# Optimized Deep Encoder-Decoder Methods for Crack Segmentation

Jacob König[a,*], Mark David Jenkins[a], Mike Mannion[a], Peter Barrie[a], Gordon Morison[a]

[a]*School of Computing, Engineering and Built Environment, Glasgow Caledonian University, G4 0BA Glasgow, Scotland*

## Abstract

Surface crack segmentation poses a challenging computer vision task as background, shape, colour and size of cracks vary. In this work we propose optimized deep encoder-decoder methods consisting of a combination of techniques which yield an increase in crack segmentation performance. Specifically we propose a decoder-part for an encoder-decoder based deep learning architecture for semantic segmentation and study its components to achieve increased performance. We also examine the use of different encoder strategies and introduce a data augmentation policy to increase the amount of available training data. The performance evaluation of our method is carried out on four publicly available crack segmentation datasets. Additionally, we introduce two techniques into the field of surface crack segmentation, previously not used there: Generating results using test-time-augmentation and performing a statistical result analysis over multiple training runs. The former approach generally yields increased performance results, whereas the latter allows for more reproducible and better representability of a methods results. Using those aforementioned strategies with our proposed encoder-decoder architecture we are able to achieve new state of the art results in all datasets.

*Keywords:* Crack Segmentation, Convolutional Neural Network, Deep Learning, Semantic Segmentation

## 1. Introduction

Regular usage, ageing and environmental conditions are some of the main factors which contribute to the deterioration of road surfaces. In combination with the ever expanding urbanization and development of public infrastructure, the surface area of roads grows and so does the need for maintenance to keep those structures in usable condition. One of the main defects that can appear

---

*Corresponding Author

*Email address:* `jacob.konig2@gcu.ac.uk` (Jacob König)

on road surfaces are cracks [1, 2]. If not treated within a reasonable timeframe these cracks can increase in size, impair the structural integrity of roads and pavements, facilitate further defects such as potholes as well as lead to material damage and required maintenance downtimes [2, 3]. Therefore it is important to regularly survey those surfaces and fix faults before they escalate. Due to the developments in information technology and image processing techniques it is now possible to enable the partly automation of the surface inspection process to detect faults such as cracks. Proposed systems can include a variety of sensors such as cameras or laser [2, 4, 5] as well different platforms on which those sensors are deployed. These may be groundborne, such as specialised vehicles [2] or airborne as in UAVs [5]. However, independent of the platform or the sensors used, a large amount of data is collected which needs to be analysed, either manually, which is unduly time consuming, or by the use of automatic algorithmic methods.
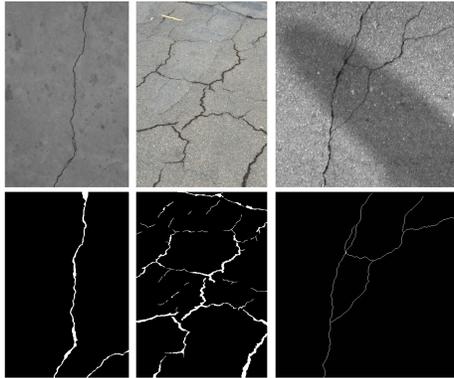


Figure 1: Sample crack images (top row) and their annotated ground truth labels (bottom row) from three datasets ([6, 7, 3]), showing some of the challenges of the crack segmentation task: non-uniform backgrounds, obstructions such as shadows, dirt and stains, differing crack widths as well as changing camera angles.

Figure 1 shows the diversity of surface cracks, they may traverse into the horizontal and vertical directions, appear in grid like patterns, vary in width along the same crack and appear atop of different surface structures. In addition to that, other anomalies may cover the appearance of cracks in images, including road markings, dirt patches, oil stains and vegetation, all of which are common occurrences. This makes this task increasingly more difficult when conducted manually, in addition to any subjectivity an annotator might have [1, 8].

In the computer vision domain, deep learning based convolutional neural networks (CNN) provide state of the art performance in a variety of tasks such as classification, detection and segmentation [9, 10, 11]. In some tasks they have even surpassed human performance [12]. Other work has also shown that when working on images with cracks, deep learning based methods outperform other techniques using handcrafted features or further machine learning approaches, making them well suited for the analysis of cracks in images [3, 13, 7, 5, 14].

In this work we follow the popular approach of using semantic segmentation to analyse surface images, automatically annotating an image on a pixel level and stating for each pixel if it either contains a crack or not. Other approaches have been proposed to classify whether images or patches contain cracks [15, 14] or locate cracks through creation of bounding boxes [16], however segmentation provides a better level of detail to assess the severity of cracks and allows for tracking of potential crack changes over time.

Encoder-decoder based CNN architectures have achieved significant results, from segmentation of city-scenes [17, 18], to medical image segmentation [19] as well as for crack segmentation [3, 13, 20, 21]. This is a popular approach as the encoder learns denser features and the decoder then increases the size of the features again for creation of a segmented output. Connecting features of the encoder and decoder parts in an architecture, not only at the bottleneck such as in SegNet[17], but at their interim stages generally leads to improved results. This is shown in state of the art results in the crack domain, such as in DeepCrack [3], which fuses multi-scale features at the different levels of the encoder-decoder architecture, FPCNet [13] which sums skipped feature maps in the decoder upsampling process or as in the architecture in [21], using concatenated shortcut connections between encoder and decoder as introduced in the U-Net [19] architecture for medical image segmentation.

Several recent encoder-decoder segmentation architectures have been designed around using an encoder backbone which has been derived from popular image classification architectures [7, 22, 23, 24]. As these encoders are designed for image classification they provide strong feature extraction capabilities and are therefore ideal for use as an encoder part in segmentation architectures.

Based on those observations we present an universal decoder design which can be added to a variety of popular feature-encoding backbones such as VGG [25], ResNet [26] and EfficientNet [10] to effectively segment surface cracks.

Whilst many of the previous works only report their results on a small amount of datasets [13, 21, 27, 7], we show that our decoder-design in conjunction with a pretrained backbone achieves high performance across a multitude of datasets.

In summary, our proposed optimized deep encoder-decoder methods consist of the following contributions:

1. We design a novel decoder part for encoder-decoder based CNN architectures to segment surface cracks in images. This decoder can be added to different encoder architectures, completing an U-Net like shape for the task of semantic segmentation. The components of this decoder architecture have been selected to optimize the performance on the task of crack segmentation.

2. The performance of this decoder is then evaluated on five relevant crack datasets. It is shown that using this decoder, in conjunction with an *EfficientNet B5* encoder backbone, new state of the art results are achieved in all datasets. To the best of our knowledge this is one of the first comparisons across a such a large number of crack segmentation datasets from

3

different works.

3. During evaluation, we apply two techniques which previously have not been used in the crack-segmentation domain; Firstly, we discover that test-time-augmentation, through resizing of images, yields an additional boost in performance and secondly, reporting of results using a single training run is error-prone and does not accurately present the performance of a model. We therefore propose to report results after performing a statistical analysis and do so by showing the performance of our model by reporting the average and standard deviation over ten training runs

The reminder of this paper is organized as follows: related work in regards to crack segmentation is reviewed in Section 2. In Section 3 we outline the details of our proposed encoder-decoder architecture and state the training configurations used. Section 4 describes the datasets and the experiments which have been carried out. This Section also reports the results and contains the experiments which have been carried out to justify the components of our architecture as well as our methods generalization abilities and limitations. This work is then concluded in Section 5.

## 2. Crack Segmentation Background

Making use of the growing capabilities in image capturing and processing, automatic crack detection and segmentation approaches have widely been studied in recent years [28].

### 2.1. Traditional Methods

Early approaches in crack segmentation exploited photometric- and geometrical characteristics of cracks, e.g. differing colors and shapes of cracks compared to their background [29]. These include crack segmentation through thresholding and mathematical morphology [30, 31, 32, 33, 34]. Further approaches make use of different filtering techniques like using Wavelet transform [35, 36, 37], Gabor filters for segmentation [1] as well as edge detection algorithms who have also been exploited for segmentation of cracks such as in [38] where the Canny edge detector [39] is used.

However, these early approaches are limited in their performance due to the varying appearance of cracks and the amount of tuning required to make them work on various datasets. Hence, machine learning methods have started to gain traction, providing the ability to learn based on their input features and increase the performance over those previous methods [40, 6].

In [6], Random Structured Forests are used to segment image patches and Support Vector Machines (SVM) or k-nearest neighbor (KNN) classifiers are used to classify whether this patch contains a crack. The work in [41] uses a SVM to segment cracks based on previously extracted graph based features.

4

## 2.2. Deep Neural Network Methods

Accompanied by the general rise of using deep learning based methods for computer vision tasks, several works have also shown that using deep learning based techniques on crack images outperform those previous methods by very large margins [3, 14, 42, 43, 27, 40].

The work in [40] compares a multitude of traditional crack segmentation methods with a CNN architecture to segment tunnel wall defects. This work shows that using CNN to segment defect areas outperforms previous methods such as SVM and KNN as well as classification trees.

In [43] a CNN architecture is introduced which classifies whether the central pixel in an image patch contains a crack. Its results are then averaged for each pixel in an image to generate a segmentation map. The work in [27] uses a similar approach with dilated receptive fields to create crack segmentation maps, in addition to using input augmentation and model-weight-sharing to predict the orientation of cracks. In [44] it is shown that an ensemble of CNN models can not only produce accurate segmentation maps but also provide crack-measurements. However, the drawback of this ensemble method is the lack of an end-to-end training approach.

Several recent deep learning architectures for crack segmentation are based on an encoder and decoder basis [42, 16, 3, 20, 13]. Inside the encoder, convolutional and downsampling operations create denser and spatially smaller feature maps which are then upsampled in the decoder to create an output segmentation maps, which matches the spatial dimensions of the input. In [42, 21] architectures based on U-Net [19] are presented. U-Net encompasses shortcut connections, between the encoder and decoder part, with the goal of retaining spatial features lost due to pooling. The DeepCrack architecture in [3] is based on SegNet [17], but it contains a multi scale fusion component, extracting and making use of features from multiple scales of the feature pyramid to generate the segmentation map. In comparison with several other baseline CNN segmentation architectures such as U-Net, SegNet and hollistically-nested edge detection (HED) [45], this architecture achieved superior results when trained on one, and evaluated on three custom datasets. FPCNet [13] combines an encoder-decoder architecture with a module in the bottleneck which uses dilated convolutions to extract features of multiple scale in feature maps before using an adaptive upsampling approach which adds encoder features to the feature maps of the decoder and channelwise attention through the squeeze-and-exitation [46] method before further processing. The work in [20] expands upon the HED method [45] adding a feature pyramid module [47] as well as a component to adaptively reweight features of different levels of those feature pyramid before generating a segmentation result. Another architecture, also called DeepCrack, in [7], uses a VGG backbone [25] and a deeply supervised [48] approach to upscale and fuse the feature maps from all levels of the backbone, before applying guided filtering [49] using the fused feature maps as well as a side output of the first convolution stage of the backbone, to create a segmentation output. The crack segmentation method in [22] uses a ResNet backbone

feature encoder and decodes the features maps through a multi-dilation module which uses dilated-convolutions to extract crack-features from different scales. It also employs multi-scale fusion, merging interim sparse feature maps with dense ones from later layers.

## 3. Crack Segmentation Architecture

The following section introduces our architecture which is built following an U-Net [19] like, encoder-decoder, shape and can make use of different pretrained feature encoders. To those encoders we add a decoder part, optimized for crack segmentation performance. An illustration of our architecture is shown in Figure 2.
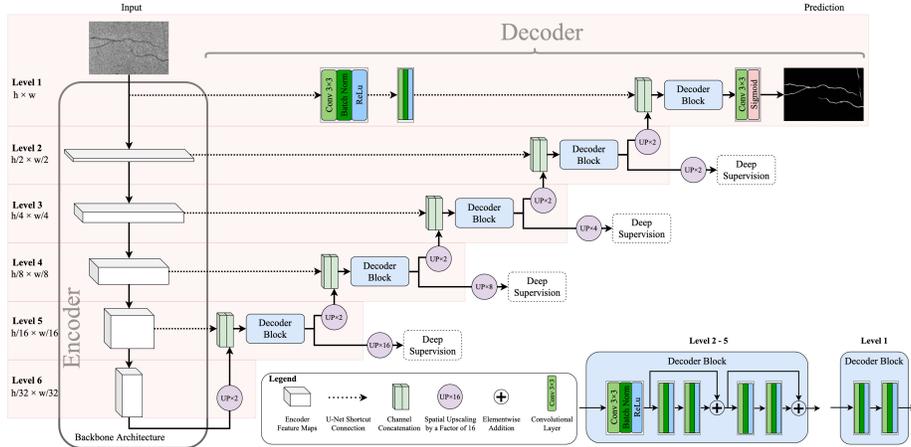


Figure 2: Our proposed Crack Segmentation Architecture. Different encoder backbones can be used from which the U-Net shortcut connections are extracted by using the feature maps right before each downsampling operation. As illustrated, the decoder block in the first level differs from the remaining ones. Upsampling between levels is done through nearest-neighbor upsampling, whereas each upsampling operation before a deep supervision operation is chosen to be bilinear upsampling. The filters in the decoder are fixed as: [16, 32, 64, 128, 256], increasing with the level size.

### 3.1. U-Net Basis

U-Net based architectures contain an encoder-decoder like structure with skip-connections. Generally an input is propagated through successive convolutional, activation and pooling layers, reducing its feature size up to a bottleneck, therefore forming the encoder. In the following decoder, instead of pooling layers, upsampling layers are used to increase the feature size back to the spatial dimensions of the input. These sequences of using convolutions, activations and pooling allows the network to extract more condensed features, with the drawback of reducing the spatial size of the feature map. To counteract the loss

of spatial information after this encoder part, U-Net therefore introduces skip-connections. These connections combine interim features of the encoder part with the output of the upsampling operation of the same spatial size through channel-concatenation, hence reintroducing spatial features. This is followed by feeding those concatenated features into the decoder layer-sequences.

### 3.2. Decoder Design

Our proposed decoder can be added to a variety of encoder backbones. However, encoders which aim to be used with our proposed decoder, need to have six levels, each distinguishable by their different spatial size. These are addressed with level numbers $l=1$ to $l=6$, each having half the spatial size of its previous one. Hence, an encoder needs to have five downsampling layers.

For each but the last level of the encoder, we add a decoder block. The decoder blocks are thus situated in levels $l=1$ to $l=5$ and are connected to the encoder through U-Net style skip connections. To increase the spatial size inbetween different levels of the decoder, we make use of nearest neighbor upsampling. The input into the decoder sequence is the upsampled output of the sixth level of the encoder backbone. After each upsampling operation, the upsampled feature maps are concatenated with the output of the encoder at the specific level before being propagated to the decoder block.

Some backbones, such as ones based on EfficientNet or ResNet, apply a downsampling convolution right after their input layer. This means that the resulting output at level $l=1$ of the encoder equals that of the input image.

To therefore still extract features at this spatial size, we expand the skip connection at that first level, to include two $3 \times 3$ convolutional layers (Conv), each followed by Batch Normalization (BN) and a ReLU activation on an input.

The decoder blocks in level $l=2$ to $l=5$ are residual blocks, meaning they contain residual connections inbetween convolutions. Adding those residual connections for semantic segmentation tasks in convolutional blocks has shown to improve the performance in U-Net based architectures [50, 51]. Every residual block applies two $3 \times 3$ Conv-BN-ReLU sequences on an input. This input is then subsequently added to the output of the Conv-BN-ReLU sequence, thus creating a residual connection. Each decoder block in those levels contains a standard Conv-BN-ReLU sequence followed by two residual blocks. The number of filters used in all decoder blocks is based on its level $l$ in the architecture, leading to the following number of filters per layer: $[16_{l=1}, 32_{l=2}, 64_{l=3}, 128_{l=4}, 256_{l=5}]$. The decoder block at level $l = 1$ applies two Conv-BN-ReLU sequences. A final segmentation output is then generated through running the resulting feature maps of the decoder block at level $l=1$ through a $3 \times 3$ convolutional layer with 1 filter and a Sigmoid activation function.

To enforce the learning of more robust feature this architecture also deeply supervises [48] the outputs of the decoder blocks on levels $l=2$ to $l=5$. This is achieved through a $1 \times 1$ convolution with one filter, followed by a Sigmoid activation and bilinear upsampling to the original spatial size.

### 3.3. Encoder-Backbone Architectures

Commonly, the encoders in encoder-decoder structures are either tailored and designed for the specific image segmentation task [13, 19] or use the same architectures to ones that have been used for image classification tasks without the fully connected components [17, 7, 24]. In this work we focus on the latter, by considering three different types of backbone architectures: VGG, ResNet and EfficientNet. The weights, having been pretrained on ImageNet, provide the encoder with a strong feature extraction baseline and enable a faster training convergence. For all backbones, we omit the fully connected components and use the output of the last operation before the classification part as the output.

The VGG architecture contains basic successions of convolutional and activation layers, with max pooling layers used to reduce the spatial feature sizes. There are several variants such as *VGG 16*, with 16- , or *VGG 19* with 19-weight layers. It has successfully been used as an encoder backbone for general segmentation architectures such as SegNet as well as for crack segmentation in the DeepCrack architectures in [7, 3].

ResNet is an architecture design which makes use of skip-connections, allowing for networks with more layers as well as better gradient flow. Several sizes of those networks have been introduced ranging from variants with 26 layers (*ResNet 26*) up to variants with over 1000 layers (*ResNet 1001*). Its variants are also commonly used as encoder-backbones for segmentation [23, 52, 22].

The EfficientNet architecture introduced a new way of designing CNNs. It is proposed that instead of scaling the depth, width and image-resolution separately, they should be scaled equally. Based upon this proposal, a baseline architecture (*EfficientNet B0*) which was found using neural architecture search [53], is scaled up seven times (*EfficientNet B1 to EfficientNet B7*). The consequently resulting scaled up architectures perform much better and faster on ImageNet [54] and other datasets such as CIFAR [55] whilst using much less parameters in comparison to other competitive approaches. EfficientNet based architectures have also been previously used for segmentation, such as in [24].

Our results are reported using an *EfficientNet B5* backbone. In later experiments we validate the performance of this and other backbones and justify our choice.

### 3.4. Training Configuration

Our proposed architecture is trained on patches of size $288 \times 288$ pixels similar to the implementation in [13]. Compared to other segmentation domains, crack segmentation can be seen as a binary classification problem, as pixels can only be assigned a probability of containing a crack or not. Let $X$ be a training dataset consisting of $N$ samples $x$, $X = \{x_1, x_2 \ldots x_N\}$ as well as $Y$ being the respective ground truth dataset to $X$ consisting of $N$ samples $y$, $Y = \{y_1, y_2 \ldots y_N\}$ then each ground truth pixel $i$ of any given sample $y$ is $y_i \in [0, 1]$.

Recall that this architecture contains a total of five outputs, four interim deeply supervised outputs at level $l = 2$ to $l = 5$ and the final segmentation output. During the training process, the sum of the binary cross entropy and

the Dice loss is applied on each output and the total sum of the losses of all outputs is used as the training loss. The cross entropy loss can be used for semantic segmentation, however it may be affected when the classes are highly unbalanced such as in the case of crack- to no-crack pixels, where a trained architecture might prefer to segment background pixels. A common approach to counteract this is issue is to either weight classes differently, or use a weighted cross entropy loss [7, 45]. We however add the Dice loss [56], which measures the overlap of the prediction and the ground truth so that the total loss takes into account both the segmented shape as well as smaller details. The loss $L$ for a prediction $\hat{y}$ consisting of $P$ pixels at a specific output of the network is therefore calculated as:

$$L = -\sum_i^P y_i \big( \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i) \big) + \frac{2 \sum_i^P y_i \hat{y}_i + 1}{\sum_i^P y_i + \sum_i^P \hat{y}_i + 1} \tag{1}$$

## 4. Experiments and Results

This section describes the experiments which have been carried out on different datasets as well as further studies on how the different components of our method affect its performance.

### 4.1. Datasets

This work makes use of five datasets, CrackForest [6], DeepCrack-DB [7], CrackTree260, CRKWH100 and Stone331, in which the latter three were introduced in [3].

- The *CrackForest* dataset[1] (CFD) consists of a total of 118 images of size $480{\times}320$ pixels. The images have been taken from road surfaces in Beijing. We omit one image as it contains a wrong ground truth[2]. This dataset is then split into 71 training, and 46 testing images[3] similar to [43, 27, 13]. The images are supplied in RGB, however due to the small dataset size and little difference in color channels we convert them to grayscale.

- *CrackTree260* (CT260) is a dataset containing 260 grayscale road pavement images of different sizes ($800 \times 600$ and $960 \times 720$ pixels). Similar to the approach in [3] this dataset is exclusively used for training the architectures.

- The *CRWKH100* dataset consists of 100 grayscale images road pavement images of size $512{\times}512$ pixels, all of which are used exclusively for testing.

---

[1]Available at `https://github.com/cuilimeng/CrackForest-dataset`

[2]The filename of the omitted image is *042.jpg*

[3]For reproducibility we use images *001.jpg* to *072.jpg* for training and the remainder for testing

- *Stone331* consists of 331 grayscale images of stone surfaces, again of size $512 \times 512$ pixels. However, not the whole surface of the image contains the stone, therefore masks are provided to exclude predictions outside the stone area. This dataset is also used exclusively for testing.

- A total of 537 RGB images are contained in *Deepcrack-DB*[4][7].The dataset split is given as 300 training-, and 237 testing images, all of size $544 \times 384$ pixels.

In [3] another dataset, CrackLS315 is used, however the authors were not able to provide us with this dataset. CT260, CRKWH100 and Stone331 can be seen as one large dataset, as algorithms trained on CT260 are tested on CRKWH100 and Stone331. It is observed that the CT260 and CRKWH100 images are of similar-looking nature, whereas Stone331 images appear different. This leads to generally worse performance in Stone331 compared to CRKWH100 (as our following experiments indicate) however it may also show the generalization performance of algorithms.

As it can be seen, the number of images in each dataset is limited compared to larger datasets such as ImageNet. Therefore we perform augmentation on the training images. At first we enhance the number of images by rotating each training image for $[0°, 90°]$ as well as mirroring it across the $[x, y, y{=}x, y{=}{-}x]$ axis. This means for each training image 8 variations of it are created, equaling the Dihedral group D4.

As previously mentioned our architecture is trained on images of size $288 \times 288$. This is achieved using an augmentation policy of random cropping, brightness and contrast shift, the addition of noise as well as random zooming in or zooming out, which results in an image of the selected size. Due to the fully convolutional nature of our architecture, as it does not contain any fully connected layers with a fixed number of connections which may restrict images sizes, evaluation can then be performed on the full sized testing images. The detailed augmentation policy is outlined in Algorithm 1. Sample images generated by our augmentation policy are shown in Figure 3.

*4.2. Metrics*

In this work, results are reported using three metrics. All of them are based on the harmonic mean: the $F1$ score of Precision $PR = \frac{TP}{TP+FP}$ and Recall $RE = \frac{TP}{TP+FN}$. These can be calculated based on True Positive $TP$, False Negative $FN$ and False Positive $FP$ predictions. The $F1$ is calculated as $F1 = 2 \cdot \frac{PR \cdot RE}{PR+RE}$. In the crack segmentation domain the *F1* metric is applied three ways, all of which make use of an oracle. This oracle contains the solutions as well as all predictions with the aim to chose the best possible result for each metric.

---

[4]Due to a naming overlap with the work in [3] we refer to this dataset as *Deepcrack-DB* as the authors have not given it a name

---

**Algorithm 1** Detailed augmentation policy for the training process of our architecture.

---

*The bold character in brackets states if augmentation is applied to:*
- *only the image* **(i)**
- *same augmentation to the image and mask* **(i,m)**

*p states the probability with which each augmentation step is applied*

**FOR EACH training *image* (i) and corresponding *mask* (m) :**

1. Random crop of size $288 \times 288$ pixels **(i,m)**
2. With $p$=0.5 apply Random brightness shift in range $[-10\%, +10\%]$ **(i)**
3. With $p$=0.5 apply Random contrast shift in range $[-10\%, +10\%]$ **(i)**
4. Apply one of:
   (a) Random Additive Gaussian Noise with a standard deviation sampled in range of $[0, 2.55]$ **(i)**
   (b) Random Multiplicative Noise in range $[75\%, 125\%]$ **(i)**
5. With $p$=0.3 apply one of:
   (a) Random spatial reduction of the image in range $[0.75, 1]$ with mirror padding back to $288 \times 288$ pixels **(i,m)**
   (b) Random sized crop, minimum size of $144 \times 144$ pixels with nearest-neighbor upsampling back to $288 \times 288$ **(i,m)**
6. **OUTPUT (i,m)**

---

1. Reporting the aggregated $F1$ metric when an oracle extracts the best $F1$ score of all confidence thresholds of each image yields the *OIS* (Optimal Image Scale) metric [57].

2. The *ODS* (Optimal Dataset Scale) metric reports the aggregation of $F1$ scores from all images when an oracle has set the best possible fixed threshold across all images in the dataset [57].

3. A cumulative ODS, *cODS*, as used in [42, 43, 13, 7], reports the $F1$ when $PR$ and $RE$ have been calculated across the whole dataset and not as an aggregation of results on each image. We assume that when no confidence threshold cutoff is reported in those works, the one that achieves the highest results extracted through an oracle has been used. Therefore we are referring to this metric as *cumulative* ODS.

The formulas for each metric are derived as follows:

$$OIS = \frac{1}{N_{img}} \sum_{i}^{N_{img}} max\Big\{ F1_t^i : \forall t \in \{0.01, ..., 0.99\} \Big\} \tag{2}$$

$$ODS = max\Big\{ \big\{ \frac{1}{N_{img}} \sum_{i}^{N_{img}} F1_t^i \big\} : \forall t \in \{0.01, ..., 0.99\} \Big\} \tag{3}$$

$$cODS = max\Big\{ F1_t^d : \forall t \in \{0.01, ..., 0.99\} \Big\} \tag{4}$$
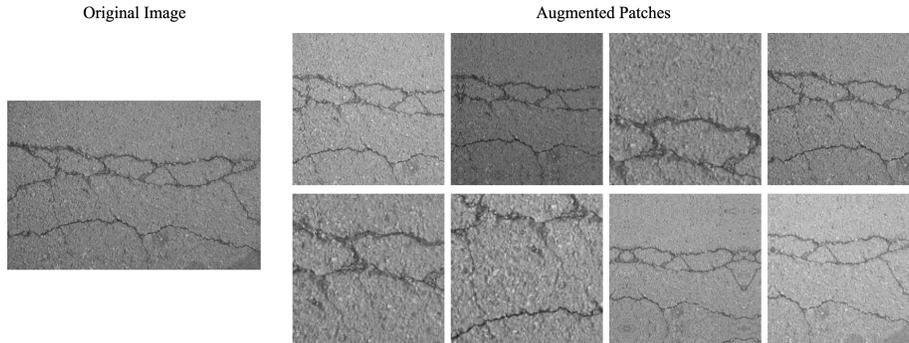
11

Figure 3: Augmentation image samples when using an image of the CFD dataset. This augmentation policy creates quadratic image patches of 288 pixels size. The brightness and contrast changes emulate different backgrounds and crack-colourings, whereas the size-augmentation exposes the network to more crack data of different widths during training.

with $t$ denoting the confidence threshold at which a predicted pixel is counted as a $TP$, $N_{img}$ the number of images in the testing dataset $d$ and $i$ each specific image.

We only report the $cODS$ metric to be able to compare with other implementations. We propose to use the $ODS$ and $OIS$ metrics in further crack segmentation works, as $cODS$ does not accurately represent the real performance of an algorithm. Due to the way the $F1$ score is calculated there, it may skew results as a small proportion of images in a dataset may contain the majority of crack pixels.

As in previous works for crack segmentation [42, 3, 13, 44, 21, 27], in the experiments on the CFD, CT260, CRKWH100 and Stone331 datasets, $TP$ pixels are accounted for if they lie within two pixels of the corresponding ground truth, leading to a relaxed distance threshold. This distance threshold also changes the calculation of $FN$ and $FP$. A $FP$ pixel is counted if it is a wrong prediction farther away from a ground truth pixel than a two-pixel distance and a $FN$ pixel is counted as a missed prediction if the ground truth is farther away than two pixel from a prediction.

This was adopted as the labelling of cracks may be slightly inaccurate. The experiments on DeepCrack-DB do not use such two-pixel threshold to be able to provide an accurate comparison with [7].

### 4.3. Experimental Configuration

Our network architecture is trained using stochastic gradient descent as the optimizer using a degrading learning rate following the formula $lr_e = lr_{e=0} * 0.96^e$ with $lr_e$ denoting the learning rate at a specific epoch $e$ and $lr_{e=0}$ denoting the initial learning rate at the first epoch $e{=}0$ with a value of 0.01. The weights in all layers are initialized using the method in [12]. The implementation of this work was achieved using the Keras-API [58] with the Tensorflow [59] backend and training and evaluation were performed on a NVIDIA Titan XP GPU. Parts

of our architecture implementation make use of code from [60]. On all training datasets the architecture is trained for 120 epoch, using a batch size of 8 and a weight decay of 1e-5. As the works introducing those datasets did not provide a subpart of the dataset for validation, we estimated the hyperparameters and training settings based on previous works [13, 3] and always used the weights of the last epoch for evaluation. The results in the following section are reported without and with test-time-augmentation (TTA). Our results show that similar to the findings in [61, 52] scaling the images to different sizes and aggregating the different predictions improves the segmentation results significantly. Due to the different image sizes, the CFD testing images are scaled in the factors [0.6, 0.8, 1.0, 1.2, 1.4] and in CRKWH100/Stone331 to [0.5, 0.75, 1.0, 1.25, 1.5]. These scaling factors differ as the spatial dimensions of the image still need to be dividable by $2^5$, based on the five reductions of the spatial size of the feature maps by a factor of two in the encoder going from an input image of spatial size $h \times w$ to $\frac{h}{2^5} \times \frac{w}{2^5}$ in level 6 of the encoder.

On the DeepCrack-DB testing set we also use TTA, however due to the height-to-width ratio and the limitation of being dividable by $2^5$ we scale the images to the fixed sizes of $[288 \times 192 , 416 \times 288, 544 \times 384, 672 \times 480, 832 \times 576]$. Thereby each spatial dimension is scaled to the closest possible divisor of the range [0.5, 0.75, 1.0, 1.25, 1.5].

In addition to that we also found that repeatedly running each experiment, even with fixed data sampling seeds, do not yield consistent results due to non-deterministic GPU operations. This is similar to the findings in [62]. Therefore we report our results as averages with their standard deviation over 10 runs. Across all different tested configurations, each experiment has the same fixed seed for selection and augmentation of training data, for each specific run. As we will see in the following sections, averaging those results is important as outliers may detract from the actual performance of those algorithms, therefore not representing the general model performance.

The approximate training run time of our proposed method using an *EfficientNet B5* backbone for a single run with the aforementioned configuration is 1:20h for CFD, 4:50h for CT260 and 5:30h for DeepCrack-DB.

Sample results of using our architecture are shown in Figure 4. We note that due to our reporting approach using averaged results and the results lying in range of >85% *OIS/ODS*, visible differences between our approaches and other implementations are very small, hence we chose to omit those visual comparisons and stick to reporting using the chosen metrics.

### 4.4. Results

#### 4.4.1. Results on CFD

The results on the CFD dataset are shown in Table 1. As it can be seen previous results report their performance in the *cODS* metric. We show that our architecture achieves a slightly better average performance of 0.04% when run without TTA. When testing the images on multiple scales we can report a total increase of 0.10% in the *cODS* metric. We also report the results for

the *OIS* and *ODS* metrics, in which it also can be seen that TTA improves the results.

Table 1: Results on CFD. Our results are reported as (Avg)±(Stdev) over 10 runs.

| Method | *OIS* | *ODS* | *cODS* |
|---|---|---|---|
| Fan *et al.* [43] | - | - | 92.44 |
| Fan *et al.* (ensemble) [44] | - | - | 95.33 |
| Inoue and Nagayoshi [27] | - | - | 95.70 |
| FPCnet [13] | - | - | 96.93 |
| Ours, no TTA | 97.36±0.26 | 96.88±0.32 | 96.97±0.23 |
| **Ours + TTA** | **97.64±0.26** | **96.91±0.36** | **97.03±0.30** |

### 4.4.2. Results on CT260, CRKWH100, Stone331

Only one other work [3] makes use of the CT260, CRKWH100 and Stone331 datasets and reports their results. Our evaluation results are shown in Table 2. Without TTA our method performs on average 2.03% on the *OIS* and 2.61% better on the *ODS* metric on CRKWH100, achieving 93.20% and 92.66% respectively. Adding TTA then increases the average results by a larger margin to 94.25% *OIS* and 93.34% *ODS*. On the Stone331 evaluation dataset, without using TTA, we also report increased results, achieving 87.58% *OIS* and 87.36% *ODS*, in comparison with 87.51% *OIS* and 85.59% *ODS* of the previous best results. However, if we then enable the use of TTA on this dataset we achieve 93.30% *OIS* and 92.17% *ODS*, outperforming the previous approach by a margin of 5.79% *OIS* and 6.58% *ODS*.

Table 2: Results on CRKWH100 and Stone331 when trained on CT260. Our results are reported as (Avg)±(Stdev) over 10 runs.

| Method | CRKWH100 | | Stone331 | |
|---|---|---|---|---|
| | *OIS* | *ODS* | *OIS* | *ODS* |
| DeepCrack [3] | 91.17 | 90.05 | 87.51 | 85.59 |
| Ours, no TTA | 93.20±0.18 | 92.66±0.18 | 87.58±0.91 | 87.36±0.83 |
| **Ours + TTA** | **94.25±0.13** | **93.34±0.07** | **93.30±0.52** | **92.17±0.48** |

### 4.4.3. Results on DeepCrack-DB

To the best of our knowledge, only the authors of [7] have used and reported their results on the DeepCrack-DB utilizing their proposed train-test split, hence we only compare against this other work. Our results on this dataset, as shown in Table 3 when not using TTA indicate a slight improvement over the previous approach by 0.79% in the *cODS* metric. However, when TTA is used this margin increases to an improvement of 1.53% in the *cODS* metric. Using the *OIS* and *ODS* metric, we achieve 88.86% and 85.39% respectively. Whilst in the previous three datasets a two pixel threshold is used to count *TP* pixel, in this dataset the results are reported using a zero pixel threshold, meaning a *TP* pixel only counts if it lies directly on the ground truth location.

Table 3: Results on DeepCrack-DB. Our results are reported as (Avg)±(Stdev) over 10 runs.

| Method | *OIS* | *ODS* | *cODS* |
|---|---|---|---|
| Liu *et al.* [7] | - | - | 86.50 |
| Ours, no TTA | 87.43±0.45 | 85.04±0.48 | 87.29±0.34 |
| **Ours + TTA** | **88.86±0.28** | **85.39±0.44** | **88.03±0.42** |

### *4.5. Analysis of Different Architecture Components*

To achieve the optimal results using our proposed encoder-decoder design, we performed several side-studies, analyzing the impact of different components that make up our architecture.

Unless stated otherwise, all the studies in this section have been performed using an *EfficientNet B3* encoder backbone. This provides a good trade-off in representative performance, model-size and training time and we assume that the results showing increases or decreases in performance also carry over to other backbones. Additionally the baseline model variations in the following studies make use of bilinear upsampling, transfer-learning, deep supervision unless it is explicitly changed.

The results are reported after using test-time-augmentation and being averaged across 10 training runs as (Avg)±(Stdev), in which each training run from one to ten uses the same seed across all models.

### *4.5.1. Upsampling method*

Encoder-decoder based architectures can use multiple methods to increase the size of the feature maps in the decoder stages. Well established methods are bilinear interpolation [63, 23], upsampling using transposed convolutions [7, 64, 65] or upsampling using max-pooling indices [17, 3]. Our analysis compared the effects of using bilinear interpolation, nearest neighbor upsampling, as well replacing the upsampling layer with a transposed convolution of kernel size $4 \times 4$. To not significantly increase the number of parameters in this transposed convolution, we utilize a bottleneck approach. The number of channels going in to this transposed convolution is set to one fourth of the number of input channels into the general upsampling operation and is achieved through a $1 \times 1$ convolution. The transposed convolution is then carried out with the number of filters matching the number of channels of its input and is followed by increasing the number of channels again through another $1 \times 1$ convolution, whose number of filters matches that of the input which is upsampled. Due to not all encoder architectures making use of max-pooling, comparing the use of pooling indices for upsampling has been omitted.

Table 4 shows that across all datasets the nearest neighbor and bilinear upsampling method achieve comparable performance. Upsampling using transposed convolutions performs about equally well on CRKWH100 and CFD, however on Stone331 it trails behind the former upsampling methods. It is to note that transposed convolutions incur a computational overhead as well as extra parameters to the architecture.
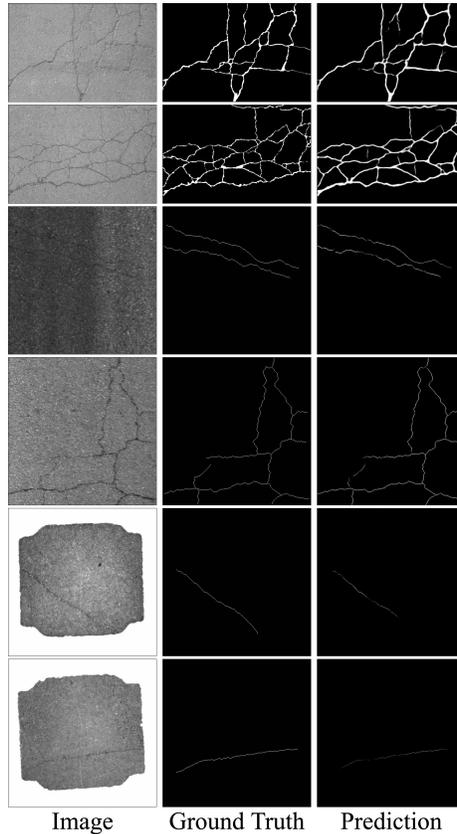
15

| Image | Ground Truth | Prediction |

Figure 4: Sample results of using our architecture on crack images from different datasets.

Due to inducing the smallest amount of extra computation, in comparison to the other methods, as well as achieving a slightly higher performance on both metrics in two separate datasets we chose to use nearest neighbor upsampling and suggest to also use this when utilizing this method on further datasets.

### 4.5.2. Transfer Learning and Deep Supervision

We compare training our proposed encoder-decoder architecture from scratch against utilizing pretrained weights in the encoder. In this transfer learning approach, the weights of the EfficientNet implementation that has been pretrained on the ImageNet [54] dataset for classification, are used. In addition to that we also show the impact that deep supervision has on the performance of the model. As our results in Table 5 indicate, using the model without deep supervision yields slightly worse performance across all datasets. When comparing using a transfer learning approach with initializing the encoder weights from scratch it is shown that pretrained encoder-weights increase the performance significantly in most of the datasets but DeepCrack-DB, where only a minor increase

Table 4: Results showing the impact of using different upsampling strategies within our proposed decoder. The column header in italic highlights the method chosen in our architecture, whereas bold results indicate the highest performance in a particular category.

| Dataset | Metric | *Nearest Upsampling* | Bilinear Upsampling | Transposed Convolution |
|---|---|---|---|---|
| CT260 + CRKWH100 | OIS | 93.93±0.10 | 93.94±0.18 | **93.97±0.18** |
| | ODS | 93.06±0.12 | 93.12±0.18 | **93.12±0.13** |
| CT260 + Stone331 | OIS | 89.30±3.36 | **89.72±2.16** | 87.98±4.20 |
| | ODS | 88.48±3.17 | **88.80±2.16** | 86.47±4.73 |
| CFD | OIS | **97.73±0.19** | 97.69±0.13 | 97.68±0.15 |
| | ODS | **97.07±0.28** | 97.00±0.19 | 97.05±0.21 |
| DeepCrack-DB | OIS | **88.91±0.16** | 88.90±0.19 | 88.88±0.18 |
| | ODS | **85.49±0.33** | 85.42±0.46 | 85.36±0.30 |

is achieved. We argue that this is due to the strong feature extraction baseline which backbones trained on ImageNet provide.

Table 5: Results showing the impact of using deep supervision (DSV) and transfer learning (TL). The column header in italic highlights the method chosen in our architecture, whereas bold results indicate the highest performance in a particular category.

| Dataset | Metric | *Ours + DSV + TL* | Ours without DSV | Ours without TL |
|---|---|---|---|---|
| CT260 + CRKWH100 | OIS | **93.93±0.10** | 93.78±0.39 | 92.24±0.88 |
| | ODS | **93.06±0.12** | 92.74±0.42 | 91.35±1.08 |
| CT260 + STONE331 | OIS | **89.30±3.36** | 81.76±5.93 | 62.01±39.74 |
| | ODS | **88.48±3.17** | 80.28±6.18 | 61.00±39.50 |
| CFD | OIS | **97.73±0.19** | 97.28±0.11 | 95.49±0.95 |
| | ODS | **97.07±0.28** | 96.49±0.22 | 94.54±1.14 |
| DeepCrack-DB | OIS | **88.91±0.16** | 88.64±0.17 | 88.34±0.15 |
| | ODS | **85.49±0.33** | 84.66±0.36 | 84.81±0.29 |

*4.5.3. Choice of Feature Encoder Backbone*

The impact on performance of utilizing our proposed decoder with a variety of different encoder backbones can be seen in Table 6. Here, we use a selection of popular image classification backbones, whose weights have been pretrained on ImageNet. By comparing EfficientNet *B1 - B5* backbones it is shown that using larger EfficientNet backbones increases results on the CRKWH100 and Stone331 validation datasets, whereas similar performance is achieved on CFD and DeepCrack-DB. When comparing the performance of the largest, *EfficientNet B5* backbone with a *VGG 19* or *ResNet 50* backbone, the EfficientNet based one generally performs better in all but, the CFD, where using a *ResNet 50* encoder yields a slightly higher performance.

Our choice to report results in Tables 1, 2 and 3 using an *EfficientNet B5* backbone is due to the large increase in performance in the CRKWH100 and Stone331 validation datasets and similar performance to the best results by other backbones in the other datasets. Comparing the performance of models with similar number of parameters and different backbones (e.g. *EfficientNet B4*

and *VGG 19* as well as *EfficientNet B5* and *ResNet 50*) it can be seen that the models with EfficientNet based backbones generally perform better in three out of the four datasets. We also report a high fluctuation in results when trained on CT260 and tested on Stone331 when using the smallest EfficientNet backbone, *B0*. However, when it is tested on CRKWH100 there are more consistent results. We attribute this to overfitting the small EfficientNet backbone on the training dataset and assume that it is because the CRKWH100 test-set and the CT260 training set are of similar nature, whereas the Stone331 dataset appears different as described in section 4.1.

Table 6: Results showing the impact of using different encoder-backbones with our decoder. The backbone in italic highlights the method chosen in our architecture, whereas bold results indicate the highest performance in a particular category.

| Backbone | #Model Params | CT260 + CRKWH100 | | CT260 + Stone331 | |
|---|---|---|---|---|---|
| | | *OIS* | *ODS* | *OIS* | *ODS* |
| EfficientNet B0 | 12.5M | 90.12±1.49 | 88.55±1.65 | 26.46±22.53 | 24.64±21.85 |
| EfficientNet B1 | 15.0M | 92.90±1.32 | 92.03±1.45 | 87.43±3.44 | 86.48±3.35 |
| EfficientNet B2 | 16.7M | 93.50±0.28 | 92.56±0.21 | 89.32±2.18 | 88.53±1.99 |
| EfficientNet B3 | 20.2M | 93.93±0.10 | 93.06±0.12 | 89.30±3.36 | 88.48±3.17 |
| VGG 19 | 26.7M | 93.11±0.19 | 92.59±0.25 | 91.56±0.67 | 90.05±0.94 |
| EfficientNet B4 | 28.1M | 93.70±0.16 | 93.05±0.21 | 91.09±2.25 | 89.96±3.17 |
| ResNet 50 | 35.0M | 93.60±0.06 | 93.16±0.13 | 86.17±1.46 | 86.04±1.45 |
| *EfficientNet B5* | 39.8M | **94.25±0.13** | **93.34±0.07** | **93.30±0.52** | **92.17±0.48** |
| Backbone | #Model Params | CFD | | DeepCrack-DB | |
| | | *OIS* | *ODS* | *OIS* | *ODS* |
| EfficientNet B0 | 12.5M | 97.64±0.10 | 96.98±0.21 | 88.88±0.18 | 85.51±0.50 |
| EfficientNet B1 | 15.0M | 97.60±0.16 | 96.90±0.15 | 89.03±0.17 | **85.61±0.31** |
| EfficientNet B2 | 16.7M | 97.75±0.07 | 97.22±0.15 | 89.02±0.17 | 85.34±0.40 |
| EfficientNet B3 | 20.2M | 97.73±0.19 | 97.07±0.28 | 88.91±0.16 | 85.49±0.33 |
| VGG 19 | 26.7M | 96.53±0.29 | 95.75±0.31 | 88.10±0.15 | 84.65±0.22 |
| EfficientNet B4 | 28.1M | 97.86±0.10 | 97.17±0.14 | **89.05±0.17** | 85.49±0.38 |
| ResNet 50 | 35.0M | **97.92±0.09** | **97.37±0.18** | 88.68±0.17 | 84.73±0.39 |
| *EfficientNet B5* | 39.8M | 97.64±0.26 | 96.91±0.36 | 88.86±0.28 | 85.39±0.44 |

### 4.6. Impact of Test-Time-Augmentation

As shown in Sections 4.4.1, 4.4.2 and 4.4.3 the application of TTA leads to a mixed performance improvement ranging from 0.24% OIS on CFD to 5.72% OIS on Stone331. We hypothesize that the difference in dataset appearances as well as crack sizes may be causing this variation. Whereas CFD has a high diversity in crack sizes, a majority of the cracks in Stone331 are only a single pixel wide. It is assumed that TTA facilitates a better segmentation of very thin cracks as enlarging them may make them easier to detect for the algorithm.

### 4.7. Generalization Abilities

The results in the previous sections show, that the model achieves high performance on all datasets tested. Additionally, to gain further insight into the model performance, we have also studied the generalization abilities of our proposed encoder-decoder model. We show the results in Table 7. Here, we

used the models which have been trained on one specific dataset and evaluate them on all other testing-datasets.

Overall our model achieves high generalization performance on all but three configurations: Training on CFD or DeepCrack-DB and testing on Stone331 leads to a low performance and training on CT260 and testing on DeepCrack-DB also leads to a low generalization performance in both metrics. For the first two configurations we attribute this low performance on the difference in background textures inbetween the CFD/DeepCrack-DB and Stone331 datasets. For the third configuration we assume that the crack widths in CT260, with a large majority of cracks only being a single pixel wide, differs too much from the data appearing in DeepCrack-DB, where the crack width can span up to several dozen pixels.

Therefore, to achieve the best possible performance using this model, the training data should include cracks of different widths, ranging from one- to several pixels width, and a large diversity in backgrounds textures.

Table 7: Generalization results inbetween different datatsets. Cells with gray background indicate the original train-test dataset split, whereas cells with a white background highlight the generalization results.

| Test / Train | CRKWH100 | | Stone331 | |
|---|---|---|---|---|
| | $OIS$ | $ODS$ | $OIS$ | $ODS$ |
| CT260 | 94.25±0.13 | 93.34±0.07 | 93.30±0.52 | 92.17±0.48 |
| CFD | 94.02±0.80 | 89.08±1.55 | 36.79±11.42 | 28.13±10.35 |
| DeepCrack-DB | 89.37±1.25 | 85.44±1.28 | 15.97±2.98 | 10.98±1.83 |

| Test / Train | CFD | | DeepCrack-DB | |
|---|---|---|---|---|
| | $OIS$ | $ODS$ | $OIS$ | $ODS$ |
| CT260 | 96.12±0.56 | 96.04±0.58 | 28.86±0.82 | 28.86±0.82 |
| CFD | 97.64±0.26 | 96.91±0.36 | 79.70±3.44 | 75.50±3.49 |
| DeepCrack-DB | 94.02±0.78 | 93.81±0.75 | 88.86±0.28 | 85.39±0.44 |

*4.8. Segmentation Faults*

To outline the limitations of our proposed method we have included samples where our algorithm failed to generate adequate segmentation maps. These are shown in Figure 5. Samples were selected as they have the the worst OIS scores (averaged over all 10 runs) per dataset and give some insight into where our algorithm fails to segment.

As it can be seen on the first DeepCrack-DB sample, the algorithm has difficulties distinguishing very dark linear areas from cracks. Additionally, the second sample shows low performance on an image where the crack has very low contrast compared to the background. In the following two CRKWH100 samples, this method achieves a low performance due to the cracks being bright, which is in contrast to the majority of crack samples in the training set appearing darker. The Stone331 samples show no crack-detection at all, leading to an $OIS$ score of 0 on those images. We argue that is because the cracks appear very similar to their background. This problem also appears in the CFD samples, where crack regions which are very similar to their background are barely detected.

19

Generally, the limitations of this method lie in the detection of cracks that appear very different from the ones that appear in the training dataset, as well ones which have a low contrast between the background and the cracks. Both of those limitations can be mitigating by providing more diverse training data, as exposure to more different samples will lead to better testing-performance.
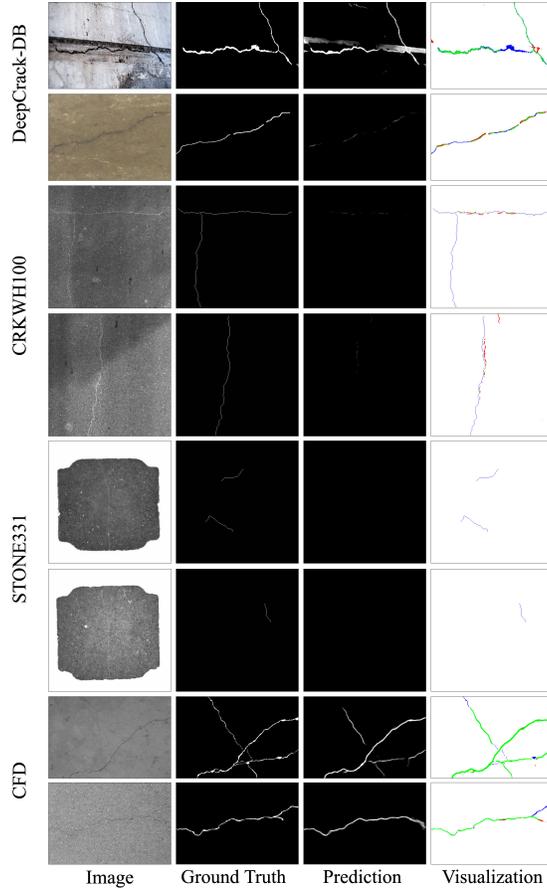


Figure 5: Sample images where our proposed model does not generate adequate segmentation results. Two images from each dataset (DeepCrack-DB, CRKWH100, Stone331 and CFD, from top to bottom) are shown with their respective ground truth, a sample prediction of our proposed method and a results visualization at the OIS. Here, green pixels are $TP$, blue pixels $FN$ and red pixels $FP$. Note that in the CFD, CRKWH100 and Stone331 samples we use the pixel distance threshold of 2, meaning predictions are $TP$ pixels if they lie within 2 pixels to a ground truth, $FP$ are incorrect predictions more than 2 pixel away from a ground truth and $FN$ are missed ground truth pixels if they are further than 2 pixels away from a predicted pixel. (Best viewed digitally and in color).

*4.9. Model Efficiency and Complexity*

When utilizing an *EfficientNet B5* backbone our model has a total of 39.8 Million parameters, with the decoder consisting of 11.3 million parameters and the backbone containing 28.5 million parameters. The number of FLOPS (multiply-adds) of this model is 42.9 billion for a single channel grayscale image of size $512 \times 512$ pixel. During inference, without using test-time-augmentation, on our hardware, the model architecture can create predictions with a speed of 11.7 frames per second on images of the previously mentioned size. When utilizing test-time-augmentation, the inference speeds drops to 1.9 frames per second, due to creating and merging of multiple segmentation maps per image. Real-time applications might therefore not benefit from that approach and it seems to be better suited for batch-processing. The approximate training time of our proposed method using an *EfficientNet B5* backbone for a single run with the aforementioned configuration is 1:20h for CFD, 4:50h for CT260 and 5:30h for Deepcrack-DB.

## 5. Conclusion

In this work we propose novel surface crack segmentation methods using an encoder-decoder based deep learning architecture. We introduce a decoder design which can be added to existing feature encoder backbone architectures such as ResNet, VGG or EfficientNet to complete the U-shape of the network. We also propose the use of test-time-augmentation and performing a statistical analysis over multiple training runs. Whilst both techniques may previously have been used in other fields, this is the first application in the crack-segmentation domain. It is shown that the use of test-time-augmentation generally increases the models performance and performing this statistical analysis allows for better reproducibility, as it captures a models performance better rather than reporting only one results, which is common in this field [3, 13, 7, 27]. We show that our optimized deep encoder-decoder methods using a pretrained *EfficientNet B5* backbone as an encoder with our proposed decoder, new state of the art results in four different crack segmentation datasets are achieved. Namely we outperform the previous best approaches: FPCNet [13] on the CFD [6] dataset, DeepCrack [3] on the CRKWH100 and Stone331 datasets [3] and another DeepCrack [7] on DeepCrack-DB [7]. By presenting our results on multiple datasets we show that our approach is well suited for crack segmentation. Additionally we also highlight the generalization abilities by cross-testing inbetween all datasets and outline the limitations which align with the results of the generalization-study: Our methods performance is limited when the appearance of cracks and their backgrounds in the training- and testing data is very different and advise for better performance to use more diverse datasets.

Further experiments show that that EfficientNet based backbones of similar size to either *ResNet 50* and *VGG 19* perform better on a majority of the datasets, exhibiting a better generalization performance. We also show that adding transfer learning with pretrained ImageNet weights and using deep supervision leads to an additional performance improvement.

In future work, we aim to study the performance of using this method in other domains as well as examine the ability to denoise crack images. We also aim to overcome the limitations of this approach by generating a much larger dataset and train this method on it. Additionally, further research may also study the direct impact of TTA and the underlying reason of different performance improvements in the various datasets.

**Acknowledgments**

**References**

[1] M. Salman, S. Mathavan, K. Kamal, M. Rahman, Pavement Crack Detection using the Gabor Filter, in: Conference on Intelligent Transportation Systems, IEEE, 2013, pp. 2039–2044.

[2] M. Eisenbach, R. Stricker, D. Seichter, K. Amende, K. Debes, M. Sesselmann, D. Ebersbach, U. Stoeckert, H.-M. Gross, How to Get Pavement Distress Detection Ready for Deep Learning? A Systematic Approach, in: International Joint Conference on Neural Networks, IEEE, 2017, pp. 2039–2047.

[3] Q. Zou, Z. Zhang, Q. Li, X. Qi, Q. Wang, S. Wang, DeepCrack: Learning Hierarchical Convolutional Features for Crack Detection, IEEE Transactions on Image Processing 28 (3) (2019) 1498–1512.

[4] K. Yamaki, K. Matsushima, O. Takahashi, Road Deformation Detection Based Sensor Fusion, in: International Conference on Information Technology and Electrical Engineering, IEEE, 2017, pp. 1–6.

[5] R. Fan, J. Jiao, J. Pan, H. Huang, S. Shen, M. Liu, Real-Time Dense Stereo Embedded in a Uav for Road Inspection, in: Conference on Computer Vision and Pattern Recognition Workshops, IEEE, 2019.

[6] Y. Shi, L. Cui, Z. Qi, F. Meng, Z. Chen, Automatic Road Crack Detection using Random Structured Forests, IEEE Transactions on Intelligent Transportation Systems 17 (12) (2016) 3434–3445.

[7] Y. Liu, J. Yao, X. Lu, R. Xie, L. Li, DeepCrack: A Deep Hierarchical Feature Learning Architecture for Crack Segmentation, Neurocomputing 338 (2019) 139–153.

[8] H. Oliveira, P. L. Correia, Automatic Road Crack Detection and Characterization, IEEE Transactions on Intelligent Transportation Systems 14 (1) (2013) 155–168.

[9] G. Lin, F. Liu, A. Milan, C. Shen, I. Reid, Refinenet: Multi-Path Refinement Networks for Dense Prediction, IEEE Transactions on Pattern Analysis and Machine Intelligence 42 (5) (2019) 1228–1242.

[10] M. Tan, Q. V. Le, EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks, arXiv preprint arXiv:1905.11946 (2019).

[11] M. Tan, R. Pang, Q. V. Le, EfficientDet: Scalable and efficient object detection, in: Conference on Computer Vision and Pattern Recognition, IEEE, 2020, pp. 10781–10790.

[12] K. He, X. Zhang, S. Ren, J. Sun, Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification, in: International Conference on Computer Vision, IEEE, 2015, pp. 1026–1034.

[13] W. Liu, Y. Huang, Y. Li, Q. Chen, FPCNet: Fast Pavement Crack Detection Network Based on Encoder-Decoder Architecture, arXiv preprint arXiv1907.02248 (2019).

[14] L. Zhang, F. Yang, Y. D. Zhang, Y. J. Zhu, Road Crack Detection using Deep Convolutional Neural Network, in: International Conference on Image Processing, IEEE, 2016, pp. 3708–3712.

[15] X. Wang, Z. Hu, Grid-Based Pavement Crack Analysis Using Deep Learning, in: International Conference on Transportation Information and Safety, IEEE, 2017, pp. 917–924.

[16] T. A. Carr, M. D. Jenkins, M. I. Iglesias, T. Buggy, G. Morison, Road Crack Detection using a Single Stage Detector Based Deep Neural Network, in: Workshop on Environmental, Energy, and Structural Monitoring Systems, IEEE, 2018, pp. 1–5.

[17] V. Badrinarayanan, A. Kendall, R. Cipolla, Segnet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (12) (2017) 2481–2495.

[18] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation, in: European Conference on Computer Vision, Springer, 2018, pp. 801–818.

[19] O. Ronneberger, P. Fischer, T. Brox, U-Net: Convolutional Networks for Biomedical Image Segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2015, pp. 234–241.

[20] F. Yang, L. Zhang, S. Yu, D. Prokhorov, X. Mei, H. Ling, Feature Pyramid and Hierarchical Boosting Network for Pavement Crack Detection, IEEE Transactions on Intelligent Transportation Systems 21 (4) (2019) 1525–1535.

[21] J. König, M. D. Jenkins, P. Barrie, M. Mannion, G. Morison, A Convolutional Neural Network for Pavement Surface Crack Segmentation Using Residual Connections and Attention Gating, in: International Conference on Image Processing, IEEE, 2019, pp. 1460–1464.

[22] W. Song, G. Jia, H. Zhu, D. Jia, L. Gao, Automated Pavement Crack Damage Detection Using Deep Multiscale Convolutional Features, Journal of Advanced Transportation 2020 (2020).

[23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs, IEEE Transactions on Pattern Analysis and Machine Intelligence 40 (4) (2018) 834–848.

[24] K. Qian, Automated Detection of Steel Defects Via Machine Learning Based on Real-Time Semantic Segmentation, in: International Conference on Video and Image Processing, ACM, 2019, pp. 42–46.

[25] K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, in: International Conference on Learning Representations, 2015.

[26] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, in: Conference on Computer Vision and Pattern Recognition, IEEE, 2016, pp. 770–778.

[27] Y. Inoue, H. Nagayoshi, Deployment Conscious Automatic Surface Crack Detection, in: Winter Conference on Applications of Computer Vision, IEEE, 2019, pp. 686–694.

[28] A. Mohan, S. Poobal, Crack Detection Using Image Processing: A Critical Review and Analysis, Alexandria Engineering Journal 57 (2) (2018) 787–798.

[29] S. Chambon, J.-M. Moliard, Automatic Road Pavement Assessment with Image Processing: Review and Comparison, International Journal of Geophysics 2011 (2011).

[30] N. Tanaka, K. Uematsu, A Crack Detection Method in Road Surface Images Using Morphology, in: Workshop on Machine Vision Applications, IAPR, 1998, pp. 154–157.

[31] Y. Maode, B. Shaobo, X. Kun, H. Yuyao, Pavement Crack Detection and Analysis for High-grade Highway, in: International Conference on Electronic Measurement and Instruments, IEEE, 2007, pp. 4–548–4–552.

[32] L. Peng, W. Chao, L. Shuangmiao, F. Baocai, Research on Crack Detection Method of Airport Runway Based on Twice-Threshold Segmentation, in: International Conference on Instrumentation and Measurement, Computer, Communication and Control, IEEE, 2015, pp. 1716–1720.

[33] Y. Fujita, Y. Mitani, Y. Hamamoto, A Method for Crack Detection on a Concrete Structure, in: International Conference on Pattern Recognition, Vol. 3, IEEE, 2006, pp. 901–904.

[34] H. Oliveira, P. L. Correia, Automatic Road Crack Segmentation Using Entropy and Image Dynamic Thresholding, in: European Signal Processing Conference, IEEE, 2009, pp. 622–626.

[35] P. Subirats, J. Dumoulin, V. Legeay, D. Barba, Automation of Pavement Surface Crack Detection using the Continuous Wavelet Transform, in: International Conference on Image Processing, IEEE, 2006, pp. 3037–3040.

[36] K. C. P. Wang, Q. Li, W. Gong, Wavelet-Based Pavement Distress Image Edge Detection with a Trous Algorithm, Transportation Research Record 2024 (1) (2007) 73–81.

[37] J. Zhou, P. Huang, F.-P. Chiang, Wavelet-Based Pavement Distress Classification, Transportation Research Record 1940 (1) (2005) 89–98.

[38] H. Zhao, G. Qin, X. Wang, Improvement of Canny Algorithm Based on Pavement Edge Detection, in: International Congress on Image and Signal Processing, Vol. 2, IEEE, 2010, pp. 964–967.

[39] J. Canny, A Computational Approach to Edge Detection, IEEE Transactions on Pattern Analysis and Machine Intelligence PAMI-8 (6) (1986) 679–698.

[40] K. Makantasis, E. Protopapadakis, A. Doulamis, N. Doulamis, C. Loupos, Deep Convolutional Neural Networks for Efficient Vision Based Tunnel Inspection, in: International Conference on Intelligent Computer Communication and Processing, IEEE, 2015, pp. 335–342.

[41] K. Fernandes, L. Ciobanu, Pavement Pathologies Classification using Graph-Based Features, in: International Conference on Image Processing, IEEE, 2014, pp. 793–797.

[42] M. D. Jenkins, T. A. Carr, M. I. Iglesias, T. Buggy, G. Morison, A Deep Convolutional Neural Network for Semantic Pixel-Wise Segmentation of Road and Pavement Surface Cracks, in: European Signal Processing Conference, IEEE, 2018, pp. 2120–2124.

[43] Z. Fan, Y. Wu, J. Lu, W. Li, Automatic Pavement Crack Detection Based on Structured Prediction with the Convolutional Neural Network, arXiv preprint arXiv:1802.02208 (2018).

[44] Z. Fan, C. Li, Y. Chen, P. D. Mascio, X. Chen, G. Zhu, G. Loprencipe, Ensemble of Deep Convolutional Neural Networks for Automatic Pavement Crack Detection and Measurement, Coatings 10 (2) (2020) 152.

[45] S. Xie, Z. Tu, Holistically-Nested Edge Detection, in: Conference on Computer Vision and Pattern Recognition, IEEE, 2015, pp. 1395–1403.

[46] J. Hu, L. Shen, G. Sun, Squeeze-and-Excitation Networks, in: IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7132–7141.

[47] T.-Y. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, S. J. Belongie, Feature Pyramid Networks for Object Detection, IEEE Internation Conference on Computer Vision and Pattern Recognition (2017) 936–944.

[48] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, Z. Tu, Deeply-Supervised Nets, in: Artificial Intelligence and Statistics, 2015, pp. 562–570.

[49] K. He, J. Sun, X. Tang, Guided Image Filtering, in: European Conference on Computer Vision, Springer, 2010, pp. 1–14.

[50] Z. Zhang, Q. Liu, Y. Wang, Road Extraction by Deep Residual U-Net, IEEE Geoscience and Remote Sensing Letters 15 (5) (2018) 749–753.

[51] M. Z. Alom, M. Hasan, C. Yakopcic, T. M. Taha, V. K. Asari, Recurrent Residual Convolutional Neural Network based on U-Net (R2U-Net) for Medical Image Segmentation, arXiv preprint arXiv:1802.06955 (2018).

[52] M. Kampffmeyer, N. Dong, X. Liang, Y. Zhang, E. P. Xing, ConnNet: A Long-Range Relation-Aware Pixel-Connectivity Network for Salient Segmentation, IEEE Transactions on Image Processing 28 (5) (2018) 2518–2529.

[53] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, Q. V. Le, Mnasnet: Platform-Aware Neural Architecture Search for Mobile, in: Conference on Computer Vision and Pattern Recognition, IEEE, 2019, pp. 2820–2828.

[54] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, Others, ImageNet Large Scale Visual Recognition Challenge, International Journal of Computer Vision 115 (3) (2015) 211–252.

[55] A. Krizhevsky, Hinton, Geoffrey, Learning Multiple Layers of Features from Tiny Images, Tech. rep., University of Toronto (2009).

[56] F. Milletari, N. Navab, S.-A. Ahmadi, V-net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation, in: International Conference on 3D Vision, IEEE, 2016, pp. 565–571.

[57] P. Arbelaez, M. Maire, C. Fowlkes, J. Malik, Contour Detection and Hierarchical Image Segmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence 33 (5) (2011) 898–916.

[58] F. Chollet, Keras, https://keras.io/ (2015).

[59] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, Others, TensorFlow: A System for Large-Scale Machine Learning, in: Operating Systems Design and Implementation, Vol. 16, 2016, pp. 265–283.

[60] P. Yakubovskiy, EfficientNet, https://github.com/qubvel/efficientnet (2019).

[61] G. Li, Y. Xie, L. Lin, Y. Yu, Instance-Level Salient Object Segmentation, in: Conference on Computer Vision and Pattern Recognition, IEEE, 2017, pp. 2386–2395.

[62] S. Marrone, S. Olivieri, G. Piantadosi, C. Sansone, Reproducibility of Deep CNN for Biomedical Image Processing Across Frameworks and Architectures, in: European Signal Processing Conference, IEEE, 2019, pp. 1–5.

[63] O. Oktay, J. Schlemper, L. L. Folgoc, M. Lee, M. Heinrich, K. Misawa, K. Mori, S. McDonagh, N. Y. Hammerla, B. Kainz, Others, Attention U-Net: Learning Where to Look for the Pancreas, arXiv preprint arXiv:1804.03999 (2018).

[64] J. Long, E. Shelhamer, T. Darrell, Fully Convolutional Networks for Semantic Segmentation, in: Conference on Computer Vision and Pattern Recognition, IEEE, 2015, pp. 3431–3440.

[65] J. König, M. D. Jenkins, P. Barrie, M. Mannion, G. Morison, Segmentation of Surface Cracks Based on a Fully Convolutional Neural Network and Gated Scale Pooling, in: European Signal Processing Conference, IEEE, 2019.