DETR improve global context but do not address redundancy within attention heads. By introducing adaptive head-level gating while preserving DETR's end-to-end design, the proposed model reduces parameters to 18.07M and achieves 96.2% mAP@0.50, demonstrating improved representation quality, faster convergence, and better inference efficiency.

IV. CONCLUSION AND FUTURE ENHANCEMENT

This work introduced a Differentiable Gating Mechanism integrated into the Detection Transformer (DETR) encoder to adaptively regulate attention head contributions for road scene understanding. By learning head-wise importance scores, the proposed model suppresses redundant attention computations and enhances contextual feature extraction. Experiments on a custom dataset of over 4,700 road images showed notable gains over the baseline DETR, with mean Average Precision (mAP)@0.50 improving from 82.9% to 96.2%, mean Intersection over Union (mIoU) from 0.79 to 0.84, and F1-score improving by 4.7%, indicating more precise and balanced detections. The model also achieves a 4.58× faster per image

inference speed and a 56% reduction in trainable parameters (41.28M \rightarrow 18.07M), confirming the computational efficiency introduced by the gating module. Overall, adaptive gating improves attention efficiency, convergence speed, and detection robustness in complex road scenarios. Future work will focus on integrating gating-enabled transformers with lightweight edge-deployment models for real-time use, exploring multimodal fusion with LiDAR or depth data, and incorporating temporal consistency for enhanced reliability in autonomous-driving applications.

REFERENCES

- [1] Y. Safyari, M. Mahdianpari, and H. Shiri, "A Review of Vision-Based Pothole Detection Methods Using Computer Vision and Machine Learning," *Sensors*, vol. 24, no. 17, Sept. 2024, Art. no. 5652, <https://doi.org/10.3390/s24175652>.
- [2] A. K. Bhatt *et al.*, "Advancements in pothole detection techniques: a comprehensive review and comparative analysis," *Discover Artificial Intelligence*, vol. 5, no. 1, Oct. 2025, Art. no. 255, <https://doi.org/10.1007/s44163-025-00297-7>.
- [3] T. Li and G. Li, "Road Defect Identification and Location Method Based on an Improved ML-YOLO Algorithm," *Sensors*, vol. 24, no. 21, Nov. 2024, Art. no. 6783, <https://doi.org/10.3390/s24216783>.
- [4] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," in *Proceedings of the 29th International Conference on Neural Information Processing Systems*, Montreal, Canada, 2015, pp. 91–99.
- [5] W. Liu *et al.*, "SSD: Single Shot MultiBox Detector," in *14th European Conference on Computer Vision*, Amsterdam, Netherlands, 2016, pp. 21–37, https://doi.org/10.1007/978-3-319-46448-0_2.
- [6] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal Speed and Accuracy of Object Detection," arXiv, Apr. 23, 2020, <https://doi.org/10.48550/arXiv.2004.10934>.
- [7] N. Bhavana, M. M. Kodabagi, B. M. Kumar, P. Ajay, N. Muthukumar, and A. Ahilan, "POT-YOLO: Real-Time Road Potholes Detection Using Edge Segmentation-Based Yolo V8 Network," *IEEE Sensors Journal*, vol. 24, no. 15, pp. 24802–24809, Aug. 2024, <https://doi.org/10.1109/JSEN.2024.3399008>.
- [8] P. Mutabarura, N. Muchuka, and D. Seger, "Comparative Evaluation of YOLO Models on an African Road Obstacles Dataset for Real-Time Obstacle Detection," *Engineering, Technology & Applied Science Research*, vol. 15, no. 1, pp. 19045–19051, Feb. 2025, <https://doi.org/10.48084/etasr.9135>.
- [9] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal Loss for Dense Object Detection," in *2017 IEEE International Conference on Computer Vision*, Venice, Italy, 2017, pp. 2999–3007, <https://doi.org/10.1109/ICCV.2017.324>.
- [10] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-End Object Detection with Transformers," in *16th European Conference on Computer Vision*, Glasgow, UK, 2020, pp. 213–229, https://doi.org/10.1007/978-3-030-58452-8_13.
- [11] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable DETR: Deformable Transformers for End-to-End Object Detection," in *9th International Conference on Learning Representations*, Virtual Event, Austria, 2020.
- [12] D. Meng *et al.*, "Conditional DETR for Fast Training Convergence," in *2021 IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, 2021, pp. 3631–3640, <https://doi.org/10.1109/ICCV48922.2021.00363>.
- [13] F. Li, H. Zhang, S. Liu, J. Guo, L. M. Ni, and L. Zhang, "DN-DETR: Accelerate DETR Training by Introducing Query DeNoising," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 4, pp. 2239–2251, Apr. 2024, <https://doi.org/10.1109/TPAMI.2023.3335410>.
- [14] H. Zhang *et al.*, "DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object Detection," arXiv, July 11, 2022, <https://doi.org/10.48550/arXiv.2203.03605>.
- [15] C. Zuo, N. Huang, C. Yuan, and Y. Li, "Pavement-DETR: A High-Precision Real-Time Detection Transformer for Pavement Defect Detection," *Sensors*, vol. 25, no. 8, Apr. 2025, Art. no. 2426, <https://doi.org/10.3390/s25082426>.
- [16] H. Zhao *et al.*, "Improved object detection method for unmanned driving based on Transformers," *Frontiers in Neurobotics*, vol. 18, May 2024, Art. no. 1342126, <https://doi.org/10.3389/fnbot.2024.1342126>.
- [17] Y. Ye, Q. Sun, K. Cheng, X. Shen, and D. Wang, "A lightweight mechanism for vision-transformer-based object detection," *Complex & Intelligent Systems*, vol. 11, no. 7, May 2025, Art. no. 302, <https://doi.org/10.1007/s40747-025-01904-x>.
- [18] Y. Zhang and C. Liu, "Vision-enhanced multi-modal learning framework for non-destructive pavement damage detection," *Automation in Construction*, vol. 177, Sept. 2025, Art. no. 106389, <https://doi.org/10.1016/j.autcon.2025.106389>.
- [19] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov, "Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy, 2019, pp. 5797–5808, <https://doi.org/10.18653/v1/P19-1580>.
- [20] P. Michel, O. Levy, and G. Neubig, "Are Sixteen Heads Really Better than One?," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Vancouver, Canada, 2019, pp. 14037–14047.
- [21] R. Child, S. Gray, A. Radford, and I. Sutskever, "Generating Long Sequences with Sparse Transformers," arXiv, Apr. 23, 2019, <https://doi.org/10.48550/arXiv.1904.10509>.
- [22] Y. Tay, D. Bahri, D. Metzler, D.-C. Juan, Z. Zhao, and C. Zheng, "Synthesizer: Rethinking Self-Attention for Transformer Models," in *Proceedings of the 38th International Conference on Machine Learning*, Online, 2021, pp. 10183–10192.
- [23] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal Visual Object Classes (VOC) Challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, June 2010, <https://doi.org/10.1007/s11263-009-0275-4>.
- [24] T.-Y. Lin *et al.*, "Microsoft COCO: Common Objects in Context," in *13th European Conference on Computer Vision*, Zurich, Switzerland, 2014, pp. 740–755, https://doi.org/10.1007/978-3-319-10602-1_48.