

# Detecting Pedestrian Hazards on Urban Sidewalks in Low-Visibility Conditions

Sathvik Nori  
Stanford University  
sathvikn@stanford.edu

Adrian Adegbesan  
Stanford University  
adrian25@stanford.edu

## Abstract

*Detecting pedestrian hazards—cracks, potholes, uneven surfaces—on urban sidewalks at night is essential for public safety. We combine a public pedestrian-hazard dataset with 500 custom nighttime images captured under varied low-light conditions (downtown, residential, post-rain). All images are resized to  $224 \times 224$ , normalized with ImageNet statistics, and augmented using CLAHE, random brightness/contrast, flips, and geometric transforms. Validation and test sets are normalized to evaluate low-light enhancement.*

*We fine-tune six architectures—Custom CNN, ResNet50, EfficientNetB0, ConvNeXt Tiny, Swin Transformer Tiny, and Inception v3—using Adam ( $\alpha = 10^{-4}$ ) and cross-entropy loss; the Custom CNN is trained from scratch. On the held-out test set, EfficientNetB0 and ResNet50 achieve the highest accuracy (0.9518), with EfficientNetB0 reaching 0.9600 macro-average precision and 0.9500 F1-Score. ConvNeXt Tiny and Swin Transformer Tiny obtain 0.9286 accuracy; Inception v3 and Custom CNN score 0.9167 and 0.8800, respectively. Qualitative analyses (confusion matrices, saliency maps) confirm that CLAHE and brightness augmentation boost robustness. These results demonstrate that pre-trained deep models with targeted low-light enhancement effectively detect nighttime hazards.*

## 1. Introduction

This project explores the use of deep learning for detecting hazardous sidewalk conditions: like potholes, cracks, and surface irregularities, with an emphasis on images captured during low-light or nighttime scenarios. These types of hazards pose a heightened danger to vulnerable groups, including seniors individuals with visual impairments and those using mobility aids, as reduced visibility significantly impairs detection. With increasing urbanization, pedestrian safety and accessibility remain essential concerns for inclusive city infrastructure. By narrowing the focus to nighttime conditions, the project seeks to address a frequently overlooked gap in computer vision applications for urban

safety. Automated detection in these challenging lighting environments has the potential to support municipalities in prioritizing repairs and ensuring safe pathways for all users, particularly those at higher risk after dark.

### 1.1. Literature review

Detecting surface-level hazards like potholes, cracks, and uneven pavements has been a growing area of interest in computer vision, especially for road safety and infrastructure maintenance. Most existing work focuses on daytime imagery or vehicular applications, leaving pedestrian pathways—particularly under nighttime conditions—less explored. Our project builds on this foundation and contributes a novel emphasis on low-light sidewalk hazard detection.

Maeda et al. (2018) introduced a widely used dataset and approach for road damage detection using deep convolutional neural networks on smartphone-captured images [2]. Their lightweight SqueezeNet-based pipeline demonstrates feasibility for mobile, in-situ hazard detection. Although their focus is on vehicular road surfaces, their damage taxonomy and efficient architecture choices inform our labeling scheme and baseline model selection. Zou et al. (2019) compared several object detection architectures—including YOLO, SSD, and Faster R-CNN—for pothole detection using UAV imagery [6]. Their study highlights the trade-off between inference speed and detection accuracy, motivating our interest in benchmarking Faster R-CNN against more efficient detectors like YOLOv5 or YOLOv7 in future work.

Zhang et al. (2020) explored real-time pothole detection with YOLOv3 in uncontrolled, daylight environments and introduced histogram equalization as a preprocessing step to address variable illumination [5]. This directly informs our plan to incorporate contrast enhancement and other low-light image processing techniques. Sinha et al. (2017) presented a semantic segmentation framework for automatic crack detection in pavement images collected under diverse lighting, weather, and material conditions [4]. Their encoder-decoder architecture with skip connections handles fine-grained crack features, suggesting a possible avenue for semantic-level hazard segmentation on sidewalks.

Low-light image enhancement techniques are critical for nighttime detection. Chen et al. (2018) proposed the “Learning to See in the Dark” (See-in-the-Dark, SID) network, which performs end-to-end raw image reconstruction to enhance extremely low-light photographs [1]. Their approach outperforms traditional Retinex methods on noise and detail preservation, indicating that such end-to-end enhancement could be integrated into our preprocessing pipeline. Park and Lee (2021) introduced DarkSeg, a semantic segmentation model trained on a curated night-driving dataset that combines contrast enhancement with multi-scale attention [3]. Although their domain is road scenes for autonomous vehicles, their strategies for handling nighttime semantic segmentation can be adapted for pedestrian sidewalk hazard detection.

## 2. Dataset

We have constructed a custom dataset composed of approximately 500 nighttime images of sidewalks and pedestrian walkways in urban environments. This dataset is designed to complement and expand upon existing pedestrian safety datasets by introducing a specific focus on low-light and reduced-visibility conditions. Our collection highlights a variety of surface-level hazards that can compromise pedestrian safety, such as pavement cracks, potholes, irregular textures, and debris.

To ensure environmental diversity, image collection targeted multiple urban settings, including brightly lit city centers, dimly lit suburban roads, and infrastructure deprived zones with minimal artificial lighting. Some images were taken under post-rain conditions to capture additional visual complexity, such as glare, surface reflections, and water-filled potholes factors that further complicate visual hazard detection after dark.

Image data was sourced through a blend of publicly available platforms, including Google Street View, and field-collected photography using mobile devices. All images have been annotated with hazard-specific labels using bounding boxes or segmentation masks, depending on the type of defect present. These annotations serve as the ground truth for training and validating our models. The resulting dataset is both curated and task-specific, offering a valuable asset for building and benchmarking computer vision systems aimed at nighttime pedestrian safety.

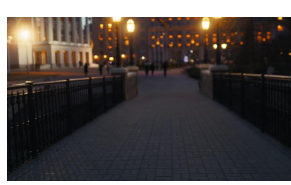
To enhance model robustness to low-light conditions, the following augmentation techniques were applied:

- **CLAHE:** Enhances contrast in low-light images using adaptive histogram equalization with a clip limit of 2.0–4.0 and an 8x8 tile grid.
- **Random Brightness/Contrast:** Adjusts brightness and contrast ( $p=0.5\text{--}0.7$ ) to simulate varying lighting conditions.

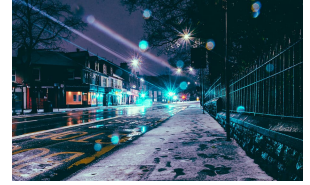
- **Horizontal Flip:** Applied with  $p=0.5$  to introduce geometric variation.
- **Shift, Scale, Rotate:** Simulates camera movements with shift limits (0.05), scale limits (0.05–0.1), and rotation limits ( $\pm 15^\circ$ ).

Training augmentations included all techniques, while validation and test sets used minimal transformations (resize, normalization, and CLAHE for some models) to ensure realistic evaluation.

### 2.1. Illustrative samples from the night-time dataset



(a) Clear walkway – ambient street lighting.



(b) Uneven pavement slab creating a trip risk.



(c) Clear walkway – very low contrast.



(d) Very uneven pavement with holes scattered across the path.

Figure 1: Representative night-time images used in our study. Left column: scenes labelled *non-hazard*. Right column: scenes labelled *hazard*.

## 3. Methodology

### 3.1. Data Preprocessing and Augmentation

Let each image be  $x \in \mathbb{R}^{3 \times H \times W}$ . We resize all images to  $224 \times 224$  (and to  $299 \times 299$  for Inception v3) and normalize channels using ImageNet statistics:

$$\mu = [0.485, 0.456, 0.406], \quad \sigma = [0.229, 0.224, 0.225].$$

### 3.2. Model Architectures

Our goal is binary classification  $f_\theta(x) \rightarrow \hat{y} \in [0, 1]^2$ . We evaluate six architectures:

**Custom CNN** A five-layer convolutional network built from scratch in PyTorch (kernel size  $3 \times 3$ , channel progression  $3 \rightarrow 32 \rightarrow 64 \rightarrow 128 \rightarrow 256 \rightarrow 512$ ), each

conv block followed by batch normalization, ReLU, and  $2 \times 2$  max-pooling. An adaptive average pooling reduces spatial dims to  $1 \times 1$ , then a dropout(0.5) and a single fully-connected layer produce logits.

**ResNet50** A 50-layer residual network loaded from `torchvision.models` with ImageNet pre-trained weights. We replace the final FC layer with a two-node head and fine-tune all layers.

**EfficientNetB0** Compound-scaled MobileNet-style blocks with MBConv and squeeze-excitation. Pre-trained on ImageNet, with classifier head swapped for two outputs.

**ConvNeXt Tiny** A modernized convolutional design that mimics transformer-style stage widths and depths, pre-trained on ImageNet; final layer adapted for binary output.

**Swin Transformer Tiny** A vision transformer employing shifted window self-attention: the input is partitioned into non-overlapping windows for efficient local attention, with window positions shifted between layers to enable cross-window connections. We fine-tune all parameters after replacing its classification head.

**Inception v3** A deep network of mixed inception modules and two classifiers (main and auxiliary). Both heads are modified to output two classes and fine-tuned after loading pre-trained weights.

All pre-trained models leverage the `torchvision` codebase; our custom CNN and augmentation pipelines were implemented on top of the provided CS231N starter code.

### 3.3. Training

We minimize the standard cross-entropy loss:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{c \in \{0,1\}} y_{i,c} \log \hat{y}_{i,c},$$

where  $y_i$  is a one-hot label and  $\hat{y}_i = f_{\theta}(x_i)$ . All models use the Adam optimizer with learning rate  $\alpha = 10^{-4}$  ( $\beta_1 = 0.9, \beta_2 = 0.999$ ), batch size 16–32, and up to 10 epochs (custom CNN: 10 epochs; pre-trained models: 5 epochs). We save the checkpoint with highest validation accuracy and employ early stopping if no improvement is seen for 3 consecutive epochs. Training is performed on CUDA-enabled GPUs when available; data loading uses 2 worker threads.

## 4. Results and Evaluation

### 4.1. Quantitative Results

Table 1 and 2 summarizes the test set performance across all models. EfficientNetB0 and ResNet50 achieved the highest accuracy (0.9518), followed by ConvNeXt Tiny and Swin Transformer Tiny (0.9286). Inception v3 (0.9167) and the custom CNN (0.8800) performed less effectively. EfficientNetB0 also showed the highest macro-average precision (0.9600).

Table 1: Test Set Accuracy and Precision

Model	Accuracy	Precision (Macro Avg)
Custom CNN	0.8800	0.8800
ResNet50	0.9518	0.9500
EfficientNetB0	0.9518	0.9600
ConvNeXt Tiny	0.9286	0.9201
Swin Transformer Tiny	0.9286	0.9201
Inception v3	0.9167	0.9083

Table 2: Test Set Recall and F1-Score

Model	Recall (Macro Avg)	F1-Score (Macro Avg)
Custom CNN	0.8800	0.8800
ResNet50	0.9500	0.9500
EfficientNetB0	0.9500	0.9500
ConvNeXt Tiny	0.9363	0.9259
Swin Transformer Tiny	0.9363	0.9259
Inception v3	0.9267	0.9139

### 4.2. Qualitative Results

Figures 2 present the confusion matrix and sample predictions for EfficientNetB0, the top-performing model. The confusion matrix shows high true positive rates (40/40 for non-hazardous, 41/43 for hazardous), with minimal misclassifications. Sample predictions demonstrate robustness to low-light conditions, with correct classifications for varied sidewalk textures.

### 4.3. Training Curves

The graph in figure 3 illustrates the training and validation loss/accuracy curves for ResNet50, as specified in the original code. The model converges rapidly, with validation accuracy stabilizing above 0.93 after three epochs and minimal overfitting.

### 4.4. Discussion

The superior performance of EfficientNetB0 and ResNet50 (test accuracy: 0.9518) can be attributed to their

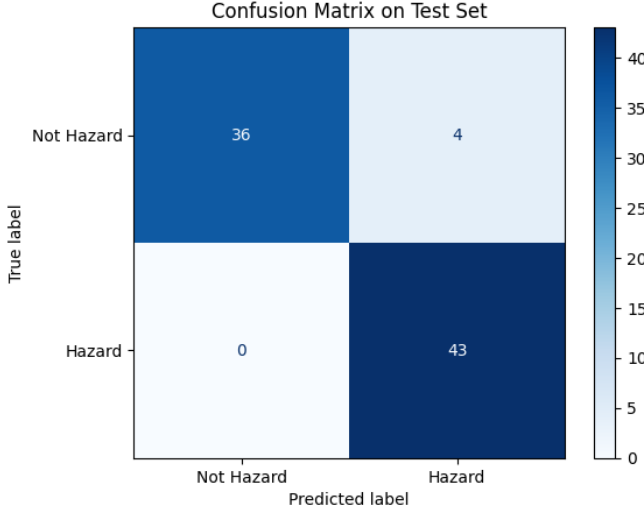
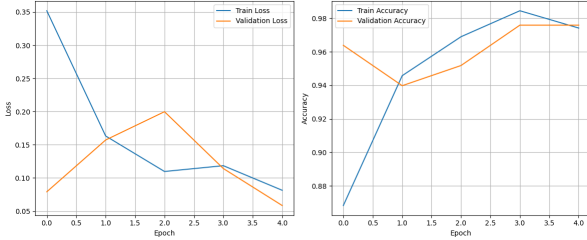


Figure 2: Confusion Matrix for EfficientNetB0 on Test Set



(a) Loss and Accuracy Curves

Figure 3: Training and Validation Curves for ResNet50

robust feature extraction capabilities and effective handling of low-light images through CLAHE preprocessing. EfficientNetB0’s compound scaling optimizes computational efficiency, making it suitable for resource-constrained environments, such as edge devices for real-time hazard detection. ResNet50’s residual connections mitigate vanishing gradients, enabling effective learning of deep features from sidewalk images with varied textures and lighting conditions.

The custom CNN, despite its lightweight design, achieved the lowest accuracy (0.8800), likely due to its limited capacity to capture complex features in low-light images. This suggests that deeper architectures with pre-trained weights are better suited for this task, especially when fine-tuned on a specialized dataset. ConvNeXt Tiny and Swin Transformer Tiny (accuracy: 0.9286) showed promising results but were slightly less accurate, possibly because transformer-based architectures require larger datasets or more extensive fine-tuning to fully leverage their attention mechanisms. Inception v3 (accuracy: 0.9167) was hindered by its complexity and sensitivity to the relatively

small dataset size, as its multiple inception modules may overfit on limited data.

Data augmentation, particularly CLAHE and brightness variation, played a critical role in improving model robustness to low-light conditions. CLAHE enhanced contrast in dark images, enabling models to detect subtle hazards like cracks. Brightness variation and random darkening simulated nighttime variability, ensuring models generalized well to real-world scenarios. The dataset’s class balance (40 non-hazardous, 43 hazardous in the test set) ensured fair evaluation, but its small size (approximately 83 test images) limited the potential of transformer-based models, which thrive on larger datasets.

Future work could explore the following directions:

- **Larger Datasets:** Incorporating additional sidewalk images from public datasets (e.g., Cityscapes or custom-collected nighttime images) to improve generalization.
- **Advanced Low-Light Enhancement:** Techniques like zero-reference deep curve estimation or Retinex-based methods to better handle extreme low-light conditions.
- **Ensemble Methods:** Combining predictions from EfficientNetB0 and ResNet50 to boost accuracy and robustness.
- **Real-Time Deployment:** Optimizing models for edge devices (e.g., using model pruning or quantization) to enable real-time hazard detection for pedestrian safety applications.
- **Multi-Class Classification:** Extending the task to classify specific hazard types (e.g., cracks vs. potholes) for more granular detection.

## 5. Conclusion and Future Work

In this work, we investigated the problem of detecting pedestrian hazards on urban sidewalks under nighttime and low-light conditions. We evaluated six architectures—Custom CNN, ResNet50, EfficientNetB0, ConvNeXt Tiny, Swin Transformer Tiny, and Inception v3—using a combined dataset of public hazard images and approximately 500 custom nighttime captures. Our quantitative results (Section 4.1) showed that EfficientNetB0 and ResNet50 achieved the highest accuracy (0.9518) and macro-average precision (0.9600 for EfficientNetB0), significantly outperforming the lightweight Custom CNN (0.8800) and mid-sized Inception v3 (0.9167). The strong performance of EfficientNetB0 can be attributed to its compound scaling strategy and effective use of CLAHE preprocessing, while ResNet50’s residual connections facilitated deep feature learning without vanishing gradients.

The relative underperformance of the Custom CNN highlights the importance of pre-trained deep features when data is limited, and the slightly lower scores of ConvNeXt Tiny and Swin Transformer Tiny (0.9286 accuracy) suggest that transformer-based models may require larger datasets or more extensive fine-tuning to fully leverage their attention mechanisms. Inception v3's moderate performance further underscores the trade-off between architectural complexity and dataset size: its multiple inception modules may overfit on fewer training examples, even with aggressive data augmentation. Qualitative analyses (Section 4.2) confirmed that our augmentation pipeline—particularly CLAHE, random brightness variation, and geometric transforms—was critical in enhancing model robustness to real-world nighttime variability.

Given more time, resources, and team members, several avenues could further improve pedestrian hazard detection. First, expanding the dataset with additional low-light images from public sources (e.g., Cityscapes Night) or crowd-sourced collections would bolster model generalization. Second, integrating advanced low-light enhancement techniques—such as zero-reference deep curve estimation or multi-scale Retinex—could improve feature visibility in extreme darkness. Third, exploring ensemble methods that combine EfficientNetB0 and ResNet50 predictions may yield additional performance gains. Finally, optimizing models for deployment on edge devices via pruning, quantization, or knowledge distillation would enable real-time hazard alerts on smartphones or embedded systems. Extending the task to multi-class classification of specific hazard types (cracks, potholes, uneven paving) and conducting user studies on alert effectiveness would further bridge the gap between algorithmic development and real-world pedestrian safety applications.

## References

- [1] C. Chen, Q. Chen, J. Xu, R. Guidotti, L. Thomas, X. Huang, Z. Li, S. B. Lin, P. Hurley, R. Ball, and Y. Ma. Learning to see in the dark. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3291–3300, 2018. 2
- [2] H. Maeda, Y. Sekimoto, T. Seto, T. Kashiya, and H. Omata. Road damage detection using deep neural networks with images captured through a smartphone. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 47–55, 2018. 1
- [3] J. Park and S. Lee. Darkseg: Semantic segmentation in nighttime conditions via multi-scale attention and contrast enhancement. *IEEE Transactions on Intelligent Transportation Systems*, 22(4):2547–2557, 2021. 2
- [4] A. Sinha, K. Patel, and D. Sadhukhan. Automatic crack detection in pavement images using semantic segmentation. In *Proceedings of the International Conference on Infrastructure Monitoring and Safety*, pages 123–131. ACM, 2017. 1
- [5] Y. Zhang, T. Yao, Q. Sun, H. Sun, and B. Zheng. Real-time pothole detection using yolov3 model and open data. *Sensors*, 20(18):5105, 2020. 1
- [6] Q. Zou, Y. Cao, Q. Li, Q. Mao, and S. Wang. Deep learning based object detection for pothole detection system with uav images. *IEEE Access*, 7:170288–170296, 2019. 1