# Predicting Damage of Dutch Road Markings

Amanda Jansen<sup>1,2</sup>, Vivienne Jansen<sup>1,2</sup>, Jurjen Helmus<sup>1</sup>, and Arnoud Visser<sup>2</sup>

Abstract. Road markings play a crucial role in road safety by guiding traffic and ensuring visibility. As markings deteriorate over time, their effectiveness diminishes, necessitating timely maintenance. This paper studies two methods to classify road-marking damage from recorded images, in accordance with the Dutch CROW guidelines. The first is a model based approach, which first uses a regression model to estimate the marking damage, and then applies the thresholds in the CROW guidelines to classify the damage class. In contrast, a data-driven approach is used, classifying directly the damage class with a YOLOv8 classifier. The data-driven approach achieves an F1-score of 0.97 for the binary-classification task and 0.75 for the multiclass classification task. Compared to other international studies, this is a competitive result.

**Keywords:** computer vision  $\cdot$  segmentation  $\cdot$  classification

### 1 Introduction

As cities grow and mobility increases, the pressure on public infrastructure and the need for efficient maintenance strategies intensifies [1]. Road markings, which include painted lines, symbols, and patterns on the road surface, play a key role in managing road safety. These markings help warn road users and ensure smooth traffic flow [17]. However, they degrade over time and currently rely on manual inspections that are time-consuming, costly, and often inconsistent [18].

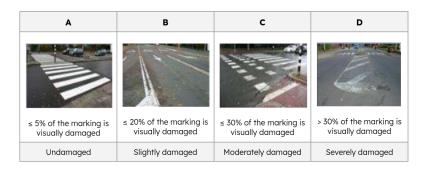


Fig. 1: CROW guidelines per severity category (classes A, B, C and D) [7].

Maintenance Lab, Applied University of Amsterdam, Amsterdam, The Netherlands
 Intelligent Robotics Lab, University of Amsterdam, Amsterdam, The Netherlands

<sup>\*</sup> Corresponding author: a.visser@uva.nl

Therefore, Velotech<sup>3</sup>, in collaboration with Amsterdam University of Applied Sciences, has developed a Smart Bikes project that uses artificial intelligence (AI) to automate the inspection of urban infrastructure. By equipping bikes with edge AI, a technology that processes data locally on the device, the system captures and analyzes road condition data in real-time. This reduces latency, minimizes dependence on external servers and aligns with municipal privacy standards, as sensitive visual data never leaves the bike [10]. This integrated approach allows municipalities to efficiently assess the condition of road markings and prioritize repairs.

In the Netherlands, such maintenance decisions are guided by CROW guidelines, which serve as the national standard for evaluating road infrastructure. These guidelines categorize road markings into four classes, from A to D, based on the damage severity [7]. Class A represents road markings in excellent condition, while class D indicates markings in poor condition that require repair (see Fig. 1). To automate this classification process, this research uses both a model-based and a data-driven approach. For the model-based approach, only binary results that indicate whether a road marking is damaged (class B, C or D) or undamaged (class A) are presented, whereas for the data-driven method, both binary and multiclass results are included.

# 2 Related Research

Datasets focused on road marking damage are scarce. A study of Iparraguirre et al. [4] combined two datasets from Japan and Spain. They added 971 new labeled images for Spanish roads. They performed binary classification with three different convolutional neural networks. Their best result used EfficientDet v1 D0 [13], and achieved an F1-score of 0.93, which was a large improvement compared to the previous result on the Japanese dataset (F1-score of 0.72) [9]. This study reports comparible results for the binary classification (see Sec. 4.2).

To estimate the level of damage, annotated data of the severity of the damage is needed. This information was available in a dataset from the USA. Recent work of Antariska et al. used a data-driven approach based on YOLOv8 [2]. This method was trained on 865 images collected along New Jersey State routes. The dataset concentrated on a subset of road markings used in this study, namely the center line. Instead of four damage classes, three damage classes were used (good, moderate, poor). This system achieved a macro-averaged precision of 0.51, considerably lower than the results reported in this study (see Sec. 4.2). Partly this can be contributed to the smaller dataset (they only annotated 865 images from the 15.536 available images). Their study was also limited by the resolution of the images collected along the New Jersey routes. The road markings cover only part of the image, so you have to zoom in at the road marking and make (implicit or explicit) an estimate of the damage ratio. In that case high-resolution images, as provided by the ZED-X stereo camera used in this study, can make the difference.

<sup>3</sup> https://velotech.ai/

### 3 Method

Estimation of the road-marking damage is performed in this study in two different ways; first by a regression method which estimates the amount of damaged paint followed by a decision-model based on the thresholds in the CROW guidelines. Second, because the regression model tends to overpredict damage severity (see Sec. 4.1), this method is compared with a fully data-driven approach.

#### 3.1 Dataset

The data that was used in this study was collected by Velotech with a ZED-X stereo camera that was developed by Stereolabs<sup>4</sup>. The camera was mounted on both bicycles and cars to simulate real-world mobile inspection scenarios [8]. Data acquisition took place in two distinct urban environments in the Netherlands: Geertruidenberg, a small municipality characterized by relatively calm residential streets, and Amsterdam Oud-Zuid, a densely populated urban district with a high volume of traffic, varied infrastructure, and complex road markings. This geographical variation introduces diversity in lighting conditions, road surface materials, marking styles, and levels of wear, and ensures that the dataset reflects a broad range of real-world conditions that are relevant to road marking assessment (see Fig. 2). The dataset is private but may be made available by Velotech upon reasonable request.

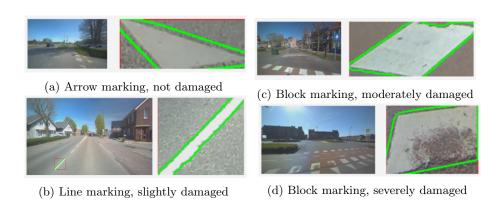


Fig. 2: Examples Velotech recordings and the individual marking selected for classification.

The combined recordings from Geertruidenberg and Amsterdam Oud-Zuid resulted in a dataset consisting of 15.745 high-resolution images of road markings. The images were each accompanied by annotations that were stored in a structured JSON file. Each entry described an individual road marking and included a polygon delineating its shape, a bounding box, a manually assigned

<sup>4</sup> https://www.stereolabs.com/

CROW severity class (A, B, C, or D), and a 'superclass' indicating the broader category of the road marking (e.g., line, arrow, or zebra crossing). Polygon and bounding box annotations were automatically generated using a YOLOv8 segmentation model developed by Velotech [12]. This segmentation model was trained on a smaller dataset, containing 6.648 annotated road markings, recorded in Amsterdam, Haarlem, Tilburg and Diemen. After removing images with missing severity or polygon annotations, the final dataset used in this study contained 15.723 annotated road markings.

The majority of markings fall into severity classes A (#5.002) and B (#6.572), while fewer instances are observed in classes C (#1.907) and D (#2.242). The dataset is strongly dominated by the line category (#13.720), with considerably fewer examples of other types such as giveawayrow (#807), block (#368), zebrawalking (#314), and smaller categories such as arrow (#63) and stopping (#48) marking. Those markings were manually annotated by multiple individuals. Velotech has implemented quality control measures to get consistent annotations between the annotators.

Because the distribution is imbalanced, care had to be taken for the data-driven approach to ensure that this imbalance did not affect the training process. For this reason, the dataset was divided into training (70%), test (15%), and validation (15%) subsets using stratified splitting.

# 3.2 Regression Method

The regression method starts with a high-resolution image recorded by Velotech. Velotech has an accurate YOLOv8-detection model (98% pixel accuracy) [12] which localizes road markings and generates a polygon and bounding box around the road marking. The polygon is used here to get a binary mask. The polygon is not intended to be tight around the road marking, so part of the road surface is visible at the edges. To prevent that these edges contribute to estimation of the amount of damaged paint, an erosion algorithm with a kernel of  $7 \times 7$  is applied to be able to concentrate on the core of the road marking (see Fig. 3).

An estimate of the color of the road surface (in grayscale) surrounding the road marker is also important, because when the road marking is damaged the road surface shines through. Yet, the color of the road surface is not always the same, nor the lighting conditions, so an outlier mask (see Fig. 3e) can be used to estimate the color of the road surface near the road marking.

To distinguish between intact and damaged areas within the marking a dynamic threshold is applied, based on Otsu's method [11]. Otsu's method used the histogram of the image to define two clusters of bright and dark pixels as classes and maximizes the between-class variance to find the optimal threshold. By calculating the proportion of dark/damaged pixels inside the eroded mask the damage ratio can be calculated, which can directly be mapped to the CROW guidelines.

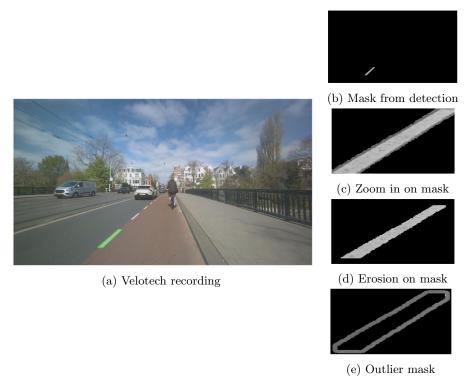


Fig. 3: Preprocessing method

#### 3.3 Data-Driven Method

The road markings are already detected and localized in the images with a YOLOv8-based detector [12], so it is logical to see how well a YOLOv8-based classification would work on this problem. You Only Look Once (YOLO) is a real-time object detection algorithm known for its speed and accuracy in identifying and classifying visual elements within images [14]. YOLOv8 is an algorithm that is slightly easier to fine-tune on new types of objects than YOLOv9 [16]. An alternative would be the most recent YOLOv12 [15], although the attention model relies on FlashAttention for optimal speed. FlashAttention is only supported on relatively modern GPU architectures and is less suitable for edge-computing.

Two models were trained with YOLOv8; one for a binary (damaged/undamaged) classification task and one for a multiclass (A/B/C/D) classification task. The models were initialized with pre-trained weights and trained on road-marking images. These images were obtained by cropping the original images to the bounding boxes of the road markings. Training was conducted for 100 epochs with an input image resolution of  $640\times640$  pixels, a batch size of 32, and 8 data-loader workers. Early stopping was applied with a patience of 10

epochs to prevent overfitting. Only the default data-augmentation settings from YOLOv8 were used.

To evaluate the learning behavior and generalization capability of both models during training, the progression of the training and validation loss was monitored. The trained YOLOv8 models were both evaluated on the 15% validation set, consisting of 2.359 previously unseen road marking instances, each labeled with one of the four CROW-defined severity classes (A-D). Fig. 4 shows the loss curves for the binary classification model, and Fig. 5 presents those for the multiclass model.

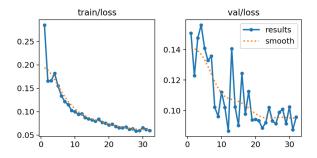


Fig. 4: Training and validation loss functions for the YOLOv8 binary (damaged and undamaged) model classification model. Training loss (left) and validation loss (right) curves over 32 training epochs.

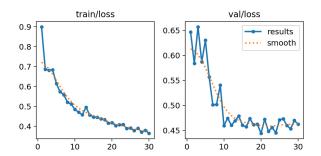


Fig. 5: Training and validation loss functions for the YOLOv8 multiclass (classes A, B, C and D) model. Training loss (left) and validation loss (right) curves over 30 training epochs

For both models, the training loss steadily decreases, demonstrating that the model is effectively learning to minimize the error on the training dataset. Simultaneously, the validation loss shows a similar downward trend and closely tracks the training loss, indicating that the model generalizes well to unseen data. The absence of any increase or divergence in validation loss gives confidence that overfitting did not occur (yet).

### 4 Results

Both the regression as well as the data-driven method were studied extensively, including multiclass classification with the regression method and analysis of the model performance per road marking type. For more details, see the theses [5,6].

### 4.1 Regression Results

When the regression method described in Sec. 3.2 is applied to a binary classification task, distinguishing between undamaged (class A) and damaged (class B,C,D) markings, the precision for the undamaged class (0.90) and the recall of the damaged class (0.99) are quite good. So, from the 2.356 instances in the testset (one instance is removed in this testset because of a missing polygon annotation), 99% of the damaged markings were correctly identified and when a road marking was predicted to be undamaged, it is highly likely to be correct.

Class	Precision	Recall	F1-score	Instances
Undamaged (class A)	0.90	0.26	0.40	749
Damaged (class $B/C/D$ )	0.74	0.99	0.85	1607
Macro-average	0.82	0.62	0.62	2356
Accuracy		_	0.75	2356

Yet, as can be seen from Table 1, the model tends to over-detect damage and frequently misclassifies undamaged markings as damaged. The undamaged recall is only 0.26 and the precision on damaged markings is only 0.74, leading to a macro-average F1-score of 0.62. This makes this model only useful as pre-filtering tool to reduce the volume of markings requiring manual inspection within the maintenance workflow of Velotech. Since these binary-classification results are not promising enough for multiclass classification, they are not further discussed in this study. Nevertheless, preliminary multiclass-classification results are available [5].

# 4.2 Data-Driven Results

The data-driven approach described in Sec. 3.3 improves these results both for the undamaged recall and the precision on recognizing damaged markings (Table 2). Although there is still a slight bias towards flagging damage, the false negatives are so low that it approaches the operational goal of Velotech.

Table 2: Binary Classification with YOLOv8 (test set)

Class	Precision	Recall	F1-score	Instances
Undamaged (class A)	0.98	0.95	0.96	749
Damaged (class $B/C/D$ )	0.97	0.99	0.98	1608
Macro-average	0.97	0.97	0.97	2357
Accuracy		_	0.97	2357

Based on these preliminary numbers, it becomes interesting to look at the multiclass results, if the data-driven approach makes the distinction between slightly damaged, moderately damaged, severely damaged (CROW classes B/C/D).

Table 3: Multiclass Classification with YOLOv8 (test set)

				`
Class	Precision	Recall	F1-score	Instances
A	0.95	0.95	0.95	749
В	0.85	0.89	0.87	986
С	0.54	0.40	0.46	288
D	0.70	0.76	0.72	334
Macro-average	0.76	0.75	0.75	2357
Accuracy	_	_	0.83	2357

The results for road markings in good condition or only slightly deteriorated (class A/B) are good, as can be seen in Table 3. In contrast, the performance is notably lower for moderately or severely damaged road markings (class C/D).

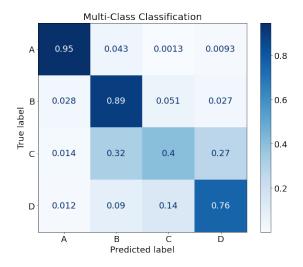


Fig. 6: Confusion matrix of overall severity classification on the test set. The confusion matrix displays the number of predicted labels versus ground truth labels across the four CROW severity classes. The diagonal cells represent correct predictions.

Most errors originate from class C. As can be seen from the confusion matrix in Fig. 6, where road markings annotated as 'moderately damaged' are predicted in 32% of the cases as 'slightly damaged', 40% of the cases as 'moderately damaged' and in 27% of the cases as 'severely damaged'. The boundaries between class B/C (20% damaged surface) and between class C/D (30% damaged surface) seem to be hard to estimate. Human annotators also have difficulty making this distinction, as can be seen in Sec. 5.1.

Classes C and D have the fewest instances in the dataset, which could potentially lead to class imbalance effects favoring classes A and B. However, classes C and D are roughly the same size, yet the model performs much better on class D than on class C. This indicates that the number of examples for class C was sufficient, and the relatively low performance for class C is therefore not solely due to a lack of training data, but rather due to the inherent difficulty in distinguishing this class from its neighbors.

# 5 Discussion

# 5.1 Damage Severity

The manual annotations of the severity classes were outsourced by Velotech and carried out by multiple individuals. Velotech has implemented quality control measures to get consistent annotations between the annotators. Still, differentiating between severity levels B ( $\leq$  20% damage), C ( $\leq$  30% damage) and D (> 30% damage) requires annotators to estimate the proportion of damaged surface area by eye, which is a task prone to individual interpretation and variability.

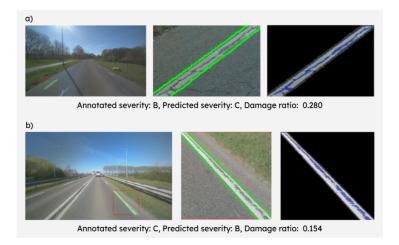


Fig. 7: Examples in which the model provides a reasonable damage ratio other than the severity classification from the manual annotations.

Consequently, even experienced annotators may produce inconsistent or inaccurate labels, particularly when the extent of damage is near the boundary between two severity classes. In cases where the model correctly detects damage to a road marking, the damage ratio given by the regression model appears to provide a more objective and consistent assessment of damage severity compared to manual evaluations of the severity class (see Fig. 7).

The regression model does not consistently detect damage correctly. Especially damage class C was hard to classify, as shown in Section 4.2. Two common failure cases can be easily demonstrated here with two examples. The first failure case is due to partial shadows. Because the algorithm depends on brightness, a sharp shadow can be easily judged to be a damaged area (see Fig. 8).

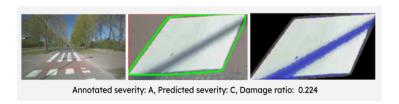


Fig. 8: Example of a road marking partially covered by shadow. The marking is labeled as class A (undamaged) in the ground truth, but the model predicts it as class C due to the presence of a partial shadow cast by a streetlight. The reduced brightness in the shadowed region is incorrectly interpreted by the model as surface damage.

Another common failure case originates from complex road marking shapes. For some road marking symbols it is difficult to clearly define the inside and outside regions. When the mask is just a square bounding box, a lot of the road surface is still visible, which could result in unrealistic high damage ratio estimates (see Fig. 9).



Fig. 9: Examples of highly imprecise polygon annotations that resemble bounding boxes rather than accurate segmentations, shown for the bicycle symbol.

Yet, with many borderline cases where class boundaries are open for interpretation, it is hard to train a data-driven approach like YOLO on learning the decision boundaries more precisely. Because the boundaries between the four

categories are hard to distinguish, both for humans and data-driven approaches, it could be interesting to try a fuzzy-logic based approach on this application, as demonstrated in [3].

#### 5.2 Information Leakage

Although care has been taken to prevent information leakage inside this study, two sources of information leakage still remain. A first source of leakage is in the pre-processing; some of the images used to train the segmentation model (the ones recorded in Amsterdam), are also used to train the classification model.

A second source of information leakage is the stratified splitting of the dataset in a train/test/validation set. Because images are recorded at different places on the same road, a sample from the same road could be present in both the train/test/validation set.

### 6 Conclusion

The model-based approach showed mixed results for different types of road-markings. Simple road markings such as lines and blocks with clear boundaries could be classified quite well, but for complex road markings such as bicycle symbols, a pixel-wise segmentation might be required. In addition, this method was sensitive to lighting conditions, such as the occurrence of partial shadows on the road markings.

In contrast, the data-driven approach worked well under all circumstances, although it suffered from class imbalance in the training dataset. It is a difficult classification task because the thresholds defined by the CROW guidelines are quite strict, so the YOLOv8 classification could easily confuse the 'moderately damaged' class with the 'slightly damaged' or 'severely damaged' classes. Human annotators also had difficulty making this distinction.

In conclusion, manual verification of the damage is still required. Automatic classification could still benefit maintenance operations, by excluding clearly undamaged road markings and allowing them to give priority to severely damaged road markings.

#### Conflict of Interest Statement

Considering the involvement of a commercial partner, Velotech, we disclose that the company provided data and YOLOv8 segmentation model results. The research was conducted independently, and Velotech had no influence over the analysis, interpretation, or reporting of the results.

### Acknowledgement

We would like to express our sincere gratitude to Velotech for their valuable collaboration and for providing us with access to their data. Their support and the availability of their dataset have been fundamental to this research.

### References

- Ai, D., Jiang, G., Lam, S.K., He, P., Li, C.: Computer vision framework for crack detection of civil infrastructure—a review. Engineering Applications of Artificial Intelligence 117, 105478 (Jan 2023)
- Antariksa, G., Chakraborty, R., Somvanshi, S., Das, S., Jalayer, M., Patel, D.R., Mills, D.: Comparative analysis of advanced ai-based object detection models for pavement marking quality assessment during daytime. preprint arXiv:2503.11008 (Mar 2025)
- Eric Manongga, W., Chen, R.C.: Road marking sign damage detection and classification using deep learning and fuzzy method. IEEE Access 13, 92943–92952 (May 2025)
- Iparraguirre, O., Iturbe-Olleta, N., Brazalez, A., Borro, D.: Road marking damage detection based on deep learning for infrastructure evaluation in emerging autonomous driving. IEEE Transactions on Intelligent Transportation Systems 23(11), 22378–22385 (Jul 2022)
- 5. Jansen, A.: Classifying the damage severity of road markings using computer vision. Bachelor thesis, University of Amsterdam (Jul 2025)
- Jansen, V.: Automated assessment of the damage severity of road markings using yolo. Bachelor thesis, University of Amsterdam (Jul 2025)
- 7. Kennisplatform CROW: Kwaliteitscatalogus openbare ruimte 2023. published online (Aug 2023)
- 8. Koomen, D.: Determining the geolocation of road markings based on stereo camera vision and a Bird's Eye View. Bachelor thesis, University of Amsterdam (Jun 2024)
- Maeda, H., Sekimoto, Y., Seto, T., Kashiyama, T., Omata, H.: Road damage detection and classification using deep neural networks with smartphone images. Computer-Aided Civil and Infrastructure Engineering 33(12), 1127–1141 (Jun 2018)
- 10. Meuser, T., et al.: Revisiting edge ai: Opportunities and challenges. IEEE Internet Computing **28**(4), 49–59 (Aug 2024)
- 11. Otsu, N.: A threshold selection method from gray-level histograms. IEEE Transactions on Systems, Man, and Cybernetics 9(1), 62–66 (Jan 1979)
- Schimmelpennink, M.: Segmentation of road markings in the Netherlands for automatic inspections, using YOLOv8-seg on an edge system. Bachelor thesis, University of Amsterdam (Jul 2024)
- Tan, M., Pang, R., Le, Q.V.: Efficientdet: Scalable and efficient object detection.
  In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10778–10787 (Aug 2020)
- Terven, J., Córdova-Esparza, D.M., Romero-González, J.A.: A comprehensive review of yolo architectures in computer vision. Machine Learning and Knowledge Extraction 5(4), 1680–1716 (Nov 2023)
- 15. Tian, Y., Ye, Q., Doermann, D.: Yolov12: Attention-centric real-time object detectors. preprint arXiv 2502.12524 (Feb 2025)
- Wang, C.Y., Yeh, I.H., Mark Liao, H.Y.: Yolov9: Learning what you want to learn using programmable gradient information. In: European conference on computer vision. pp. 1–21. Springer (Oct 2024)
- Wu, J., Liu, W., Maruyama, Y.: Street view image-based road marking inspection system using computer vision and deep learning techniques. Sensors 24(23), 7724 (Nov 2024)
- 18. Younesi Heravi, M., Dola, I.S., Jang, Y., Jeong, I.: Edge ai-enabled road fixture monitoring system. Buildings 14(5) (Apr 2024)