

1 **DEEP LEARNING PIPELINE FOR MODELING PAVEMENT CRACKS WITH AN
2 IMBALANCED DATASET (MOPAC)**

3

4

5

6 **Andrew Scouten**

7 Texas State University

8 601 University Drive, San Marcos, Texas 78666-4684

9 Email: yzb2@txstate.edu

10

11 **Haitao Gong, Ph.D.**

12 Texas State University

13 601 University Drive, San Marcos, Texas 78666-4684

14 Tel: 512-245-1826; Email: h_g153@txstate.edu

15

16 **Jelena Tešić, Ph.D.**

17 Texas State University

18 601 University Drive, San Marcos, Texas 78666-4684

19 Tel: 512-245-3409; Email: jtesic@txstate.edu

20

21 **Feng Wang, Ph.D., Corresponding Author**

22 Texas State University

23 601 University Drive, San Marcos, Texas 78666-4684

24 Tel: 512-245-5258; Email: f_w34@txstate.edu

25

26

27

28 Word Count: 4317 words + 7 table(s) × 250 = 6067 words

29

30

31

32

33

34

35 Submission Date: August 1, 2024

1 ABSTRACT

2 Automated pavement crack detection through automatic scanning and processing of pavement
3 surface images is a critical task for efficient assessment of pavement conditions and informed
4 decision-making for pavement maintenance activities. Recently, the adoption of Deep Learning
5 has significantly improved pavement distress detection performance. However, current datasets
6 and studies in this field predominantly focus on standard distress classes, such as longitudinal and
7 transverse cracks, while rare classes or detailed categorizations are often neglected. This study
8 introduces the MoPac algorithm: a novel approach that leverages deep neural networks (DNNs)
9 through a cascade modeling technique combined with low-shot learning, with a particular focus on
10 imbalanced distress classes, to enhance the accuracy and efficiency of pavement crack detection.
11 Our proposed method utilizes a multi-stage DNN cascade to progressively refine crack detection,
12 starting from a broad classification stage and narrowing down to precise localization and severity
13 assessment. The integration of low-shot learning facilitates practical training on limited datasets
14 by incorporating transfer learning and data augmentation strategies. This hybrid approach not
15 only improves detection performance on diverse pavement conditions but also reduces the need for
16 extensive annotated data. Empirical results demonstrate significant advancements in the accuracy
17 and robustness of MoPac detection compared to state-of-the-art methods.

18

19 *Keywords:* Automated pavement crack detection, Deep learning, YOLO, Imbalanced Dataset

1 INTRODUCTION

2 A pavement management system is crucial for maintaining highway pavements in good shape at
3 a minimum cost—exposures to weathering and traffic loading lead to cracks in the pavements.
4 The detection of pavement cracks ensures on-time repair and prolongs the life of the road. For
5 example, in Texas, there are close to 73 million miles of roads and 18.3 million licensed drivers.
6 The Texas drivers drive an average of over 15.2 thousand miles a year, so they cover close to 230
7 billion miles a year. That means, on average, a single patch of Texas road is driven on more than
8 3100 times a year. That is a lot of usage, and efficient pavement crack detection and repair is cru-
9 cial to keep drivers and passengers safe. Currently, automated pavement crack detection systems
10 are employed to identify and address pavement issues. These systems utilize various sensors and
11 image-processing techniques to detect cracks and other forms of distress. However, achieving ac-
12 curate detection remains challenging due to several factors. The primary challenge arises from the
13 complexities of various distress types, different pavement types, and different pavement conditions
14 related to age, traffic load, environment, and other factors. Pavement distress encompasses a wide
15 range of issues, from minor hairline cracks to severe structural defects. Roads are constructed us-
16 ing different materials, such as asphalt, concrete, and composite materials, each exhibiting unique
17 background characteristics. The existing condition can further complicate the attributes of both
18 distress and the pavement itself. Sensor acquisition also complicates the detection process. The
19 diverse types of sensors used, ranging from high-resolution cameras to less sophisticated devices,
20 contribute to inconsistencies in image quality and data richness. These discrepancies impact the
21 ability of automated systems to perform consistently across different settings and conditions. Addi-
22 tionally, the distribution of distress classes varies in a wide range, which necessitates sophisticated
23 detection methods capable of accurately classifying and assessing a wide range of crack types.

24 There has been a plethora of work in the field in the past decade, as described in Section 3.
25 However, the proposed approaches fail to scale to real-life applications due to the following road-
26 blocks: (1) the lack of open-source datasets and a unified significant benchmark for evaluating
27 detection technologies, and (2) the inadequacy of models in detecting the full range of pavement
28 distress classes. Without widely available datasets that represent diverse and realistic scenarios,
29 it is challenging for researchers to train and test their models under varied conditions. This lack
30 of standardization also makes it difficult to compare different approaches and technologies effec-
31 tively, slowing down progress in the field. Furthermore, many current models fail to detect rare
32 distress classes or provide detailed categorization, which can result in missing critical defects or
33 miss-classifying the severity of the cracks, leading to inadequate maintenance responses.

34 The images used in pavement crack detection are distinct from consumer images captured
35 in everyday settings, and we cannot expect out-of-the-box Deep Neural Network (DNN) solutions
36 to work unless it is a highly controlled experiment. We have found that state-of-the-art models
37 often struggle with pavement image variability, leading to incomplete or inaccurate assessments
38 (1). There is a pressing need for advanced automated solutions that can address the specific de-
39 mands of pavement crack detection in the images captured under realistic conditions. Thus, we
40 propose the novel MoPac algorithm, short for the **Modeling Pavement Cracks**. MoPac provides
41 a more robust, adaptable, and standardized system. MoPac is a cascaded deep learning algorithm
42 that incorporates improvements for pavement imagery to detect and identify cracks. The algorithm
43 was named after a nickname for the Loop 1 Expressway in Austin, Texas, MoPac. MoPac sets a
44 skeleton baseline and introduces an innovative framework for practical deployment in road main-
45 tenance and safety assessment. MoPaC will enhance the efficiency of road maintenance processes,

- 1 help extend the lifespan of infrastructure, and ultimately ensure safer and more reliable roadways
- 2 for all users.

3 RELATED WORK

4 In the past couple of years, pavement crack detection research has solely focused on adapting and
5 improving state-of-the-art DNN systems to identify cracks accurately and efficiently. The work
6 stems from the significant advances in the field of computer vision in the past decade in object de-
7 tection using DNN. Several of those DNN networks were shown to be quite reasonable in detecting
8 cracks: *YOLO* (You Only Look Once) revolutionized the field with its unified approach, predicting
9 bounding boxes and class probabilities directly from full images in a single forward pass, enabling
10 real-time processing while maintaining high accuracy, developed by Redmon et al. (2). Subsequent
11 versions of *YOLO*, such as *YOLOv2* (3), *YOLOv3* (4), *YOLOv4* (5), and beyond have incorporated
12 advanced techniques to achieve better performances, including multi-scale prediction, improved
13 network architectures, and data augmentation methods. Among the *YOLO* family, *YOLOv5* is
14 widely adopted for applications and serves in many studies as the base model for developing new
15 approaches. The latest version was proposed by Wang et al. (6), which proposed the concept of
16 Programmable Gradient Information (PGI) for more reliable training. The *Mask R-CNN* utilizes a
17 Region Proposal Network (RPN) to generate potential object regions, which are then refined by a
18 Fast *R-CNN* detector. This technique achieves remarkable precision in real-time object detection
19 and predicts segmentation masks in parallel with bounding box recognition He et al. (7). *Cas-
20 ccade R-CNN* further enhances detection performance by employing a multi-stage object detection
21 architecture that progressively improves the quality of detected objects through a series of detec-
22 tors trained with increasing IoU thresholds introduced by Cai and Vasconcelos (8). The *U-Net*'s
23 encoder-decoder architecture enables precise localization Ronneberger et al. (9), and the *DeepLab*
24 employs convolution and pyramid pooling modules to capture multi-scale context Chen et al. (10).

25 Many studies have been conducted to evaluate the effectiveness of state-of-the-art object
26 detection methods for crack detection. These studies focus on the evaluation of detection perfor-
27 mance over a specific dataset, with little or no modification to the methods themselves. Majid-
28 ifard et al. (11) applied both *YOLO* v2 and *Faster R-CNN* on a dataset developed with Google
29 Street View images. According to the results, *YOLO* v2 outperformed *Faster R-CNN*, achieving
30 an F1 score of 0.84 compared to *Faster R-CNN*'s score of 0.65. Among the studies based on
31 another dataset that includes pavement images collected by smartphones in three different coun-
32 tries (12), *YOLO*-based models have also been found to outperform other methods, such as *Faster*
33 *R-CNN* and *EfficientDet* (13). By comparing the studies using datasets from only one country
34 and all three countries (13)(14), it is noticed that the performance dropped significantly with the
35 larger dataset, which includes all three countries. A study using high-resolution ($1,800 \times 1,200$)
36 three-dimensional (3D) asphalt and concrete pavement images reported that the *Faster R-CNN* and
37 *YOLO* v3 delivered the same level of performance for the automated distress detection and classi-
38 fication task, with both methods achieving around 90% accuracy (15). A study applied *YOLO* v5l
39 to a two-dimensional (2D) asphalt concrete pavement set of 3,001 images (16), and they reported
40 similar average precision. However, all types of distress are considered as one class in this study.
41 Neither of the datasets has been made public, so we can only summarize the findings reported.
42 Some researchers combined the two-dimensional (2D) gray-scale images and three-dimensional
43 (3D) laser scanning data and demonstrated the benefit of both, including the 2D and 3D infor-
44 mation in the process (17). A significant trend observed in the literature is the concentration on

1 the most common types of pavement distresses, such as longitudinal and transverse cracks. These
2 types of cracks are the most prevalent on a highway network scale, leading to their dominance in
3 datasets. However, this focus creates a bias towards the more frequent crack types and neglects
4 other critical distress types that might be less common. Moreover, models that cannot detect the
5 full range of pavement distress classes have limited practical application.

6 In the area of deep neural network (DNN) modeling based on the characteristics of pave-
7 ment cracking, Safaei et al. proposed the segmentation with an adaptive threshold approach, show-
8 ing promising results in a small setup (130 images) for medium and large cracks (18). Nguyen et al.
9 investigate the bigger picture of deep learning methods for crack identification in asphalt pavement
10 (19) as they compare the state-of-the-art for crack classification, crack object detection, pixel-level
11 crack segmentation, generative adversarial networks (GANs) for crack segmentation, and crack
12 identification using unsupervised learning. They propose using a high-resolution 2D image dataset
13 with 10,000 images to evaluate the 25 deep-learning methods and found the Faster R-CNN per-
14 forms best for crack object detection. Cheng et al. show that the improved crack extraction method
15 has an accuracy of the crack classification on a minimal dataset of close to 90% (20). Gao et al.
16 propose to leverage low-rank representation in the deep learning pipeline to discriminate images
17 (or frames) with cracks (21). They suggest a deep convolutional neural network for crack detection
18 leveraging multilevel features and atrous spatial pyramid pooling (ASPP) (21). One study incorpo-
19 rates the latest transformer module with *YOLO* v5, aiming to enhance crack detection by capturing
20 long-range dependencies and learning context information of crack regions. A model trained on
21 a public dataset, including pavement images from India, the Czech Republic, and Japan, achieved
22 F1 scores of 0.6739 and 0.6650 on two online test sets (22). Wu et al. introduce a novel detection
23 model based on *YOLO* v5, enhanced with a lightweight crossed feature pyramid network (CFPN)
24 and an improved loss function (23). The model was tested on a dataset of 7,076 images covering
25 four common pavement distress types. Results show the model excels in challenging conditions
26 like shadows and overlapping objects, achieving a mean average precision (mAP) of 69.3%.

27 In summary, most datasets in pavement crack detection studies primarily include standard
28 recurrent distress classes, such as longitudinal and transverse cracks, and the rare classes or more
29 detailed and specific categorizations are often neglected. Most studies focus on improving the
30 overall performance of detection methods rather than addressing particular distress classes, result-
31 ing in a lack of sufficient attention to the detection of rare distress types. This imbalance leads
32 to models that are proficient at identifying common distresses but cannot accurately detect and
33 address rarer and critical pavement issues, limiting their practical application. In this study, we
34 propose a novel approach that leverages deep neural networks (DNNs) through a cascade model-
35 ing technique combined with low-shot learning, with a particular focus on underrepresented and
36 complex crack types, to enhance the overall effectiveness and practical applicability of pavement
37 crack detection models.

38 METHODOLOGY

39 The increasing role of AI for pavement condition assessment requires building a practical yet
40 robust DNN-based distress detection model that considers data in the wild. To this end, we first
41 collect visually different data under different pavement conditions. Second, we define the labels
42 as visually meaningful to both human experts and computer vision algorithms and label the data.
43 Third, we develop a novel object recognition system based on the deep learning You Only Look
44 Once (*YOLO*) algorithm and propose improvements based on the unique nature of the pavement

1 imagery.

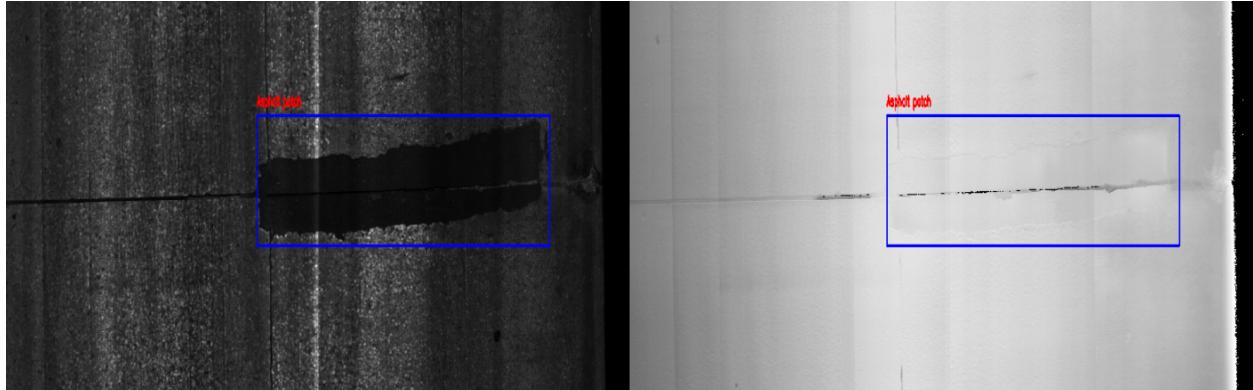


FIGURE 1 A sample of 2D (left) and 3D (right) pavement surface images with an Asphalt Patch annotated.

2 Step 1: Data in the Wild Acquisition

3 Pavement surface images used in this study are 2D/3D image data collected in Texas with a specialized
 4 3D laser camera sensor. The 3D laser sensor is mounted on the top back of a van and employs
 5 a downward-pointing orientation to project a transverse laser line across the pavement as the van
 6 moves. A line camera within the sensor captures the image of this laser line, focusing on the laser's
 7 wavelength to minimize ambient light interference. Using the triangulation method, the sensor cal-
 8 culates the relative depth and surface profile of the pavement from the deformation of the laser line
 9 on uneven surfaces. This technique allows the construction of a detailed 3D transverse profile of
 10 the pavement, significantly enhancing the accuracy and consistency of surface imaging compared
 11 to traditional cameras, especially under variable lighting conditions. With high-resolution 2d/3D
 12 images, we can focus on addressing the complexities of the crack and pavement conditions, which
 13 remain challenging tasks. Each pavement surface image is roughly one lane in width and about 5
 14 meters in length per AASHTO standard recommendations(24).

15 Figure 1 illustrates a sample of the 2D and corresponding 3D images. In the 2D image,
 16 each pixel value represents the intensity of the corresponding area on the pavement surface, while
 17 in the 3D image, each pixel value indicates the range or relative height. A total of 12 highway
 18 sections with Joint Concrete Pavement (JCP) surfaces were selected for this study. Each image has
 19 a resolution of 900×1536 pixels, which covers approximately 12 feet in the transverse direction
 20 and 6-20 feet in the longitudinal direction, depending on the driving speed.

21 Step 2: Crack Region Labeling

22 A total of 8,016 images were collected. These surface images were manually annotated based
 23 on the distress classification defined in the Pavement Rater's Manual by the Texas Department of
 24 Transportation (TxDOT) and separated into 13 visual classes. Each distress identified is labeled
 25 with a bounding box and the corresponding distress class. The bounding box is supposed to cover
 26 the extent of a complete cracking within an image, with the background included as tight as possi-
 27 ble. The number of annotated instances of distress, the intensity image with the labeled region, and
 28 the corresponding 3D laser image representatives are summarized in Table 1 for recurring labels
 29 and in Table 2 for rare labels. Deep learning networks require a lot of training data to localize
 30 cracks reliably. Thus, we separated our dataset into two classes: the recurring classes and the rare

1 classes.

TABLE 1 Recurring pavement crack label abbreviation, number of annotated instances of each class in the entire set, TxDOT explanation of the type of the crack, and a 2D (intensity image) and 3D (range image) visual class example.

Label →	Recurring classes: number of labeled instances is enough for the robust modeling						
	Joint	L-track	T-crack	A-patch	Spall	C-break	F-joint
Instances	10,286	1512	1,159	602	551	273	243
Description	Joint	Longitudinal crack	Transversal crack	Asphalt Patch	Spall	Corner Break	Failed Joint
Intensity							
Range							

TABLE 2 Rare pavement crack label abbreviation, number of annotated instances of each class in the entire set, TxDOT explanation of the type of the crack, and a 2D (intensity image) and 3D (range image) visual class example.

Label →	Recurring classes: number of labeled instances is enough for the robust modeling					
	Rare classes: number of labeled instances is not enough for robust modeling					
Instances	Popout	C-patch	D-crack	FC-patch	S-long	P-out
Description	190 Popout	152 Concrete Patch	34 D-cracking	27 Failed concrete patch	9 Sealed longitudinal	5 Punchout
Intensity						
Range						

2 The recurring classes have more than 200 annotated instances in the dataset, as summarized in
 3 Table 1, and the rare classes have less than 200 instances (annotated occurrences) in the dataset, as
 4 summarized in Table 2.

5 Step 3: Visual Deep Learning for Crack Localization and Identification

6 In the third step, we propose three different approaches to modeling: 1. the vanilla deep learning
 7 multi-class modeling, 2. the cascade modeling, and 3. binary modeling to extract features fol-
 8 lowed by traditional machine learning modeling. For the vanilla multi-class, we use YOLOv5 and
 9 YOLOv8 (25). The YOLO framework (2) treats object detection as a single regression problem,
 10 and it uses a single convolutional neural network (CNN) to predict bounding boxes and class proba-
 11 bilities simultaneously. The input image is divided into a grid. Each grid cell is responsible for pre-
 12 dicting bounding boxes and class probabilities for objects whose center falls within that cell. Each

1 grid cell predicts a fixed number of bounding boxes. For each bounding box, the network outputs
2 coordinates (center, width, height), confidence scores, and class probabilities. Each grid cell also
3 predicts the probability distribution over all possible object classes for the detected object. After
4 the network outputs the bounding boxes, multiple boxes may overlap. YOLO uses Non-Maximum
5 Suppression to remove redundant boxes and retain only the most accurate ones. YOLOv5 was
6 implemented and published in 2020, and it brought practical improvements, ease of use, and high
7 performance and has become widely used in the community (25). YOLOv8 introduced further
8 refinements and optimizations in 2023, focusing on improving accuracy and inference speed and
9 integrating advanced techniques like attention mechanisms and enhanced training strategies (25).
10 YOLOv8 incorporates advanced backbone and neck architectures to improve feature extraction
11 and information flow. This includes enhancements in the network’s depth and width, as well as
12 improvements in feature pyramid networks or similar components to handle objects at various
13 scales better. YOLOv8 aims to strike a better balance between accuracy and speed. It typically
14 achieves higher precision and recall compared to its predecessors while maintaining or improv-
15 ing inference speed. YOLOv8 uses dropout, batch normalization, and advanced loss functions to
16 prevent overfitting and improve generalization (25). For the *cascade* modeling, we combine all
17 annotated instances into one robust crack annotation and train the model as a binary model: crack
18 or no crack. This step has the advantage of removing the black box one-shot decision and pro-
19 viding an analyst with candidates where clear decisions cannot be made. The goal of this step is
20 also to help the system learn in the future, reduce the role of human-in-loop over time, and reduce
21 false positives identified by the system. Then, we only consider regions clearly labeled as generic
22 cracks to classify them into the types of cracks. So, after the binary classifier, we employ another
23 multi-label classifier that only considers regions with detected cracks. Suppose the category does
24 not have enough examples. In that case, we test the following robust pipeline that extracts deep
25 features from binary images in the training data and uses a support vector machines (SVM) tradi-
26 tional classifier to find the related classes. Suppose the crack category has only a few examples. In
27 that case, this will allow us to do domain translation and adapt the model for crack identification
28 for the previously unseen class in the dataset, as we have successfully implemented for different
29 non-consumer tasks in (26).

30 EXPERIMENTS

31 We measure the effectiveness of the models using Precision (P), Recall (R), F1 measure (F1), Av-
32 erage Precision (AP) per class, and mean Average Precision (mAP) for all classes. This provides
33 us with a set of comprehensive measures of how well a model can detect and localize cracks in
34 images. For each class, a Precision-Recall (PR) curve is plotted. Precision measures the accuracy
35 of the detections (how many of the detected objects are actually relevant), while Recall measures
36 the model’s ability to find all relevant objects. The Average Precision (AP) for each class is com-
37 puted by integrating the Precision-Recall curve. The mean Average Precision (mAP) is calculated
38 by taking the mean of these AP scores across all classes.

39 Experiment 1. Select Model and Data

40 We compare the classification performance of two different YOLO versions (5 and 8), two different
41 inputs (Intensity and Two-Channel), and two different classification outcomes (multi-label and
42 binary) to pick the most robust classification baseline. The two-channel input is the 3D scanning
43 detection result (concatenated intensity and range images), and it consistently yields better results.
44 The scanning detection results inform the 2D-3D correspondence for formulating classification

TABLE 3 Table for the average performance of different YOLO Models for all 13 classes.**Note that the Precision and Recall are reported for the highest F1.**

Method YOLO Model	Intensity				Intensity Binary				Two-Channel				Two-channel Binary			
	Prec	Rec	mAP	mAP	Prec	Rec	mAP	mAP	Prec	Rec	mAP	mAP	Prec	Rec	mAP	mAP
	vision	all	50	50-95	vision	all	50	50-95	vision	all	50	50-95	vision	all	50	50-95
5su train	0.53	0.27	0.24	0.11	0.72	0.66	0.69	0.36	0.45	0.49	0.28	0.14	0.86	0.78	0.75	0.43
8s train	0.54	0.30	0.27	0.12	0.72	0.69	0.69	0.35	0.57	0.40	0.30	0.15	0.87	0.79	0.78	0.45
5su test	0.55	0.24	0.24	0.10	0.75	0.65	0.71	0.36	0.42	0.35	0.29	0.11	0.80	0.68	0.77	0.38
8s test	0.53	0.28	0.25	0.10	0.73	0.65	0.70	0.35	0.48	0.33	0.30	0.13	0.81	0.72	0.79	0.40

- 1 problems, identifying rare crack types, and informing the final decision, as illustrated in Table 5.
 2 Table 5 also shows that a robust binary classifier performs much better than a multi-class one. The
 3 explanation can be found in Table 2. For six classes, we do not have enough training data to learn
 4 the separation from the related classes. We used this experiment as a guideline to help us decide
 5 on the YOLO version for the subsequent ablation study.

TABLE 4 Validation dataset’s mAP scores for different IoU and confidence thresholds.

Confidence ↓	Prediction IoU Thresholds					
	0.10	0.20	0.30	0.40	0.50	0.75
0.6000	0.7605	0.7590	0.7587	0.7564	0.7530	0.6255
0.5000	0.8046	0.8025	0.8019	0.7968	0.7901	0.6327
0.4000	0.8293	0.8266	0.8247	0.8177	0.8083	0.6317
0.3000	0.8481	0.8443	0.8415	0.8332	0.8210	0.6285
0.2000	0.8621	0.8577	0.8542	0.8447	0.8309	0.6256
0.1000	0.8727	0.8679	0.8636	0.8528	0.8364	0.6145
0.1500	0.8674	0.8631	0.8590	0.8489	0.8332	0.6209
0.0100	0.8775	0.8704	0.8648	0.8523	0.8293	0.5707

6 Experiment 2. Ablation study for NMS IoU, Prediction IoU, and Prediction Confidence Thresholding

8 In this experiment, we report the IoU and Threshold ablations study for the binary YOLOv8 small
 9 model on the training data, as it showed the most robust performance in experiment 1. The IoU,
 10 or Intersection over Union, is a metric commonly used in visual classification tasks, especially in
 11 object detection and image segmentation. IoU is calculated as the ratio of the area of their intersec-
 12 tion to the area of their union, where the intersection is the overlapping area between the predicted
 13 bounding box and the annotated bounding box, and the union is the total area of the predicted
 14 bounding box and the annotated bounding box combined. Thus, the IOU is defined as $\text{IoU} = \frac{\text{Area of Intersection}}{\text{Area of Union}}$. The IoU score evaluates how well the predicted bounding boxes
 15 match the actual bounding boxes. For instance, if an IoU threshold is set, predictions with IoU
 16 greater than this threshold are considered true positives, and the smaller ones are ignored. Prior,
 17 we found that the high default values of IoU can be very restrictive in non-consumer tasks like
 18 crack detection (27).

20 The confidence summarizes the model’s score on how the model is confident that the de-
 21 tected object within the bounding box belongs to a particular class. This overall confidence score
 22 helps in making a final decision on whether to accept or reject a detected object based on the con-

TABLE 5 Validation dataset’s mAP scores for different NMS and prediction IoU thresholds, and confidence threshold 0.001.

NMS IoU	Confidence = 0.001 and Prediction IoU Thresholds							
	0.10	0.20	0.30	0.40	0.50	0.75	AP	
0.7500	0.8548	0.8479	0.8423	0.8302	0.8088	0.5471	0.5573	
0.7000	0.8648	0.8575	0.8514	0.8387	0.8157	0.5471	0.5594	
0.6000	0.8769	0.8690	0.8628	0.8484	0.8226	0.5420	0.5598	
0.5000	0.8815	0.8734	0.8666	0.8511	0.8224	0.5392	0.5583	
0.4000	0.8835	0.8751	0.8670	0.8479	0.8179	0.5400	0.5571	
0.3000	0.8844	0.8742	0.8646	0.8429	0.8122	0.5427	0.5567	
0.2000	0.8792	0.8658	0.8548	0.8311	0.8006	0.5437	0.5542	
0.1000	0.8641	0.8485	0.8377	0.8159	0.7873	0.5424	0.5499	

1 fidence threshold. In Table 5, we summarize the results of our ablation study on the validation set
 2 where we vary the confidence threshold and prediction IoU threshold. Note that mAP scores vary
 3 only slightly if the confidence threshold is between 0.01 and 0.2 and IoU is between 0.1 and 0.3.
 4 Thus, we adopt the confidence score threshold of 0.01 and IoU threshold of 0.3 as the standard
 5 YOLO recommendation for the threshold is 0.01, and it aligns with our experiments. The IoU
 6 of 0.3 is still selective enough not to include false positives when the model is evaluated on the
 7 unseen dataset. Next, for the confidence 0.01 and varying IoU thresholds in prediction, we test the
 8 influence of the NMS IoU setup in Table 5. We find that the NMS IoU of 0.3 in the validation set
 9 produces a good balance between false positives and false negatives at each level.

TABLE 6 Performance of the models per recurring label on the test set. The full description of class labels is in Table 1

Recurrent Label	Mean		Joint		L-crack		T-crack		A-patch		Spall		C-break		F-joint	
	F1	AP	F1	AP	F1	AP	F1	AP	F1	AP	F1	AP	F1	AP	F1	AP
Method ↓																
YOLOv8s	0.552	0.528	0.91	0.96	0.71	0.76	0.66	0.67	0.50	0.51	0.32	0.21	0.44	0.32	0.33	0.27
YOLOv8x	0.538	0.533	0.90	0.95	0.69	0.74	0.65	0.68	0.45	0.49	0.32	0.21	0.37	0.32	0.39	0.34
v8s Bin + SVM	0.502	0.459	0.88	0.92	0.69	0.67	0.62	0.59	0.53	0.50	0.30	0.22	0.26	0.20	0.24	0.11
v8x Bin + SVM	0.465	0.444	0.88	0.92	0.68	0.69	0.62	0.63	0.52	0.48	0.29	0.17	0.23	0.16	0.02	0.05
v8s Cascade	0.476	0.448	0.88	0.92	0.69	0.70	0.62	0.63	0.50	0.46	0.26	0.14	0.24	0.17	0.15	0.12
v8x Cascade	0.489	0.455	0.83	0.85	0.62	0.63	0.61	0.59	0.64	0.60	0.31	0.17	0.32	0.19	0.10	0.15

10 Experiment 3. Recurring and Rare Class Modeling

11 In this experiment, we evaluate four different models and their performances for *each* of the class
 12 on the test set. We analyze mean average precision (mean AP) for recurring labels in Table 5, and
 13 the results paint an interesting picture. The YOLOv8x cascade does not improve F1 and AP for
 14 the most frequent classes (Joint, L-crack, and T-crack), and the baseline model produces the best
 15 results. The cascade modeling significantly improves A-patch modeling, and the AP jumps from
 16 0.51 to 0.6. D-crack and FC-patch also benefit from the cascade approach as rare labels in Table 5,
 17 and the popout is so visually distinctive compared to other classes that the baseline model predicts

1 better than expected.

TABLE 7 Performance of the model per rare class on the test set. The full description of class labels is in Table 2.

Sparse Label Method ↓	Mean		C-patch		D-crack		FC-patch		S-long		P-out	
	F1.	AP	F1.	AP	F1.	AP	F1.	AP	F1.	AP	F1.	AP
YOLOv8s	0.129	0.145	0.13	0.07	0.00	0.15	0.00	0.02	0.00	0.06	0.52	0.43
YOLOv8x	0.093	0.099	0.13	0.05	0.00	0.07	0.00	0.04	0.00	0.02	0.34	0.31
v8s Binary + SVM	0.107	0.060	0.08	0.04	0.23	0.11	0.00	0.00	0.00	0.00	0.23	0.15
v8x Binary + SVM	0.102	0.068	0.08	0.02	0.13	0.10	0.00	0.00	0.00	0.00	0.29	0.22
v8s Cascade	0.132	0.138	0.10	0.02	0.14	0.08	0.00	0.23	0.00	0.00	0.43	0.36
v8x Cascade	0.156	0.125	0.00	0.01	0.09	0.06	0.25	0.23	0.00	0.00	0.44	0.34

2 Summary of Findings

3 In this section, we have conducted three different experiments: the first two were for the parameter
4 selection. At the same time, the third one evaluates the various effects of baseline and three new
5 proposed models for the crack identification task. According to the results of experiment 1, the
6 selection of the baseline model (YOLOv8) and the input (2D/23 combined) are consistent with the
7 findings in the field. The recurrent crack classes with a more significant number of instances were
8 found to produce better performances across the board. Experiment 2 discovered that our models
9 performed better at much lower IoU thresholds, and we selected the NMS IoU of 0.3 and prediction
10 IoU of 0.3 while keeping the confidence level low at 0.001. The experiment highlights the impor-
11 tance of IoU thresholding selection and hyper-parameter choices for NMS. Cascade modeling in
12 Experiment 3 produced a mixed bag of results. For the classes with more than 1000 annotations, it
13 did not improve the performance as it was too restrictive and excluded some true positives, result-
14 ing in the drop of the AP. It did improve the performance of the A-patch, D-crack, and FC-patch
15 classifier.

16 CONCLUSION AND FUTURE WORK

17 In this paper, we propose MoPaC, a state-of-the-art crack localization and identification frame-
18 work that effectively recognizes different types of cracks in the pavement from intensity and range
19 images. First, we have made significant improvements in training data set labeling and region
20 annotation and defined categories that are visually meaningful both for human experts and the
21 automated system. We propose the 2D-3D YOLO-based system that robustly detects joints, lon-
22 gitudinal and traversal cracks, as well as asphalt patches captured over larger patches of the road.
23 MoPac’s next step is to continue developing high-quality training datasets and increase the number
24 of classes with over 1000 annotations to confirm these initial findings.

25 ACKNOWLEDGMENTS

26 The authors would like to acknowledge funding supports received from TxDOT (Research Project
27 #0-7150).

1 AUTHOR CONTRIBUTIONS

2 Scouten led the machine learning implementation and analysis, created tables, and contributed to
3 manuscript preparation and editing. Gong led the creation and manual annotation of the dataset,
4 created figures, and contributed to the literature review, manuscript preparation, and editing. Tešić
5 designed the experiments and data analysis and contributed to the literature review and the manuscript
6 preparation and editing. Wang conceived the research idea, supervised the project, and contributed
7 to the manuscript writing and editing.

8 REFERENCES

- 9 1. Gong, H., J. Tešić, J. Tao, X. H. Luo, and F. Wang, Automated Pavement Crack Detection
10 with Deep Learning Methods: What Are the Main Factors and How to Improve the
11 Performance? *Transportation Research Record*, Vol. 2677, No. 10, 2023, pp. 311–323.
- 12 2. Redmon, J., S. Divvala, R. Girshick, and A. Farhadi, You only look once: Unified, real-
13 time object detection. In *Proceedings of the IEEE conference on computer vision and*
14 *pattern recognition*, 2016, pp. 779–788.
- 15 3. Redmon, J. and A. Farhadi, YOLO9000: better, faster, stronger. In *Proceedings of the*
16 *IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.
- 17 4. Redmon, J. and A. Farhadi, Yolov3: An incremental improvement. *arXiv preprint*
18 *arXiv:1804.02767*, 2018.
- 19 5. Bochkovskiy, A., C.-Y. Wang, and H.-Y. M. Liao, Yolov4: Optimal speed and accuracy of
20 object detection. *arXiv preprint arXiv:2004.10934*, 2020.
- 21 6. Wang, C.-Y., I.-H. Yeh, and H.-Y. M. Liao, Yolov9: Learning what you want to learn using
22 programmable gradient information. *arXiv preprint arXiv:2402.13616*, 2024.
- 23 7. He, K., G. Gkioxari, P. Dollár, and R. Girshick, Mask r-cnn. In *Proceedings of the IEEE*
24 *international conference on computer vision*, 2017, pp. 2961–2969.
- 25 8. Cai, Z. and N. Vasconcelos, Cascade R-CNN: High quality object detection and instance
26 segmentation. *IEEE transactions on pattern analysis and machine intelligence*, Vol. 43,
27 No. 5, 2019, pp. 1483–1498.
- 28 9. Ronneberger, O., P. Fischer, and T. Brox, U-net: Convolutional networks for biomedical
29 image segmentation. In *Medical image computing and computer-assisted intervention–*
30 *MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, pro-*
31 *ceedings, part III 18*, Springer, 2015, pp. 234–241.
- 32 10. Chen, L.-C., G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, Deeplab: Semantic
33 image segmentation with deep convolutional nets, atrous convolution, and fully connected
34 crfs. *IEEE transactions on pattern analysis and machine intelligence*, Vol. 40, No. 4, 2017,
35 pp. 834–848.
- 36 11. Majidifard, H., P. Jin, Y. Adu-Gyamfi, and W. G. Buttlar, Pavement image datasets: A new
37 benchmark dataset to classify and densify pavement distresses. *Transportation Research*
38 *Record*, Vol. 2674, No. 2, 2020, pp. 328–339.
- 39 12. Arya, D., H. Maeda, S. K. Ghosh, D. Toshniwal, A. Mraz, T. Kashiyama, and Y. Sekimoto,
40 Transfer learning-based road damage detection for multiple countries. *arXiv preprint*
41 *arXiv:2008.13101*, 2020.
- 42 13. Arya, D., H. Maeda, S. K. Ghosh, D. Toshniwal, H. Omata, T. Kashiyama, and Y. Sekimoto,
43 Global road damage detection: State-of-the-art solutions. In *2020 IEEE International Conference on Big Data (Big Data)*, IEEE, 2020, pp. 5533–5539.

- 1 14. Mandal, V., L. Uong, and Y. Adu-Gyamfi, Automated road crack detection using deep
2 convolutional neural networks. In *2018 IEEE International Conference on Big Data (Big
3 Data)*, IEEE, 2018, pp. 5212–5215.
- 4 15. Ghosh, R. and O. Smadi, Automated Detection and Classification of Pavement Distresses
5 using 3D Pavement Surface Images and Deep Learning. *Transportation Research Record*,
6 Vol. 2675, No. 9, 2021, pp. 1359–1374.
- 7 16. Hu, G. X., B. L. Hu, Z. Yang, L. Huang, and P. Li, Pavement Crack Detection Method
8 Based on Deep Learning Models. *Wireless Communications and Mobile Computing*, Vol.
9 2021, No. 1, 2021.
- 10 17. Huang, J., W. Liu, and X. Sun, A Pavement Crack Detection Method Combining 2D with
11 3D Information Based on Dempster-Shafer Theory. *Computer-Aided Civil and Infrastructure
12 Engineering*, Vol. 29, No. 4, 2014, pp. 299–313.
- 13 18. Safaei, N., O. Smadi, A. Masoud, and B. Safaei, An Automatic Image Processing Al-
14 gorithm Based on Crack Pixel Density for Pavement Crack Detection and Classification.
15 *International Journal of Pavement Research and Technology*, Vol. 15, 2022, pp. 159–172.
- 16 19. Son Dong Nguyen, S. D., T. S. Tran, V. P. Tran, H. J. Lee, M. J. Piran, and V. P. Le, Deep
17 Learning-Based Crack Detection: A Survey. *International Journal of Pavement Research
18 and Technology*, Vol. 16, 2022, pp. 943–967.
- 19 20. Chang, J., Y. Liu, Z. Shu, H. Zhang, and H. Cao, Pavement Crack Identification Based
20 on Deep Learning and Denoising Model. In *2021 4th International Symposium on Traffic
21 Transportation and Civil Architecture (ISTTCA)*, 2021, pp. 180–185.
- 22 21. Gao, Z., X. Zhao, M. Cao, Z. Li, K. Liu, and B. M. Chen, Synergizing Low Rank Repre-
23 sentation and Deep Learning for Automatic Pavement Crack Detection. *IEEE Transactions
24 on Intelligent Transportation Systems*, Vol. 24, No. 10, 2023, pp. 10676–10690.
- 25 22. Xiang, X., Z. Wang, and Y. Qiao, An improved YOLOv5 crack detection method combined
26 with transformer. *IEEE Sensors Journal*, Vol. 22, No. 14, 2022, pp. 14328–14335.
- 27 23. Wu, P., J. Wu, and L. Xie, Pavement distress detection based on improved feature fusion
28 network. *Measurement*, Vol. 236, 2024, p. 115119.
- 29 24. AASHTO, Standard Specification for File Format of Two-Dimensional and Three-
30 Dimensional (2D/3D) Pavement Image Data. Washington, D.C., 2023.
- 31 25. Jocher, G., A. Chaurasia, and J. Qiu, *Ultralytics YOLO*, 2023.
- 32 26. David Heyse, Nicholas Warren, and Jelena Tešić, Identifying maritime vessels at multiple
33 levels of descriptions using deep features. In *Artificial Intelligence and Machine Learning
34 for Multi-Domain Operations Applications*, SPIE, 2019, Vol. 11006, pp. 423 – 431.
- 35 27. Warren, N., B. Garrard, S. E., and Tešić, J., Transfer Learning of Deep Neural Networks
36 for Visual Collaborative Maritime Asset Identification. In *2018 IEEE 4th International
37 Conference on Collaboration and Internet Computing (CIC)*, 2018, pp. 246–255.