# Pavementscapes: a large-scale hierarchical image dataset for asphalt pavement damage segmentation[1]

Zheng Tong[a], Tao Ma[a,*], Ju Huyan[a], Weiguang Zhang[a]

[a]*School of Transportation, Southeast University, Nanjing 211189, China.*

## Abstract

Pavement damage segmentation has benefited enormously from deep learning. However, few current public datasets limit the potential exploration of deep learning in the application of pavement damage segmentation. To address this problem, this study has proposed Pavementscapes, a large-scale dataset to develop and evaluate methods for pavement damage segmentation. Pavementscapes is comprised of 4,000 images with a resolution of $1024 \times 2048$, which have been recorded in the real-world pavement inspection projects with 15 different pavements. A total of 8,680 damage instances are manually labeled with six damage classes at the pixel level. The statistical study gives a thorough investigation and analysis of the proposed dataset. The numeral experiments propose the top-performing deep neural networks capable of segmenting pavement damages, which provides the baselines of the open challenge for pavement inspection. The experiment results also indicate the existing problems for damage segmentation using deep learning, and this study provides potential solutions.

*Keywords:* pavement damage dataset, supervised learning, damage segmentation, deep learning, convolutional neural network, attention-based network

## 1. Introduction

Performance of pavement surface decays over time due to various factors, such as traffic volume and weather, and, therefore, the understanding of the deterioration degree is essential to efficient maintenance, which aims to keep and improve the high performance of the surface. Pavement damage, a key characteristic of road surface deterioration degree, is evaluated with three approaches: manual, semi-automatic, and fully automatic.

In the manual approach, investigators walk or slowly drive along the pavement to inspect damages, which is subjective and time-consuming. In the semi-automated one, a fast-moving vehicle automatically collects pavement surface images and an off-line and manual

---

[1]The Pavementscapes dataset is available at `https://github.com/tongzheng1992/Pavementscapes`.
[*]Corresponding author

process of damage inspection is then performed in workstations, which is still very time-consuming. Compared to the semi-automated one, the fully automatic approach adopts some technologies of computer vision to perform the inspection using road surface images, which brings the potential of real-time data processing with a low labor cost.

In the fully automatic approach, the application of computer vision on pavement visual inspection involves three main tasks. The first one is *damage recognition* [3, 28, 61], also known as *damage classification*, in which an algorithm should indicate the category of each damage present in a 2D or 3D pavement image. Another task is *damage detection* [4, 43]: bounding boxes identify and locate one or more effective damages in a pavement image. The last one is *damage semantic segmentation* [55, 56] which splits an image into multiple sets of pixels and each set has its individual damage category. These splits, called the *segmentation mask*, are regarded as a simplified but informative representation of the original image. For example, the boundary areas of the segments in a mask provide the position information of pavement damages in an image. Segmentation results can also be used to measure the damage morphology, such as the width and length of a crack. Besides, compared to the results of damage recognition and detection, the ones from damage segmentation are more informative and useful.

The explosive development of deep neural networks [26] brings a large number of outstanding and state-of-the-art methods for semantic segmentation. Many top-performing methods are nowadays extended into the application of pavement damage segmentation and have achieved remarkable success [3, 4, 49, 55, 63, 64]. The predominant reason for the success is the availability of the large-scale and public datasets, such as ImageNet [9], PASCAL VOC [13], Cityscapes [7], and Microsoft COCO [30]. These datasets exploit the powerful capacity of deep neural networks. Unfortunately, there are a small number of public datasets for pavement damage segmentation [44]. This condition has limited the potential exploration of deep learning in the application of pavement damage segmentation. Besides, the lack of publicly available datasets makes the existing algorithms incomparable since they are reported as the state-of-the-art ones only in their own datasets.

This study introduces a large-scale dataset, called *Pavementscapes*, to solve the above-mentioned issue of pavement damage segmentation. The contributions of this study can be summarized as follows:

1. Propose a large-scale hierarchical image dataset for asphalt pavement damage segmentation. This dataset can be used to train and test the approaches for pavement damage segmentation, especially for deep neural networks. The proposed dataset consists of 4,000 real-world pavement images annotated in image, block, and pixel levels. For each level, there are six categories of visual pavement damages. The *image-level annotations* defines the damage categories shown in each image. The *block-level annotations* use bounding boxes to identify and locate each damage. In the *pixel-level annotations*, each pixel of a pavement image is labeled as one of the damage categories or background. Besides, all images have a view of top-down shooting, which allows users to measure the accurate morphology of pavement damages using segmentation masks.

2. Compare the top-performing segmentation methods of computer vision on the proposed dataset. This study investigates the state-of-the-art deep-learning algorithms based on the proposed dataset, which provides the baselines of an open challenge for pavement damage segmentation. The comparison study also indicates the existing issues and potential solutions for the segmentation task, including small damage segmentation, unbalanced training set, and over-confidence in modern neural networks.

3. Record pavement damages with non-iconic views, which are also known as non-canonical perspectives [38]). An instance has an iconic view if it is near the center of a digital image. Despite the existing gap in human performance, current algorithms segment damage fairly well on iconic views but struggle to do it otherwise – in the partially occluded and amid clutter [18], reflecting the complexity of real-world inspection projects. An ideal model should also perform well in the non-iconic views. Thus, the proposed dataset includes numerous pavement damages in non-iconic views.

The rest of the paper is organized as follows. Section 2 begins with the literature review of the public pavement datasets and state-of-the-art segmentation models based on deep learning. Section 3 describes the details of the proposed datasets, followed by the damage segmentation experiments in Section 4. Finally, Section 5 concludes this study.

## 2. Related work

### 2.1. Pavement damage datasets

Recent studies have adopted various machine learning algorithms, especially deep learning for automatic pavement damage segmentation, such as [42, 64, 65]. To develop these algorithms, several datasets of pavement images have been released and the majority of them are summarized in Table 1. Three problems can be found in these public datasets. The first one is the views of the pavement images. Most of these datasets captured the images using smartphones with wide views, such as the most successful one RDD-2020 [2]. However, many real-world projects of pavement inspection require a top-down view because pavement maintenance needs information on damage morphology, such as the dimensions of cracks and portholes. A wide view cannot accurately provide morphological information since the damage morphology is distorted in the view. In addition, many datasets present damages with an iconic view, appearing objects in a profile unobstructed near the center of a neatly composed image. However, pavement damages with non-iconic views, such as the incomplete and occluded cracks, are common in real-world inspection projects.

Another problem is the annotations. Table 1 indicates that the existing datasets only annotate one or two categories of pavement damages, e.g., crack and pothole. However, asphalt pavement inspection requires several visual damage categories, such as eight in [24]. Besides, many datasets only provide image- and block-level annotations, which cannot be used to train and test segmentation models. A few datasets include a small number of images with pixel-level annotations, which still cannot meet the requirement of developing a deep neural network for damage segmentation.

The last problem is the baseline algorithms on these datasets. Many datasets still use the results of machine learning proposed thirty years ago as their baselines, even though a few

Table 1: Summary of the existing pavement damage datasets. The algorithms in bold are the deep neural networks. "Pavementscape" in the last row is the dataset proposed in this study.

| Dataset | Images | Resolution | Data collection device | View | | Damage category | | | | Annotation level | | | Baseline algorithms |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Top-down | Wide-view | Crack | Pothole | Rut | Repair | Image | Block | Pixel | |
| Ouma and Hahn [37] | 75 | 1080 × 1920 | Galaxy S5 cammer | ✓ | | | ✓ | | | | | ✓ | Fuzzy c-means |
| CrackIT [36] | 84 | 1536 × 2048 | Optical device | ✓ | | ✓ | | | | ✓ | ✓ | | K-nearest neighbor |
| CFD [41] | 118 | 480 × 320 | Iphone 5 | ✓ | | ✓ | | | | | | ✓ | Random decision forests |
| CrackTree200 [67] | 206 | 800 × 600 | Area-arry camera | ✓ | | ✓ | | | | | | ✓ | Minimum spanning tree |
| SDNET2018 [10] | 230 | 4068 × 3456 | 16MP Nikon Digital Camera | ✓ | | ✓ | | | | ✓ | | | **AlexNet** |
| Crack500 [60] | 500 | 2000 × 1500 | GoPro 7 | ✓ | | ✓ | | | | | | ✓ | **FPHBN** |
| CrackDataset [21] (segmentation part) | 1205 | 1280 × 960 | Action camera | ✓ | | ✓ | | | | | ✓ | ✓ | **U-net** |
| GAPs v1 [12] | 1,969 | 1920 × 1080 | Professional camera | | ✓ | ✓ | ✓ | | ✓ | ✓ | | | **ASINVOS net** |
| GAPs v2 [45] | 2,468 | 1920 × 1080 | Professional camera | | ✓ | ✓ | ✓ | | ✓ | | ✓ | | **ASINVOS net** and **ResNet34** |
| RDD-2020 [2] | 26,620 | 720 × 960 | LG Nexus 5X cameras | | ✓ | ✓ | ✓ | | | | ✓ | | **MobileNet** |
| Pavementscapes | 4,000 | 2048 × 1024 | Professional camera | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **Tables 4** |

novel deep neural networks have been adopted in some datasets. This phenomenon does not explore the potential of deep learning in the application of pavement damage segmentation.

The three problems motivate us to release a new large-scale hierarchical image dataset for asphalt pavement damage segmentation and analyze the proposed dataset with state-of-the-art deep neural networks.

## 2.2. Deep neural networks for pavement damage segmentation

After the success of Long et al. [33] on semantic segmentation, a large number of deep neural networks have been proposed and achieved the state-of-the-art results. Generally, there are two main directions: convolution- and attention-based deep neural networks.

*Convolution-based deep neural network*, also known as convolutional neural network (CNN), is a neural network that uses convolution in place of general matrix multiplication in at least one of their layers. For the segmentation problems, the most widely-used CNN architecture is the *fully convolution networks (FCN)* [33], as shown in Figure 1, which only consists of locally connected layers, e.g., convolution, pooling, and upsampling layers. No fully connected layer is utilized in this networks. An FCN extracts high-dimension features from a pavement image using convolution and pooling layers and the features are then upsampled into pixel-wise feature maps by upsampling layers. The upsampled feature maps are finally imported into a softmax layer to predict the classes of all pixels in the input image. Many FCN-based models have been used for pavement damage segmentation, such as FCN-8s [55], U-net [14, 32], W-Net [17], a series of DeepLab [5, 31]. These studies have demonstrated that the CNN models have significant superiority in pavement damage segmentation, once given enough learning samples with reasonable pixel-wise annotations.

Another direction is the *attention-based deep neural networks*. Compared to an FCN directly using a full image for segmentation, an attention-based network first splits an input into a square-patch grid. Each patch is vectorized by concatenating its channels of all elements and then linearly projected to the required size. After dividing the sample, the network is agnostic to the position information about these patch vectors. Thus, learnable position embeddings are linearly added to each vector, which allows it to learn about the relative or absolute positions of the patches. These embedded patch vectors are then sequentially imported into a transformer encoder. The encoder consists of alternating layers of self-attention and multi-layer perceptron. Self-attention of an embedded patch vector is
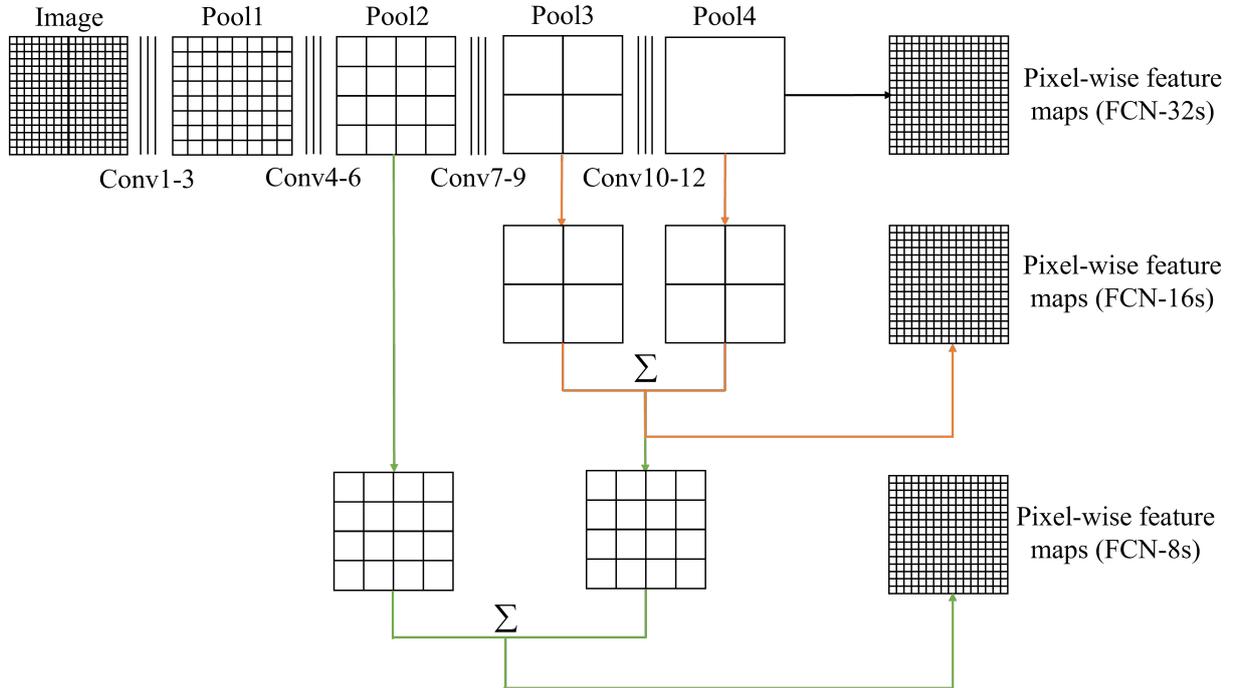
Figure 1: Overview of fully convolution networks [33].

defined as its relationship to every other vector. Feeding the embedded patch vectors sequentially, a self-attention layer computes their self-attentions as introduced in [11]. These self-attentions are then fed into a multi-layer perceptron layer to handle their dimension. The self-attention outputs of the final transformer encoder are concatenated and imported into a mask transformer. The outputs of the mask transformer, as the pixel-wise feature maps of the input image, are imported into a softmax layer for object segmentation. The processes can be summarized as Figure 2. Even though attention-based networks (e.g., transformer segmentor [46], attention U-net [35] and R2U-net [39]) have achieved remarkable results on the majority of benchmark datasets for semantic segmentation, there are only a few applications on pavement damage segmentation [48, 25, 48]. Therefore, this study should compare the performances of convolution- and attention-based deep neural networks on pavement damage segmentation, once given a new large-scale hierarchical pavement damage dataset.
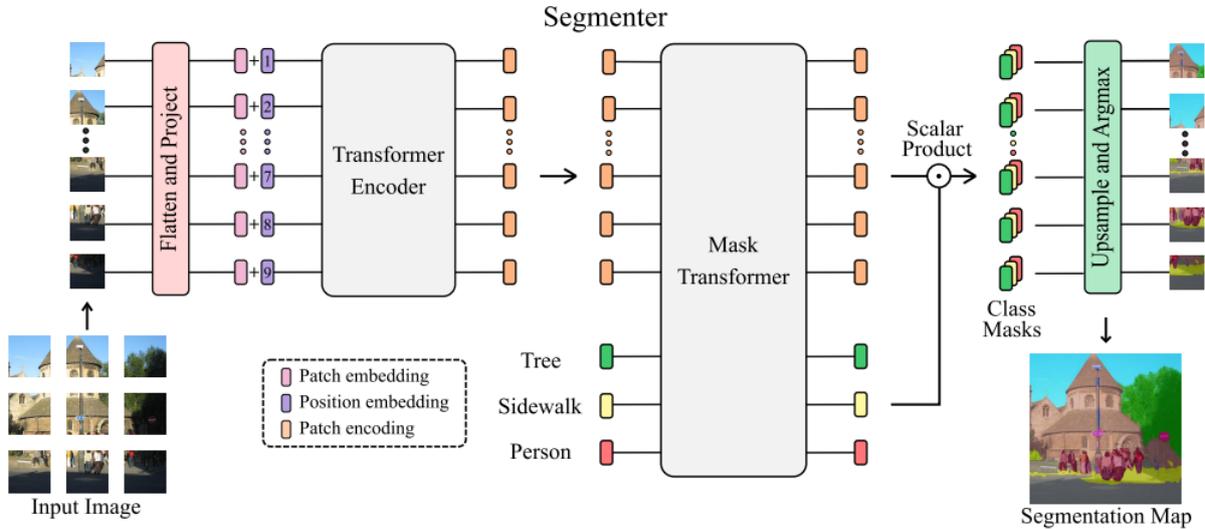
5

Figure 2: Overview of an attention-based deep neural network [46].

## 3. Proposed dataset

This section describes the proposed dataset, including image collection (Section 3.1), categories and annotations (Section 3.2), the protocol of dataset split (Section 3.3), and the statistical analysis (Section 3.4 ). The proposed dataset, named *Pavementscapes*, is comprised of 4k images with a resolution of $1024 \times 2048$, which were collected from 15 different pavements in China. A total of 8,680 damage instances in six categories are manually labeled in image, block, and pixel levels.

### 3.1. Image collection

In order to guarantee the comprehensiveness of the proposed dataset, the areas for image collection are made up of 15 pavements at different locations in China (Jiangxi, Gansu, Heilongjiang, and Xinjiang provinces), which have various service years (1-10 years), traffic volumes, weather, and surface materials (AC-13, AC-16, SMA-13, etc). Figure 3 presents the details of these pavements.

Pavement images were captured using a multi-function detection vehicle equipped with a professional camera, as shown in Figure 4. The camera was installed with a top-down shooting view of pavement surfaces, and it captured PNG images with a size of $1024 \times 2048 \times 1$ when the vehicle moved about 60 km/h. More than 500k images were gathered and 4k of them with at least one damage were used to make up the Pavementscapes dataset.

### 3.2. Categories and annotations

The Pavementscapes dataset consists of six damage categories in total, covering the majority of the visual damage categories in the Chinese Highway Performance Assessment Standards (JTG H20-2007) [24], as shown in Table 2.

The proposed dataset provides annotations at image, block, and pixel levels. Using Labelme library [58], these annotations were labeled in-house by five annotators with at least

6

Figure 3: Study area of the Pavementscape dataset in the Google map. The blue curves are the investigated pavement areas.



Figure 4: Multi-function detection vehicle equipped with a professional camera.

Table 2: Pavement damage types in the Pavementscapes dataset and [24].

| Damge type | | Included in the Pavementscapes dataset or not |
|---|---|:---:|
| Crack | Longitudinal | ✓ |
| | Lateral | ✓ |
| | Alligator | ✓ |
| Others | Pothole | ✓ |
| | Material loose | × |
| | Rut | ✓ |
| | Wave crowding | × |
| | Repair area | ✓ |

five-year experiments on pavement inspection to guarantee high quality. At the image level, annotators label the damage categories presented in each image, in which multi-class annotations have existed. At the block level, bounding boxes identify and localize damages, which are stored in Excel format. The position information of a bounding box is represented by its left-top and right-bottom coordinates in an image. At the pixel level, annotators assign each pixel to one of the six categories or "background" class by labeling these images. The "background" class has the semantic of "anything else" except the six damage categories. The pixel-level annotations are also restored in the PNG format. To guarantee the annotation quality, an annotator should take at least 10 min to label one pavement image. Figure 5 shows some pavement images and their pixel-level annotations. In practice, these images and their annotations should be transformed into a TFRecord format, which can save computer memory [1]. Even though three types of annotations have been provided, this study only focuses on pixel-level ones. This is because the pixel-level annotations can also represent the information of the image- and block-level ones. The boundary areas of the pixel-level annotations provide the position information of pavement damages in an image, while the pixel classes include the category information of pavement damages. Of course, users can only use the image- and block-level annotations if they only need to train and test a damage detection model.

*3.3. Protocol of dataset split*

The Pavementscapes dataset has been pre-split into separate training, validation, and test sets for any supervised algorithms of computer vision. The protocol of dataset split is not random, but rather in the principle that makes each split representative for various pavement surface scenarios. Specifically, each split set is made up of pavement images collected with the following properties: (i) with a real-world distribution of the damage numbers across individual categories; (ii) in the different geographic locations of China with completely different climate conditions; (iii) at the fine and poor sunshine; (iv) with the different service years. Following this scheme, this study designs a protocol of the Pavemenetscape dataset split with 2,500 training images, 500 validation images, and 1,000 testing images.

To evaluate how representative the three splits w.r.t the protocol properties, an FCN [33] was trained by the 600 images from the training set and then evaluated by the testing
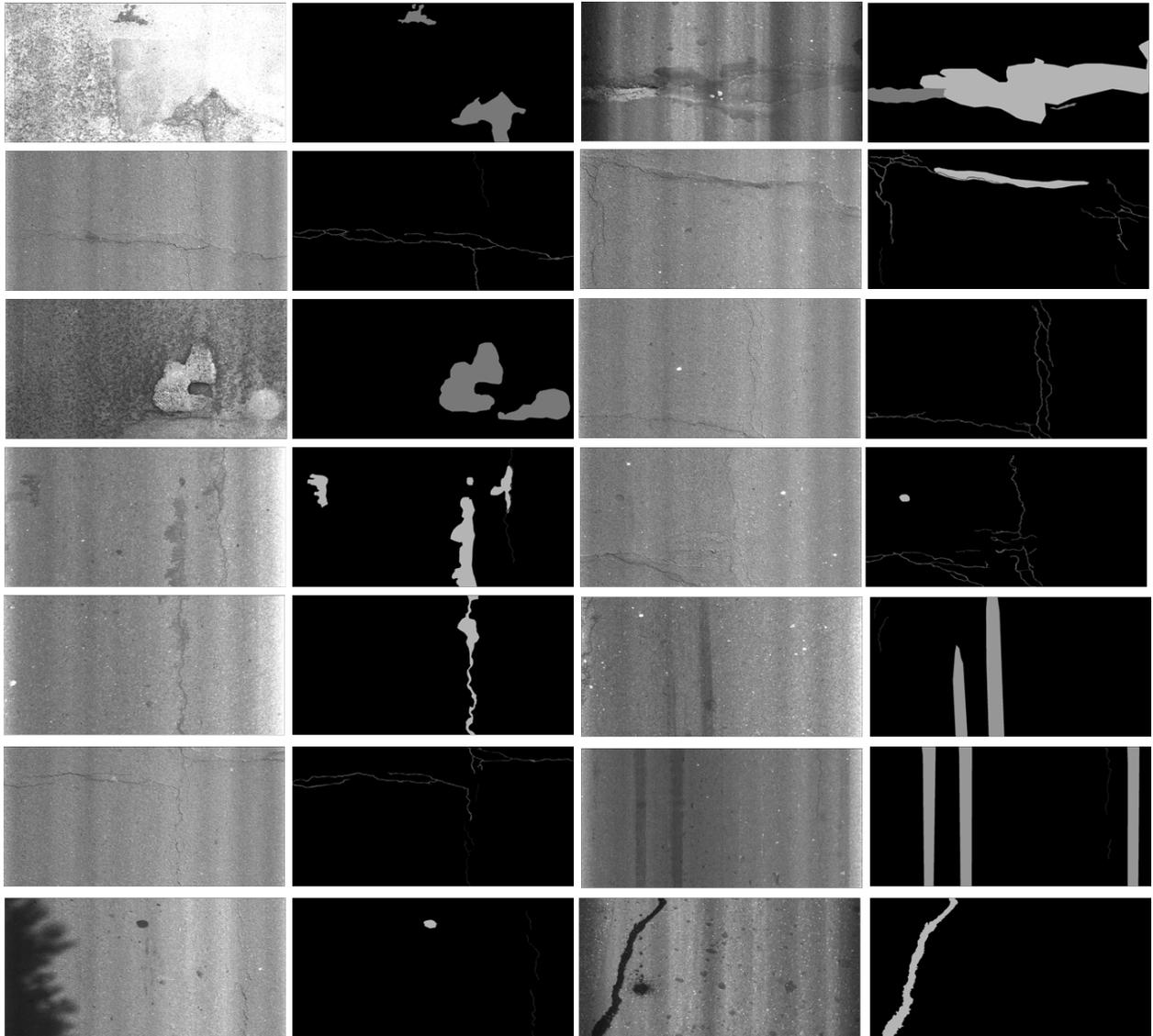
Figure 5: Examples of pavement images and pixel-level annotations from the Pavementscapes dataset. The masks with gray-scale values of 0 are the "background" pixels; other masks with different gray-scale values are the pixels belonging to different classes, such that 30, 60, 90, 120, 150, and 180 gray-scale values stand for the pixels of "longitudinal crack", "lateral crack", "alligator crack", "pothole", "rut", and "repair area", respectively.

set, as well as eight subsets of the testing set. For each subset, this study randomly selects three-eighths of the testing set. The accuracies of the testing set and its subsets are very uniform, varying less than 2.0%. Similar phenomena can also be found in the properties of geographic locations and service years. Interestingly, the performance on the fine sunshine is higher than the one on the whole test set. This is mainly because the images in lighting conditions represent damage features better than in the other conditions. To analyze this behavior in-depth, an additional test is performed by using images collected in low- or high-sunshine conditions, observing a 3.2% accuracy decrease in the low one and a 1.1% increase in the high one. Similarly, extreme training samples for one condition, such as all training images collected from the SMA pavement, improve the performance on the special testing samples but decrease the one on the whole testing set. These results highlight the comprehensiveness of the proposed dataset that should cover the majority of pavement surface scenes in the real world.

### 3.4. Statistical analysis

This section provides a statistical analysis of the proposed datasets, including (i) the distribution of visual damages, (ii) scene complexity under various real-world conditions, (iii) annotation accuracy, and (iv) non-iconic views. Regarding the first aspect, we compare the Pavementscapes dataset to Crack500 [60] and CrackDataset [21] having pixel-level annotations. Note that many other pavement damage datasets only have the image- and block-level annotations, such as the ones in Table 1. However, this study restricts this part of the analysis to those with a focus on damage segmentation because any pixel-level annotations can be easily converted into the image- and block-level annotations but the image- and block-level annotations cannot be transformed into the pixel-level ones.

*Damage distribution.* Figure 6a presents the numbers of damages across individual classes on the Pavementscapes dataset. In terms of the overall composition, the distribution of different damages is not uniform but close to the distribution of these damages in the real world, which makes deep neural networks easy to learn different damage features. Figure 6b presents the damage distributions in the training, validation, and testing sets. The distributions on the three sets are similar to the distribution of the proposed dataset. The ratio of the damage numbers in the training, validation, and testing sets is about 6:1:2.

Figure 6a also compares the Pavementscapes dataset with the Crack500 and Crack-Dataset datasets. The proposed dataset exceeds the other two datasets in the inherently different configurations. The Pavementscapes dataset involves different pavement surface damages in wide roads (at least one lane with 3.75 m), whereas the Crack500 and CrackDataset datasets are only composed of pavement crack scenes. As a result, the Pavementscapes dataset exhibits six types of pavement visual damages, while the other two datasets only include crack damages. Besides, the other two datasets do not refine the crack category into some sub-categories, such as longitudinal, lateral, and alligator cracks.

*Scene complexity.* The scene complexity is assessed on the Pavementscapes dataset, where the dataset is split based on the environmental conditions when the images were collected.
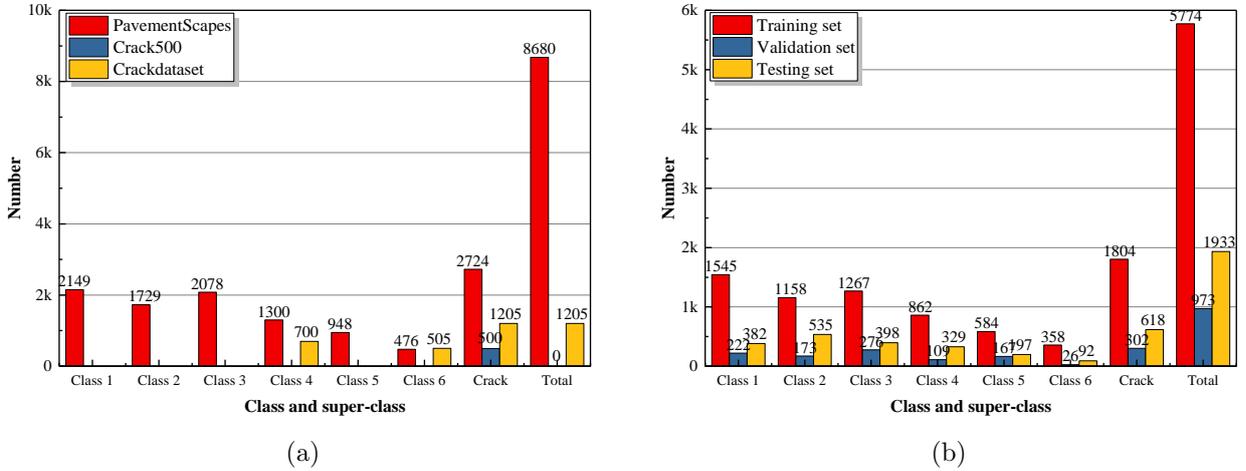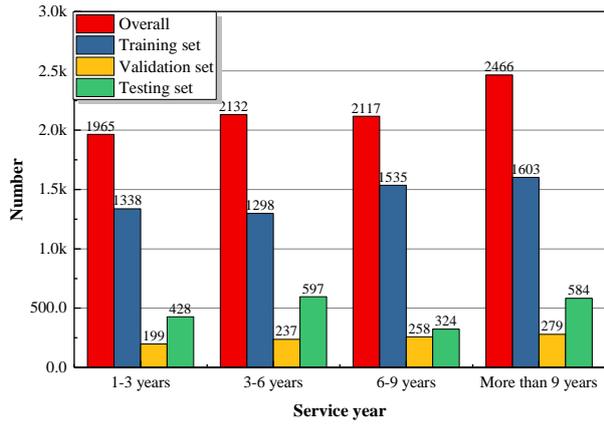
Figure 6: Damage distribution: (a) comparison of different datasets and (b) distribution in the training, validation, and testing sets of the Pavementscapes dataset. Class 1, Class 2, Class 3, Class 4, Class 5, and Class 6 stand for "pothole", "rut", "repair area", "longitudinal crack", "lateral crack", and "alligator crack", respectively. The number of the "crack" super-class are the sum of the numbers of Class 4, Class 5, and Class 6.
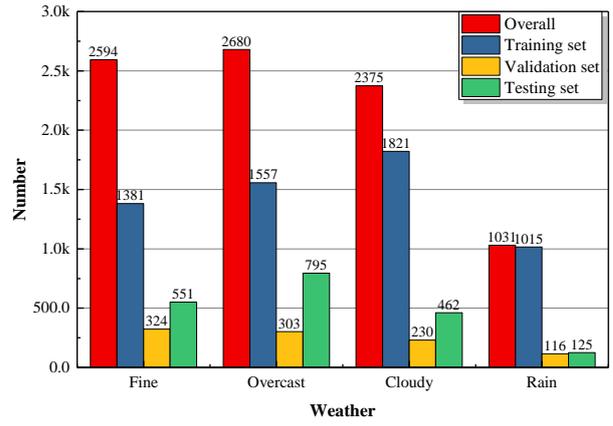
Figure 7 shows the image distributions under various real-world conditions. The distributions of the pavement images are uniform under various service years, sunlight conditions, and surface materials, which ensures the comprehensiveness of the proposed dataset. Once given the proposed dataset, deep neural networks can learn the knowledge of damage features in different environments, which can ensure the generality of deep learning models.

*Annotation accuracy.* The quality of the annotations is assessed in the study. First, 50 images were randomly selected and labeled three times by different annotators following the quality control in Section 3.2. More than 92% of pixels were labeled as the same classes. Second, the annotators were required to select a "background" label if they did not have the full certainty about the pixel class, such as some small damage instances. After excluding the "background" pixels, we recounted 95% agreement in the annotations of the 50 images. Finally, all annotations of different categories of cracks were coarsely annotated as "crack" super-category, for example, the longitudinal-crack pixels are annotated as "crack". In the 50 images, 98% pixels in the coarse annotations were defined with the same category. Therefore, the Pavementscapes dataset has annotations with good and stable quality.
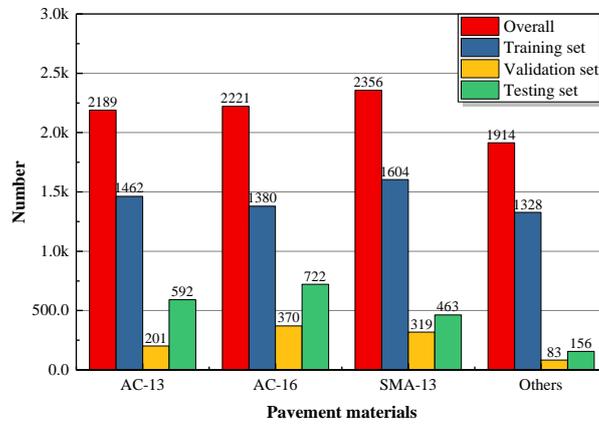
*Non-iconic views.* One goal of the proposed dataset is to collect non-iconic pavement images. Most pavement image datasets only include the images on an iconic view, such as Crack500 and CrackDataset. Besides, the current deep learning systems perform fairly well on iconic views. However, in many pavement inspection projects, many pavement images present many damages on non-iconic views, such as partially occluded cracks. Unfortunately, the current deep learning systems struggle to segment objects on the non-iconic views. The

(a)

(b)



(c)

Figure 7: Numbers of pavement damages under different real-world conditions: (a) service years, (b) weather when the images were collected, and (c) pavement materials.

Pavementscapes dataset collects 4,828 damage instances on an iconic view (e.g., repair areas in the third row and left column of Figure 5) and 3,852 damages on different non-iconic views (e.g., alligator cracks in the second row and right column of Figure 5). The proposed dataset with both iconic and non-iconic views allows to train deep neural networks with reasonable segmentation performance on different views.

## 4. Experiment

This section provides the numerical experiment that uses the Pavementscapes dataset to train and test top-performing deep neural networks. Sections 4.1 and 4.2 introduces the metrics and implementation details in the experiment, respectively. Section 4.3 discusses the performances of convolution- and attention-based deep neural networks on pavement damage segmentation. Finally, Section 4.4 provides the recommendations and future scopes on the damage segmentation task.

### 4.1. Metrics

This experiment uses five metrics to evaluate the performance of deep neural networks for damage segmentation: pixel accuracy (PA), mean intersection over union (mIoU), and expected calibration error (ECE), floating point operations (FLOPs), and network parameters.

*Pixel accuracy.* Let $\Omega = \{\omega_0, \omega_1, \ldots, \omega_m\}$ be the set of classes, where $\omega_0$ is the "background" class and $\omega_i$ is one of the damage classes, $i = 1, \ldots .m$. Given an image with $T$ pixels, the *pixel accuracy* is defined as

$$PA = \frac{1}{|T|} \sum_{j=1}^{|T|} \mathbb{1}_{\omega(j)} \left( \widehat{\omega}(j) \right), \tag{1}$$

where $\omega_*(j)$ and $\widehat{\omega}(j)$ are the labeled and predicted class of pixel $j$, and $\mathbb{1}$ is the indicator function of class $\omega(j)$. Note that the pixels belonging to the "background" class do not consider in the metric, as do in many benchmark datasets [7, 9, 13, 30].

*Mean intersection over union.* This metric measures overlap between labeled and predicted areas of a object as

$$mIoU = \frac{1}{m} \sum_{i=1}^{M} \frac{|\boldsymbol{G}^i \cap \boldsymbol{P}^i|}{|\boldsymbol{G}^i \cup \boldsymbol{P}^i|} \tag{2}$$

where $\boldsymbol{G}^i$ and $\boldsymbol{P}^i$ are the ground truth and predicted pixel set of class $i$. This experiment do not consider the IoU of the "background" class, as do in many benchmark datasets [7, 9, 13, 30].

*Expected calibration error.* In a learning system, a network should not only make correct predictions but also show when it may fail. The *confidence* of a network is defined as a mass of belief supporting the hypothesis that the prediction of a network is correct. This experiment utilizes the *expected calibration error* (ECE) [15] to measure the confidence of a network and calibrate whether its confidence matches its accuracy. First, the *prediction confidence* of pixel $j$ is defined as

$$pc(j) = \widehat{p}(\omega_*(j)), \tag{3}$$

where $\widehat{p}(\omega_*(j))$ is the predicted probability for pixel $j$ in its true class. Let $b_k$ be the bin of pixels whose prediction confidence falls into the interval $(\frac{k-1}{K}, \frac{k}{K}]$, $k = 1, \ldots, K$. The accuracy and confidence of bin $b_k$ are then computed, respectively, as

$$ac(b_k) = \frac{1}{|b_k|} \sum_{j \in b_k} \mathbb{1}_{\omega(j)}\left(\widehat{\omega}(j)\right), \tag{4a}$$

$$co(b_k) = \frac{1}{|b_k|} \sum_{j \in b_k} pc(j). \tag{4b}$$

A network is well calibrated with $ac(b_k) \approx co(b_k)$ for all bins, and the ECE[2] is defined as

$$ECE = \frac{\sum_{k=1}^{K} |b_k| \times |co(b_k) - ac(b_k)|}{\sum_{k'=1}^{K} |b'_k|}. \tag{5}$$

The ECE in this experiment does not consider the "background" pixels.

*Floating point operations and network parameters.* This experiment measures the complexity of a deep neural network using floating point operations (FLOPs) and network parameters. FLOPs are widely used to describe how many operations are required to run a single instance in a deep neural network [11, 16, 59]; calculation processes can be found in [20]. Lower values of FLOPs and network parameters always mean that an algorithm processes a new instance with fewer computation costs.

### 4.2. Implementation details

The experiment only focuses on deep neural networks for pavement damage segmentation because many previous studies [3, 4, 50, 51] have demonstrated that deep learning completely outperforms the other machine learning algorithms involving manual feature engineerings, such as support vector machine and random forest.

The experiment considers the convolution- and attention-based deep neural networks for pavement damage segmentation. For the convolution-based deep neural networks, the

---

[2]The code of ECE is available at `https://github.com/tongzheng1992/E-FCN`, which has been released by the first author in the previous study [54].

Table 3: Architecture of VGG16 network.

| Stage | Layer | Details |
|-------|-------|---------|
| | Conv 1-1 | 3×3 Conv. 16 *ReLu* with 1 strides |
| Stage 1 | Conv 1-2 | 3×3 Conv. 16 *ReLu* with 1 strides |
| | Pooling | 2×2 max-pooling with 2 strides |
| | Conv 2-1 | 3×3 Conv. 32 *ReLu* with 1 strides |
| Stage 2 | Conv 2-2 | 3×3 Conv. 32 *ReLu* with 1 strides |
| | Pooling | 2×2 max-pooling with 2 strides |
| | Conv 3-1 | 3×3 Conv. 64 *ReLu* with 1 strides |
| Stage 3 | Conv 3-2 | 3×3 Conv. 64 *ReLu* with 1 strides |
| | Conv 3-3 | 3×3 Conv. 64 *ReLu* with 1 strides |
| | Pooling | 2×2 max-pooling with 2 strides |
| | Conv 4-1 | 3×3 Conv. 128 ReLu with 1 strides |
| Stage 4 | Conv 4-2 | 3×3 Conv. 128 *ReLu* with 1 strides |
| | Conv 4-3 | 3×3 Conv. 128 *ReLu* with 1 strides |
| | Pooling | 2×2 max-pooling with 2 strides |
| | Conv 5-1 | 3 3 Conv. 256 ReLu with 1 strides |
| Stage 5 | Conv 5-2 | 3×3 Conv. 256 *ReLu* with 1 strides |
| | Conv 5-3 | 3×3 Conv. 256 *ReLu* with 1 strides |
| | Pooling | 2×2 max-pooling with 2 strides |

Pavementscapes dataset is used to train and test a series of the original FCN models (FCN-32s, FCN-16s, and FCN-8s) [33], U-net [40], and DeepLabv3+ [5]. These models use the same backbone, VGG16, as shown in Table 3. For the attention-based deep neural networks, four models was considered, including self-attention net [57], criss-cross attention (CC-attention) net [19], double-attention net [6], and segmentation transformer [62]. The patch size of the attention-based models are $32 \times 32$. Other detailed hyper-parameters of these networks is the same as their original works.

During training, all networks use the generalized dice loss function [47], which reduces the negative effects of unbalanced learning set. In this study, the unbalanced learning set means that the number of pixels belonging to different classes are very different, such that the proposed dataset includes a very small number of crack pixels and a very huge number of background pixels in the training set of the Pavementscapes dataset. The phenomenon cannot be avoided because cracks only occupy very small areas in a pavement. Given a pixel with one-hot label $\boldsymbol{y}$ and predicted probabilities $\widehat{\boldsymbol{p}}$, the generalized dice loss function is defined as

$$\mathcal{L}(\boldsymbol{y}, \widehat{\boldsymbol{p}}) = 1 - \frac{2\boldsymbol{y}\widehat{\boldsymbol{p}} + 1}{\boldsymbol{y} + \widehat{\boldsymbol{p}} + 1}. \tag{6}$$

All models are achieved based on TensorFlow 2.8 version. The input image size are set as $1024 \times 2048 \times 1$. The training batch size is 24 and the popular ADAM optimizer with momentum 0.9 and weight decay 1e-4 is used to optimize the networks for backpropagation. Note that some types of data augmentations cannot be used in the Pavementscapes dataset

15

Table 4: Testing performances of deep neural networks on the Pavementscapes dataset. GFLOPs stands for $10^9$ (giga) floating point operations and $M$ means million. The best and second results in each term are marked in bold and italics.

|  | PA/% | mIoU | ECE/% | GFLOPs | Parameter/M |
|---|---|---|---|---|---|
| FCN-32s [33] | 66.53 | 51.94 | 22.48 | - | - |
| FCN-16s [33] | 67.02 | 52.21 | 22.41 | - | - |
| FCN-8s [33] | 67.32 | 52.98 | 22.30 | **177** | 134.3 |
| U-net [40] | 69.56 | 54.71 | 22.14 | *194* | 19.4 |
| DeepLabv3+ [5] | 71.90 | 57.51 | 21.81 | 783 | 41.1 |
| Self-attention net [57] | 73.07 | 58.74 | 21.32 | 619 | *10.5* |
| CC-attention net [19] | 73.15 | 58.52 | 21.14 | 804 | 10.6 |
| Double-attention net [6] | *74.01* | *59.23* | *21.10* | 338 | **10.2** |
| Segmentation Transformer [62] | **74.50** | **59.74** | **20.95** | 340 | 10.5 |

since flips and rotations change the semantics of longitudinal and lateral cracks. The deep neural networks are trained on an Nvidia V100 GPU with 32GB memory.

### 4.3. Results of damage segmentation

Table 4 display the overall test performance of the deep neural networks. Some testing examples are shown in Appendix A. The attention-based networks exceed the convolution-based ones on PA and mIoU, even though the FLOPs and network parameters of the attention-based models are larger than the ones of the convolution-based ones. In detail, the segmentation transformer model achieves the best segmentation performance, followed by double-attention and CC-attention nets. This demonstrates that the attention-based models outperform the convolution-based ones on pavement damage segmentation. This behavior can be explained by Table 5. The two types of deep neural networks have similar and high PAs and mIoUs in the "rut", "repair area", and "pothole" classes, but the PAs and mIoUs of convolution-based models in the three crack classes are lower than these of the attention-based models. This is mainly because the attention mechanism allows the attention-based models to focus on the features of some small target objects [34], such as the cracks with thin widths. However, the convolution-based feature extraction in the convolution-based models easily ignores the small features, e.g., the boundary areas of cracks and background. Therefore, attention-based deep neural networks have the powerful potential for the accuracy improvement of pavement damage segmentation, especially for some small damages.

Table 4 indicates that the two types of deep neural networks can accurately segment the "rut", "repair area", and "pothole" instances but cannot do it well on the three types of cracks. This problem derives from the fact that the proportion of crack pixels in the training set is much lower than the ones of other classes. Unfortunately, the fact cannot be changed in the projects of pavement inspection because the cracks and some other damages only occupy very small parts of a pavement. In detail, the proportion of crack pixels is less than 1%. Thus, the training set of the Pavementscapes dataset is very unbalanced. Such an unbalanced training set makes deep neural networks tend to classify crack pixels to

Table 5: Testing IoU results of different damage classes on the Pavementscapes dataset. The best and second results in each term are marked in bold and italics.

| | Longitudinal crack | Lateral crack | Alligator crack | Pothole | Rut | Repair area |
|---|---|---|---|---|---|---|
| FCN-32s | 25.82 | 27.17 | 25.38 | 67.17 | 75.92 | 90.17 |
| FCN-16s | 25.36 | 26.51 | 25.56 | 67.32 | 76.17 | 92.34 |
| FCN-8s | 25.93 | 27.27 | 26.42 | 68.13 | 77.42 | 92.73 |
| U-net | 26.99 | 30.26 | 29.14 | 69.33 | 80.10 | 92.45 |
| DeepLabv3+ | 33.12 | 35.25 | 33.84 | 71.31 | 79.25 | 92.31 |
| Self-attention net | *35.62* | 37.25 | 34.92 | 72.82 | 78.92 | *92.89* |
| CC-attention net | 35.46 | 37.36 | 34.74 | 72.14 | 78.82 | 92.58 |
| Double-attention net | 35.00 | *38.29* | *36.12* | *73.01* | *80.22* | 92.74 |
| Segmentation Transformer | **36.20** | **39.14** | **36.42** | **73.42** | **80.42** | **92.81** |

"background" during training since the trend does not introduce a large loss, even though the generated dice loss 6 has been used to reduce the negative effect. This behavior demonstrates that more advanced loss functions should be considered in the future to train the deep neural networks for the pavement damage segmentation task.

Table 4 also shows that the two types of deep neural networks have similar ECEs, demonstrating the two types of deep neural networks are over-confident because their accuracies do not match their confidences. Figures 8 and 9, respectively, shows the pixel distribution and pixel accuracy histograms of the deep neural networks in Table 4. Note that the background pixels are not considered in the two figures. The average confidence of each network is substantially higher than its average pixel accuracy, indicating that the network is not calibrated. This problem is mainly because the deep neural networks work within the probabilistic framework, in which the features from the backbone are imported into a softmax layer to generate probabilities of the classes for decision-making. Probability theory only captures the randomness aspect of the features but neither ambiguity nor incompleteness [22, 23], which are inherent in damage features. For example, a deep neural network may extract incomplete damage features from an image with a non-iconic view. Besides, a network sometimes extracts ambiguous features from some small damages. Such uncertain features lead that multiple classes having similar probabilities. In such a case, deep neural networks in the probabilistic framework often arbitrarily assign the pixel to one and only one of the possible classes, which may result in misclassification and finally leads to over-confidence. The problem of over-confidence is common in the deep neural networks work within the probabilistic framework [15]. Section 4.4 will provide a potential way to solve the problem.

Table 4 indicates the conflict between performance and computation cost. Compared with the convolution-based models except for DeepLabv3+, the attention-based models cost more FLOPs but have larger PAs and mIoUs. Moreover, the training time of an attention-based model is twice that of a convolution-based model, demonstrating a small improvement in terms of PAs and mIoUs always requires large increases in computation costs. The costs of the attention-based networks are unbearable, even though the GPU memory and computation force has a significant increase in recent years. Therefore, attention-based deep neural networks with light weights are required for the task of pavement damage segmentation.
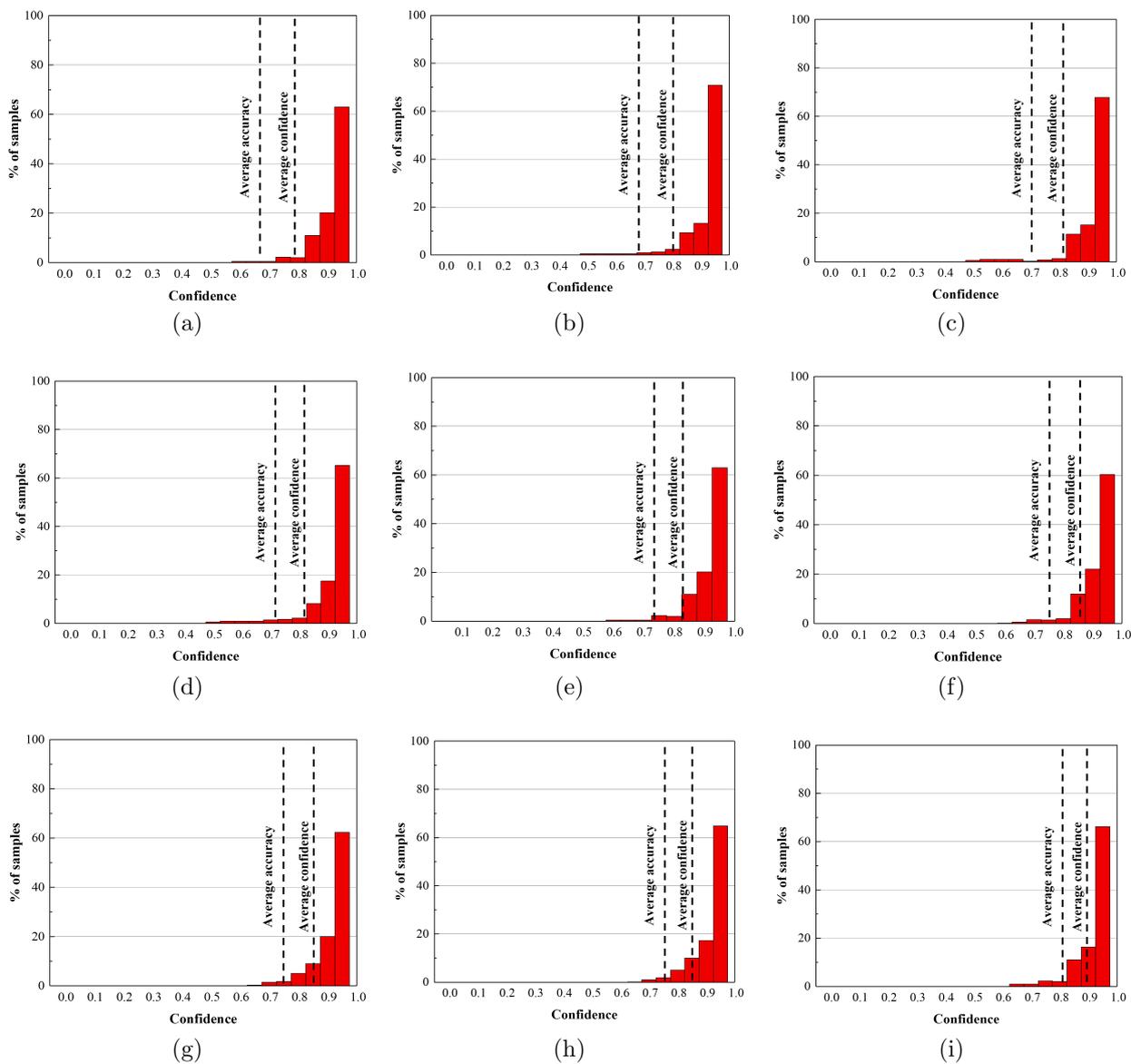
17

Figure 8: Testing pixel distribution on the Pavementscapes dataset: (a) FCN-32s, (b) FCN-16s, (c) FCN-8s, (d) U-net, (e) DeepLabv3+, (f) Self-attention net, (g) CC-attention net, (h) Double-attention net, (i) Segmentation transformer.
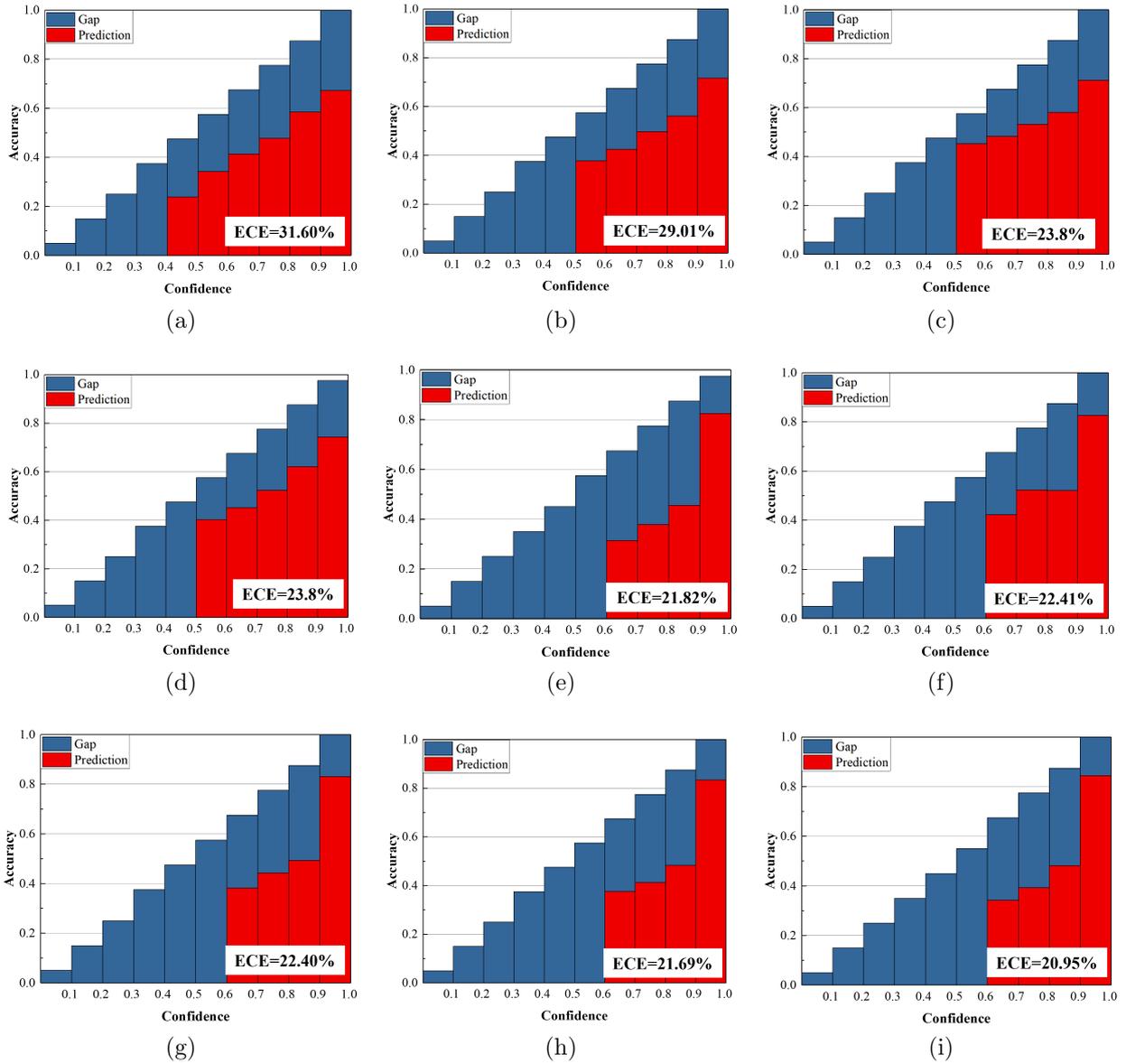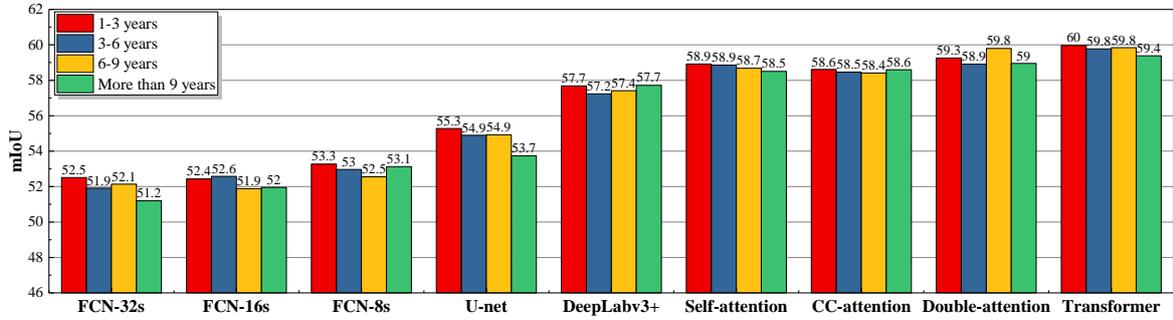
Figure 9: Testing pixel accuracy histograms on the Pavementscapes dataset: (a) FCN-32s, (b) FCN-16s, (c) FCN-8s, (d) U-net, (e) DeepLabv3+, (f) Self-attention net, (g) CC-attention net, (h) Double-attention net, (i) Segmentation transformer.

Figure 10 presents the mIoUs of the two types of deep neural networks under various real-world conditions. These networks have stable performances on different service years, weathers, and pavement materials since their IoUs and PUs do not significantly change under different conditions. This demonstrates that the two types of deep neural networks can perform damage segmentation well in the real world.

### 4.4. Recommendations and future scopes

Deep neural networks achieve some good results on the Pavementscapes dataset, which can be considered as the baseline of the PavmentScapes challenges. However, they still face several problems that are not easy to be solved by using traditional deep neural networks. This study provides the recommendations for these problems as follows.

1. *Small damage instances.* Current deep neural networks sometimes ignore cracks with thin widths, especially the convolution-based networks, though the features of these cracks are important for decision-making. This behavior makes these networks have low performances on crack segmentation. Attention-based deep neural networks show the potential capacity to solve this problem. In the future, more attention-based networks should be trained and tested by the Pavementscapes dataset, which may replace the convolution-based networks for pavement damage segmentation.

2. *Unbalanced training set.* Even though the Pavementscapes dataset has the real-world distribution of different damage instances, it still has the unbalanced problems in the numbers of different classes, such that the background pixels is much more than the sum of the damage pixels and the crack pixels are less than the pixels of other damages. This fact introduces a negative effect on learning systems that tend to assign a pixel to the background class or the damage classes with a large number of pixels. Similar behaviors are common in the segmentation task of medical and cell images [8, 27]. In the medical image segmentation, many morphology-based loss functions are used to solve the unbalanced problem, such as focal loss [29], dice loss [47] used in the experiment, and IoU Loss [66]. Therefore, these loss functions should be introduced into the deep neural networks to improve the performance of crack segmentation.

3. *Over-confidence.* Two types of deep neural networks are not calibrated well in the damage segmentation task. This problem derives from the use of the probability framework. During the last decade, many theories have been combined with deep neural networks to solve the problem, and one of the successful cases is the evidential deep neural network [52], which converts the features from the backbone of a deep neural network into Dempster-Shafer belief functions, rather than the probabilities using a softmax layer. This architecture allows the network to represent the feature uncertainty [53] and reduce the confidence of the network [54]. Such architecture should be considered to be combined with the attention-based models to make the deep neural networks well-calibrated.

Figure 10: Stability anaylsis using the Pavementscapes dataset under different (a) service years, (b) weathers, and (c) pavement materials.

## 5. Conclusions

This study has proposed a large-scale hierarchical image dataset for asphalt pavement damage segmentation, called Pavementscapes. The statistical study and the deep learning experiment provide an in-depth analysis of the dataset. The following conclusions are can be drawn.

1. The Pavementscapes dataset consists of 4,000 pavement images with a resolution of $1024 \times 2048$ and 8,680 damage instances, which were recorded from several real-world projects of pavement inspection in China. Six damage classes are included in the dataset. The proposed dataset exceeds the other public pavement dataset in the number of pavement images, damage classes, annotation levels, and shooting views.

2. The statistical analysis demonstrates that the Pavementscapes dataset has a reasonable damage distribution, complex pavement scenes, and high annotation accuracy, which ensure the completeness and comprehensiveness of the dataset. In addition, the images with different non-iconic views improve the complexity and truth of the dataset. In summary, the dataset represents the real-world pavement damages well.

3. The numerical experiment uses the Pavementscapes dataset to train and test the top-performance deep neural networks. The results demonstrate that the deep neural networks have the powerful potential to segment pavement damages once given enough good training samples. The attention-based models outperform the convolution-based ones on the segmentation task, which may be the new direction for visual damage segmentation. The experiment results can be considered as the baseline for the public damage segmentation challenge.

4. The results of the numerical experiment indicate three problems with the use of deep neural networks in pavement damage segmentation: the segmentation of small damage instances, the unbalanced training set, and the over-confidence of deep neural networks. The three problems are not easy to solve using the current state-of-the-art deep networks.

5. Future work will focus on three main aspects, corresponding to the above three problems. First, the attention mechanism should be further applied in the segmentation of small damage instances, which has the potential to improve the performance of crack segmentation. Other advanced morphology-based loss functions should be introduced into deep neural networks to solve the problem of unbalanced learning set. Finally, some uncertainty frameworks, such as Dempster-Shafer theory, should be used to overcome the over-confidence problem.

## Author Contributions

Zheng Tong: conception and study design, data collection and analysis and interpretation, drafting the article; Tao Ma: conception and study design, reviewing the article; Ju Huyan: data collection and reviewing the article.

## Acknowledgment

## References

[1] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X., 2015. TensorFlow: Large-scale machine learning on heterogeneous systems. URL: https://www.tensorflow.org/. software available from tensorflow.org.

[2] Arya, D., Maeda, H., Ghosh, S.K., Toshniwal, D., Mraz, A., Kashiyama, T., Sekimoto, Y., 2021. Deep learning-based road damage detection and classification for multiple countries. Automation in Construction 132, 103935.

[3] Cha, Y.J., Choi, W., Büyüköztürk, O., 2017. Deep learning-based crack damage detection using convolutional neural networks. Computer-Aided Civil and Infrastructure Engineering 32, 361–378.

[4] Cha, Y.J., Choi, W., Suh, G., Mahmoudkhani, S., Büyüköztürk, O., 2018. Autonomous structural visual inspection using region-based deep learning for detecting multiple damage types. Computer-Aided Civil and Infrastructure Engineering 33, 731–747.

[5] Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H., 2018a. Encoder-decoder with atrous separable convolution for semantic image segmentation, in: Proceedings of the European conference on computer vision (ECCV), pp. 801–818.

[6] Chen, Y., Kalantidis, Y., Li, J., Yan, S., Feng, J., 2018b. Aˆ 2-nets: Double attention networks. Advances in neural information processing systems 31.

[7] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B., 2016. The cityscapes dataset for semantic urban scene understanding, in: Proc. of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255.

[8] Craddock, C., Benhajali, Y., Chu, C., Chouinard, F., Evans, A., Jakab, A., Khundrakpam, B.S., Lewis, J.D., Li, Q., Milham, M., et al., 2013. The neuro bureau preprocessing initiative: open sharing of preprocessed neuroimaging data and derivatives. Frontiers in Neuroinformatics 7.

[9] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer vision and pattern recognition, IEEE. pp. 248–255.

[10] Dorafshan, S., Thomas, R.J., Maguire, M., 2018. SDNET2018: An annotated image dataset for non-contact concrete crack detection using deep convolutional neural networks. Data in brief 21, 1664–1668.

[11] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al., 2021. An image is worth 16x16 words: Transformers for image recognition at scale, in: Proceedings of the 2021 International Conference on Learning Representations, Vienna, Austria. pp. 1–21.

[12] Eisenbach, M., Stricker, R., Seichter, D., Amende, K., Debes, K., Sesselmann, M., Ebersbach, D., Stoeckert, U., Gross, H.M., 2017. How to get pavement distress detection ready for deep learning? a systematic approach, in: 2017 international joint conference on neural networks (IJCNN), IEEE. pp. 2039–2047.

[13] Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A., 2015. The pascal visual object classes challenge: A retrospective. International journal of computer vision 111, 98–136.

[14] Guan, J., Yang, X., Ding, L., Cheng, X., Lee, V.C., Jin, C., 2021. Automated pixel-level pavement distress detection based on stereo vision and deep learning. Automation in Construction 129, 103788.

[15] Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q., 2017. On calibration of modern neural networks, in: Proceedings of the 34th International Conference on Machine Learning, JMLR.org. p. 1321–1330.

[16] Guo, J., Han, K., Wu, H., Xu, C., Tang, Y., Xu, C., Wang, Y., 2021. CMT: Convolutional neural networks meet vision transformers. arXiv preprint arXiv:2107.06263 .

[17] Han, C., Ma, T., Huyan, J., Huang, X., Zhang, Y., 2021. Crackw-net: A novel pavement crack image segmentation convolutional neural network. IEEE Transactions on Intelligent Transportation Systems , 1–13.

[18] Hoiem, D., Chodpathumwan, Y., Dai, Q., 2012. Diagnosing error in object detectors, in: European conference on computer vision, Springer. pp. 340–353.

[19] Huang, Z., Wang, X., Huang, L., Huang, C., Wei, Y., Liu, W., 2019. Ccnet: Criss-cross attention for semantic segmentation, in: Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision, Seoul , South Korea. pp. 603–612.

[20] Hunger, R., 2005. Floating point operations in matrix-vector calculus. Munich University of Technology, Inst. for Circuit Theory and Signal.

[21] Huyan, J., Li, W., Tighe, S., Xu, Z., Zhai, J., 2020. Cracku-net: A novel deep convolutional neural network for pixelwise pavement crack detection. Structural Control and Health Monitoring 27, e2551.

[22] Jiao, L., Denœux, T., Pan, Q., 2015a. A hybrid belief rule-based classification system based on uncertain training data and expert knowledge. IEEE Transactions on Systems, Man, and Cybernetics: Systems 46, 1711–1723.

[23] Jiao, L., Pan, Q., Denoeux, T., Liang, Y., Feng, X., 2015b. Belief rule-based classification system: Extension of frbcs in belief functions framework. Information Sciences 309, 26–49.

[24] JTG H20—2007, 2008. Highway Performance Assessment Standards. Standard. China Communications Press. Beijing, China.

[25] Kang, D.H., Cha, Y.J., 2021. Efficient attention-based deep encoder and decoder for automatic crack segmentation. Structural Health Monitoring , 14759217211053776.

[26] LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. nature 521, 436–444.

[27] Lévy, D., Jain, A., 2016. Breast mass classification from mammograms using deep convolutional neural networks. arXiv preprint arXiv:1612.00542 .

[28] Lin, J., Liu, Y., 2010. Potholes detection based on svm in the pavement distress image, in: 2010 Ninth International Symposium on Distributed Computing and Applications to Business, Engineering and Science, IEEE. pp. 544–547.

[29] Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2017. Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, Venice, Italy. pp. 2980–2988.

[30] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L., 2014. Microsoft coco: Common objects in context, in: European conference on computer vision, Springer. pp. 740–755.

[31] Liu, C., Chen, L.C., Schroff, F., Adam, H., Hua, W., Yuille, A., Fei-Fei, L., 2019a. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation, in: CVPR, pp. 1–13.

[32] Liu, Z., Cao, Y., Wang, Y., Wang, W., 2019b. Computer vision-based concrete crack detection using u-net fully convolutional networks. Automation in Construction 104, 129–139.

[33] Long, J., Shelhamer, E., Darrell, T., 2015. Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3431–3440.

[34] Niu, Z., Zhong, G., Yu, H., 2021. A review on the attention mechanism of deep learning. Neurocomputing 452, 48–62.

[35] Oktay, O., Schlemper, J., Folgoc, L.L., Lee, M., Heinrich, M., Misawa, K., Mori, K., McDonagh, S., Hammerla, N.Y., Kainz, B., et al., 2018. Attention u-net: Learning where to look for the pancreas. arXiv preprint arXiv:1804.03999 .

[36] Oliveira, H., Correia, P.L., 2014. Crackit—an image processing toolbox for crack detection and characterization, in: 2014 IEEE international conference on image processing (ICIP), IEEE. pp. 798–802.

[37] Ouma, Y.O., Hahn, M., 2017. Pothole detection on asphalt pavements from 2d-colour pothole images using fuzzy c-means clustering and morphological reconstruction. Automation in Construction 83, 196–211.

[38] Palmer, S., 1981. Canonical perspective and the perception of objects. Attention and performance , 135–151.

[39] Qin, X., Zhang, Z., Huang, C., Dehghan, M., Zaiane, O.R., Jagersand, M., 2020. U2-Net: Going deeper with nested u-structure for salient object detection. Pattern Recognition 106, 107404.

[40] Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer. pp. 234–241.

[41] Shi, Y., Cui, L., Qi, Z., Meng, F., Chen, Z., 2016. Automatic road crack detection using random structured forests. IEEE Transactions on Intelligent Transportation Systems 17, 3434–3445.

[42] Silva, W.R.L.d., Lucena, D.S.d., 2018. Concrete cracks detection based on deep learning image classification, in: Multidisciplinary digital publishing institute proceedings, p. 489.

[43] Song, L., Wang, X., 2021. Faster region convolutional neural network for automated pavement distress detection. Road Materials and Pavement Design 22, 23–41.

[44] Stricker, R., Aganian, D., Sesselmann, M., Seichter, D., Engelhardt, M., Spielhofer, R., Hahn, M., Hautz, A., Debes, K., Gross, H.M., 2021. Road surface segmentation - pixel-perfect distress and object detection for road assessment., in: International Conference on Automation Science and Engineering (CASE), pp. 1–8.

[45] Stricker, R., Eisenbach, M., Sesselmann, M., Debes, K., Gross, H.M., 2019. Improving visual road condition assessment by extensive experiments on the extended gaps dataset, in: 2019 International Joint Conference on Neural Networks (IJCNN), IEEE. pp. 1–8.

[46] Strudel, R., Garcia, R., Laptev, I., Schmid, C., 2021. Segmenter: Transformer for semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7262–7272.

[47] Sudre, C.H., Li, W., Vercauteren, T., Ourselin, S., Jorge Cardoso, M., 2017. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations, in: Deep learning in medical image analysis and multimodal learning for clinical decision support. Springer, pp. 240–248.

[48] Sun, X., Xie, Y., Jiang, L., Cao, Y., Liu, B., 2022. DMA-Net: Deeplab with multi-scale attention for pavement crack segmentation. IEEE Transactions on Intelligent Transportation Systems .

[49] Tong, Z., Gao, J., Han, Z., Wang, Z., 2018. Recognition of asphalt pavement crack length using deep convolutional neural networks. Road Materials and Pavement Design 19, 1334–1349.

[50] Tong, Z., Gao, J., Yuan, D., 2020a. Advances of deep learning applications in ground-penetrating radar: A survey. Construction and Building Materials 258, 120371.

[51] Tong, Z., Gao, J., Zhang, H., 2017. Recognition, location, measurement, and 3d reconstruction of concealed cracks using convolutional neural networks. Construction and Building Materials 146, 775–787.

[52] Tong, Z., Xu, P., Denœux, T., 2019. ConvNet and Dempster-Shafer theory for object recognition, in: Proceedings of the 6th International Conference on Scalable Uncertainty Management, Springer, Compiégne,France. pp. 368–381.

[53] Tong, Z., Xu, P., Denœux, T., 2021a. An evidential classifier based on Dempster-Shafer theory and deep learning. Neurocomputing 450, 275–293.

[54] Tong, Z., Xu, P., Denœux, T., 2021b. Evidential fully convolutional network for semantic segmentation. Applied Intelligence 51, 6376–6399.

[55] Tong, Z., Yuan, D., Gao, J., Wang, Z., 2020b. Pavement defect detection with fully convolutional network and an uncertainty framework. Computer-Aided Civil and Infrastructure Engineering 35, 832–849.

[56] Tsai, Y.C., Kaul, V., Mersereau, R.M., 2010. Critical assessment of pavement distress segmentation methods. Journal of transportation engineering 136, 11–19.

[57] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. Advances in neural information processing systems 30.

[58] Wada, K., 2016. LabelMe: Image Polygonal Annotation with Python. `https://github.com/wkentaro/labelme`.

[59] Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K., 2017. Aggregated residual transformations for deep neural networks, in: Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Honolulu, USA. pp. 1492–1500.

[60] Yang, F., Zhang, L., Yu, S., Prokhorov, D., Mei, X., Ling, H., 2019. Feature pyramid and hierarchical boosting network for pavement crack detection. IEEE Transactions on Intelligent Transportation Systems .

[61] Ye, W., Jiang, W., Tong, Z., Yuan, D., Xiao, J., 2021. Convolutional neural network for pothole detection in asphalt pavement. Road materials and pavement design 22, 42–58.

[62] Yuan, Y., Chen, X., Chen, X., Wang, J., 2019. Segmentation transformer: Object-contextual representations for semantic segmentation. arXiv preprint arXiv:1909.11065 .

[63] Zhang, A., Wang, K.C., Fei, Y., Liu, Y., Tao, S., Chen, C., Li, J.Q., Li, B., 2018. Deep learning–based fully automated pavement crack detection on 3d asphalt surfaces with an improved cracknet. Journal of Computing in Civil Engineering 32, 04018041.

[64] Zhang, A., Wang, K.C., Li, B., Yang, E., Dai, X., Peng, Y., Fei, Y., Liu, Y., Li, J.Q., Chen, C., 2017. Automated pixel-level pavement crack detection on 3d asphalt surfaces using a deep-learning network. Computer-Aided Civil and Infrastructure Engineering 32, 805–819.

[65] Zhang, L., Yang, F., Zhang, Y.D., Zhu, Y.J., 2016. Road crack detection using deep convolutional neural network, in: 2016 IEEE international conference on image processing (ICIP), IEEE. pp. 3708–3712.

[66] Zheng, Z., Wang, P., Liu, W., Li, J., Ye, R., Ren, D., 2020. Distance-IoU loss: Faster and better learning for bounding box regression, in: Proceedings of the 34th AAAI Conference on Artificial Intelligence, New York, USA. pp. 12993–13000.

[67] Zou, Q., Cao, Y., Li, Q., Mao, Q., Wang, S., 2012. Cracktree: Automatic crack detection from pavement images. Pattern Recognition Letters 33, 227–238.

## Appendix A. Examples of segmentation results

There are eight examples in the appendix. For each example, Figures (a), (b), (c), and (d) are the original image, ground trurh, segmentation results from the DeepLabv3+, and segmentation results from the segmentation transformer, respectively. The masks with gray-scale values of 0 are the "background" pixels; other masks with different gray-scale values are the pixels belonging to different classes, such that 30, 60, 90, 120, 150, and 180 gray-scale values stand for the pixels of "longitudinal crack", "lateral crack", "alligator crack", "pothole", "rut", and "repair area", respectively.
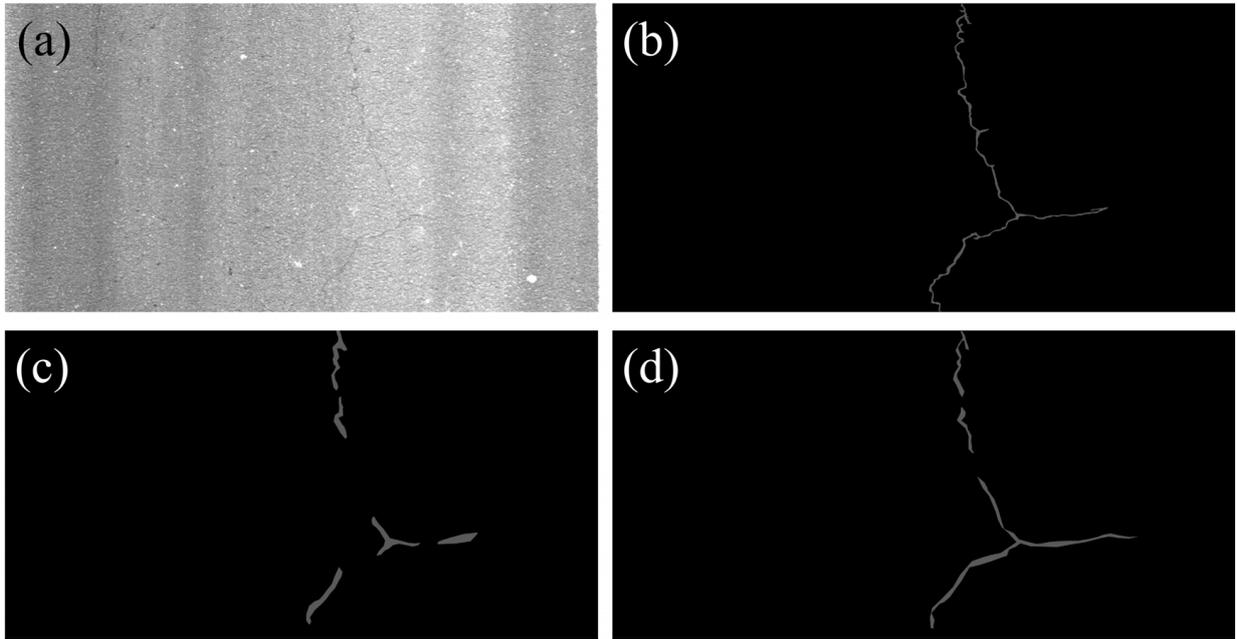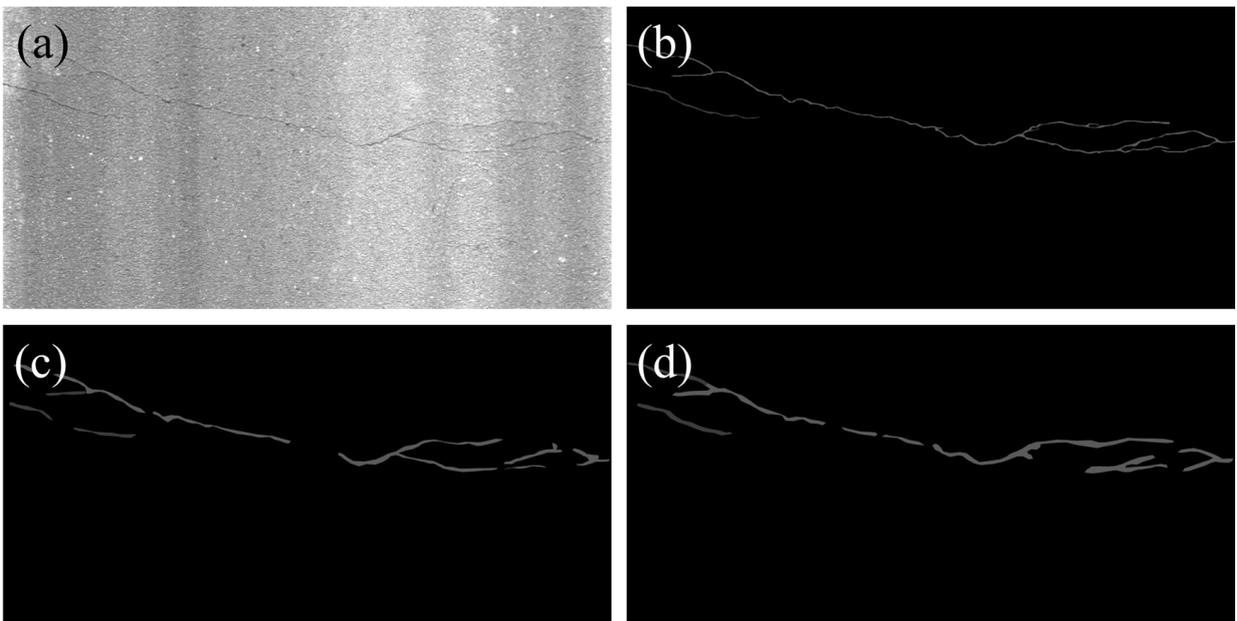
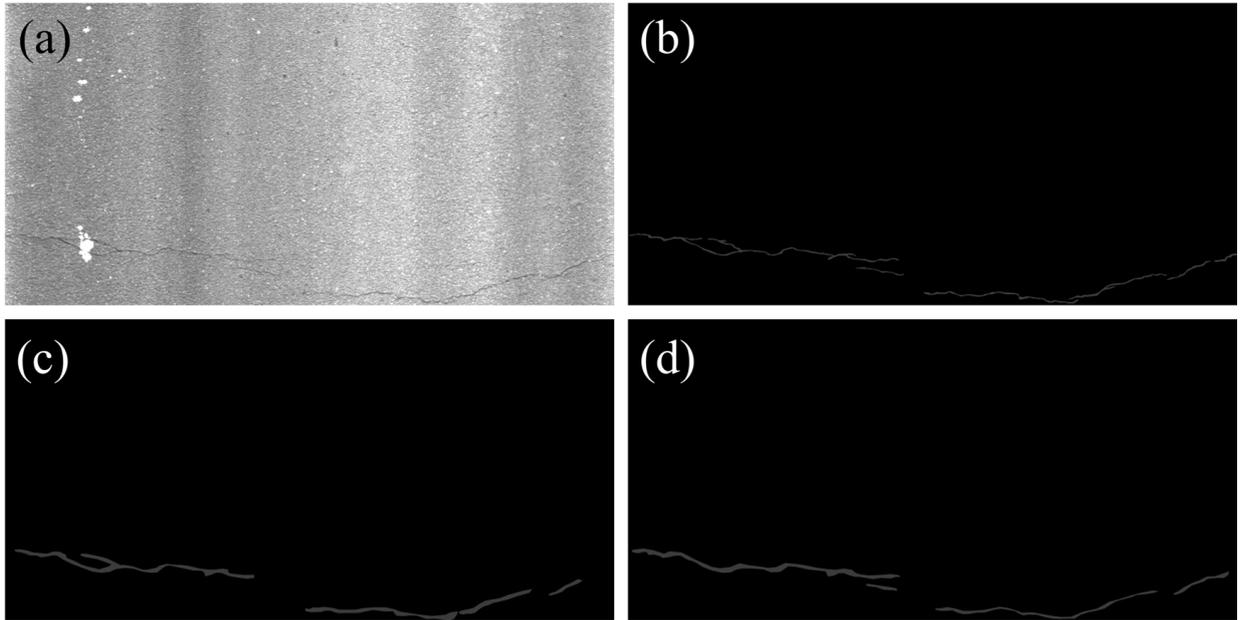Figure A.11: Example 1.



Figure A.12: Example 2.
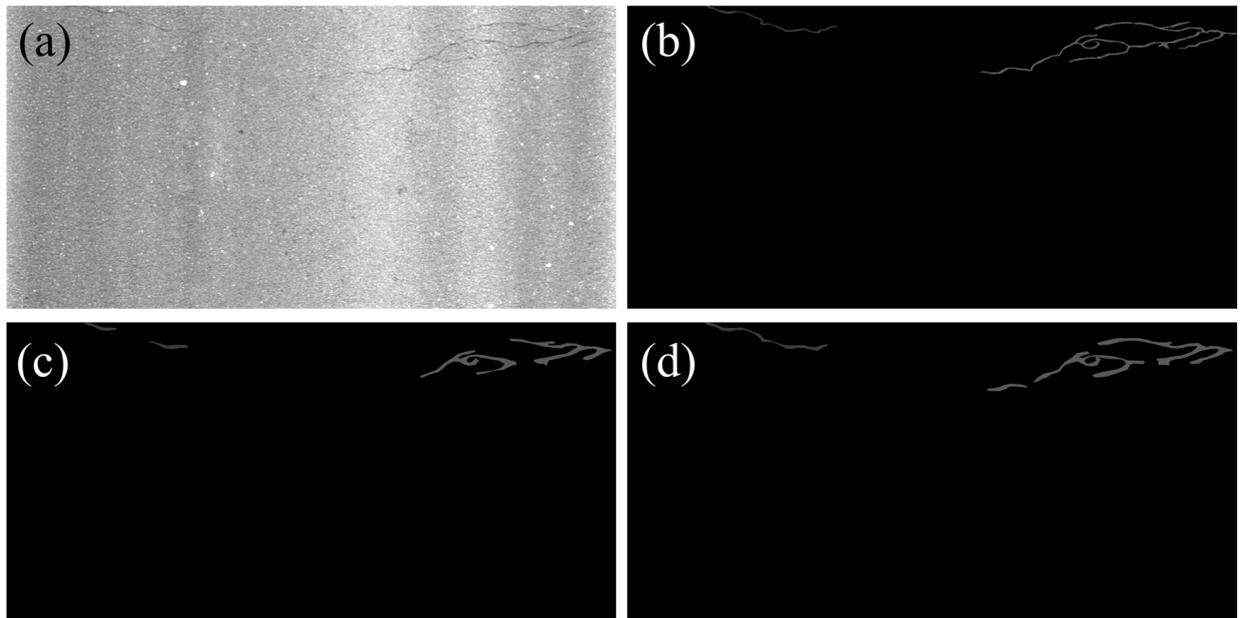
Figure A.13: Example 3.
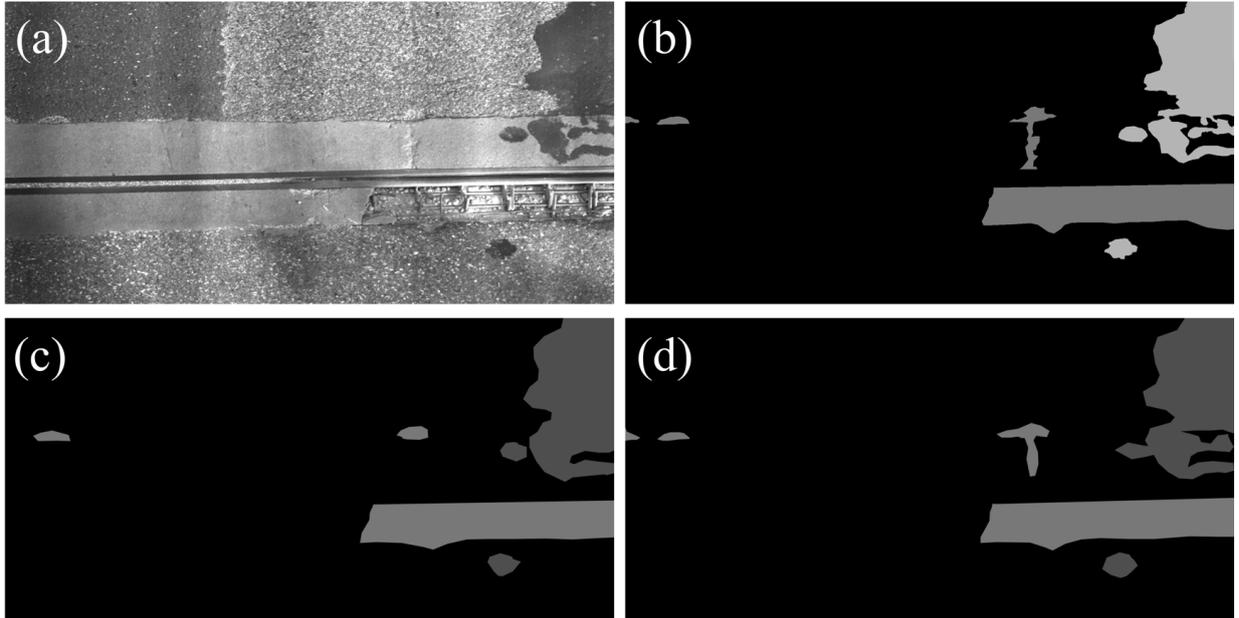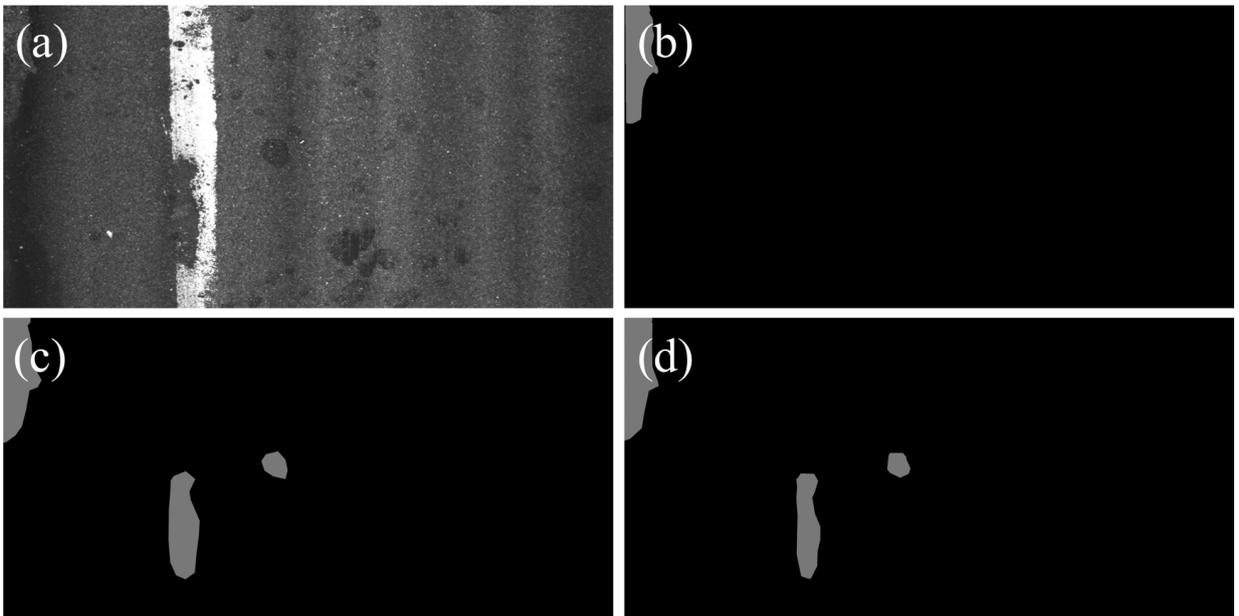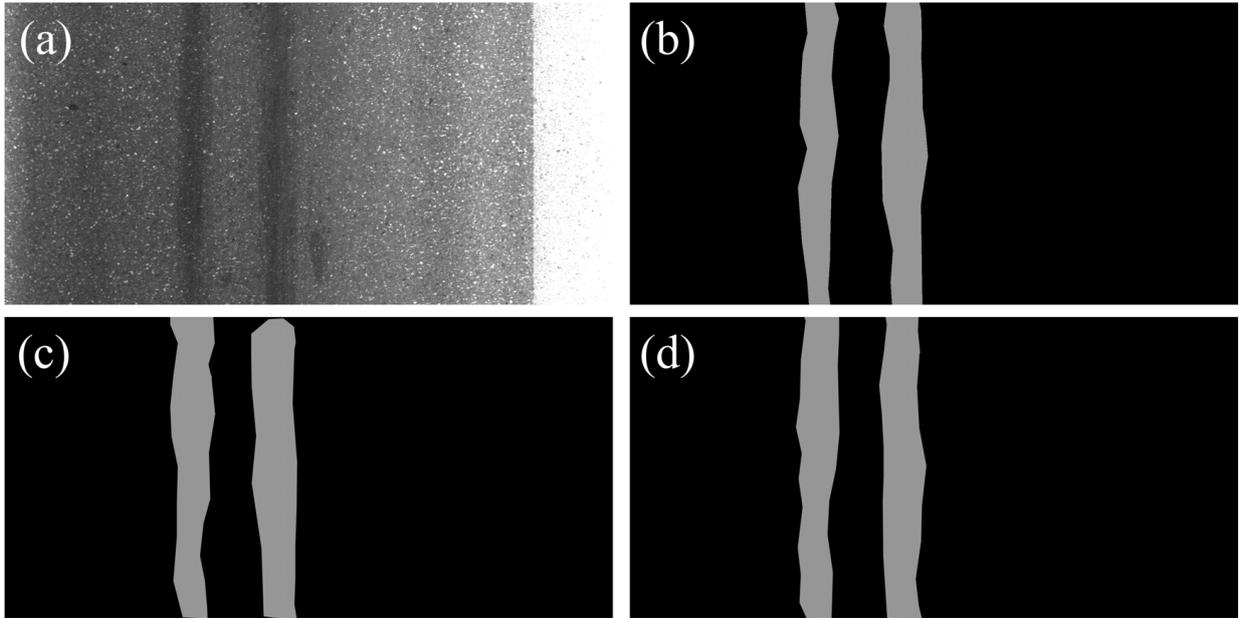


Figure A.14: Example 4.

Figure A.15: Example 5.



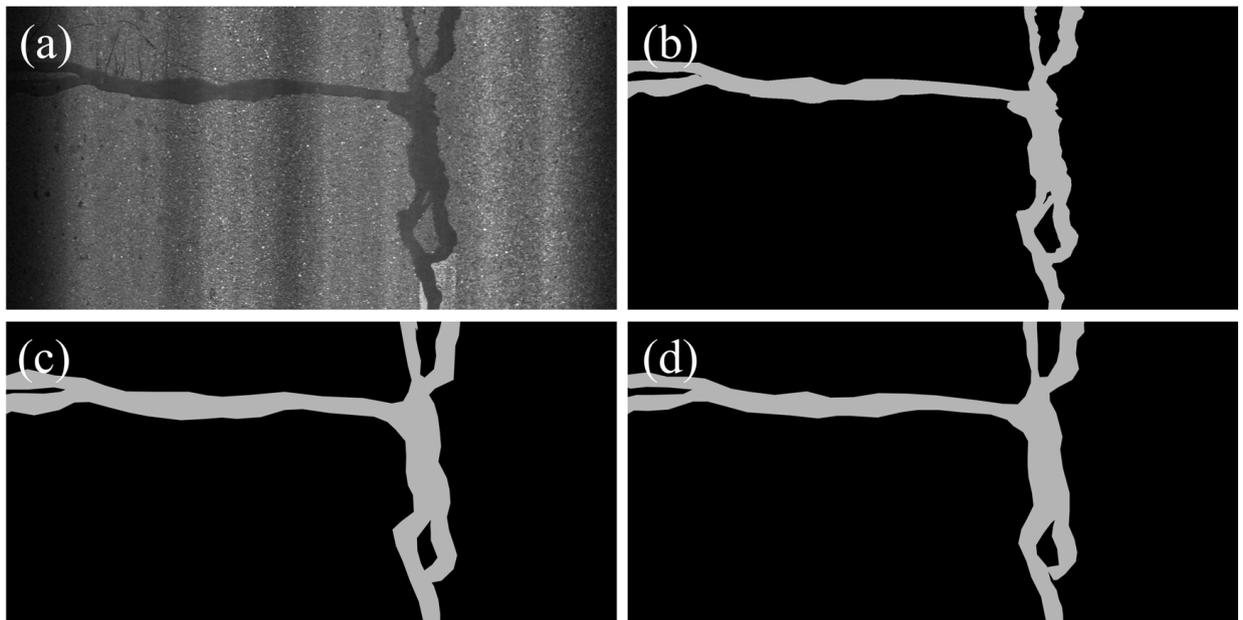Figure A.16: Example 6.

Figure A.17: Example 7.



Figure A.18: Example 8.