

Self-Supervised Multi-Scale Transformer with Attention-Guided Fusion for Efficient Crack Detection

Abstract

This study presents *Crack-Segmenter*, a fully self-supervised segmentation framework for crack detection. Unlike existing supervised methods requiring expensive pixel-level annotations, our approach eliminates this dependency completely through three integrated modules: Scale-Adaptive Embedder (SAE), Directional Attention Transformer (DAT), and Attention-Guided Fusion (AGF). These modules work together seamlessly, with SAE capturing multi-scale crack features, DAT enhancing linear crack continuity through directional attention, and AGF adaptively integrating these representations into unified segmentation outputs. Extensive experiments on ten public crack datasets demonstrated that *Crack-Segmenter* consistently outperformed 13 state-of-the-art supervised methods across all evaluation metrics, including mean Intersection over Union (mIoU), Dice score, XOR, and Hamming Distance (HM). These results demonstrate that annotation-free crack segmentation can achieve superior performance while enabling scalable infrastructure monitoring and automated maintenance decision-making, advancing the state of self-supervised learning in infrastructure applications.

Keywords: Self-supervised segmentation, Crack detection, Multi-scale transformer, Directional attention, Infrastructure monitoring, Annotation-free learning, Road maintenance automation

1. Introduction

Transportation infrastructure, particularly road networks, is critical for public safety and economic development. Roads connect communities, facilitate commerce, and ensure efficient movement of people and goods. However, the constant exposure of roads to traffic loads and weather conditions gradually weakens their structural integrity, often resulting in surface cracks. Even tiny cracks that develop on these roads can quickly grow into severe defects

1 such as potholes or large pavement failures if they are not detected and re-
2 paired early. This makes preventive maintenance very vital in prolonging
3 the lifespan of pavements. For instance, research indicates that preventive
4 maintenance on small pavement cracks can reduce future repair costs by
5 approximately 50–70%, highlighting the financial benefit of timely interven-
6 tion [1, 2]. Accurate, pixel-level crack detection maps are thus essential for
7 enabling early and low-cost pavement maintenance. These detailed maps
8 can guide agencies in prioritizing repairs, allocating budgets efficiently, and
9 reducing overall maintenance expenses. To generate such maps automati-
10 cally, researchers have developed supervised learning approaches that have
11 achieved strong performance in diverse pavement crack segmentation tasks.

12 Several fully-supervised segmentation models for pavement crack segmen-
13 tation have been proposed. For instance, Lau et al. [3] replace the U-Net
14 encoder with a pretrained ResNet-34 and report F1 scores of 96 % on CFD
15 and 73 % on Crack500, proving that full-mask supervision can yield strong
16 accuracy. [4] proposed Context-CrackNet, a supervised model using global
17 and local attention modules to accurately segment tiny and large pavement
18 cracks. [5] also developed PDSNet, a supervised deep learning framework for
19 segmenting multiple asphalt pavement distresses using 2D and 3D images,
20 achieving a high mean Intersection over Union (MIoU) of 83.7%. All these ap-
21 proaches required extensive annotations during the model training process.
22 Especially, [5] had to manually annotate 5000 images of pavement cracks
23 by drawing pixel-level labels which is labor-intensive. Given the practical
24 challenges and costs of extensive pixel-level annotations in fully supervised
25 approaches, there is growing interest in alternative learning methods that
26 reduce annotation demands while preserving segmentation performance.

27 Weakly supervised and semi supervised segmentation methods have emerged
28 as promising alternatives, significantly reducing annotation efforts while aim-
29 ing to retain segmentation accuracy. For example, Xiang et al. [6] intro-
30 duced UWSCS, a crack segmentation framework that uses limited coarse
31 labels alongside superpixel and shrink-based correction modules to train a
32 dual encoder network. Similarly, recent studies have integrated the Segment
33 Anything Model (SAM) with interactive segmentation using bounding box
34 prompts and deep transfer learning to enable semi supervised crack detec-
35 tion [7, 8]. However, these approaches still rely on manual bounding box
36 annotations, which remain costly and time consuming, particularly for large
37 scale pavement crack datasets. To further reduce annotation dependency,
38 some recent frameworks have incorporated advanced learning strategies such

1 as adversarial training [9, 10, 11], student teacher learning [12, 13, 14], and
 2 graph based modeling using graph convolutional networks [15, 16]. These
 3 techniques are often used within weakly or semi supervised architectures to
 4 enhance learning from limited annotations, but they still depend on some
 5 form of manual supervision during training. Although several recent studies
 6 describe their approaches as fully self supervised, many of them still rely on
 7 ground truth annotations at some point in the pipeline, such as during pre-
 8 training, pseudo label generation, or model calibration [17, 18]. As a result,
 9 they do not fully meet the criteria of annotation-free segmentation. These
 10 continued limitations highlight the need for a robust, efficient, and scalable
 11 self supervised framework that can accurately segment pavement cracks with-
 12 out relying on any annotated data during training and still achieve superior
 13 performance compared to current supervised and semi-supervised methods.

14 To address the limitations discussed above, we propose an efficient, fully
 15 self-supervised segmentation framework designed specifically for pavement
 16 crack segmentation. Unlike previous approaches, our method requires no
 17 manual annotations or ground truth pixel labels, thus significantly reduc-
 18 ing the cost and time associated with labeling. Our proposed architecture
 19 consists of three main modules: the Scale-Adaptive Embedder, the Directed
 20 Multi-Branch Transformer, and the Attention-Guided Fusion module. The
 21 Scale-Adaptive Embedder transforms input images into token representations
 22 at three distinct spatial scales, capturing features ranging from tiny hairline
 23 cracks to wider pavement defects. Building upon these multi-scale embed-
 24 dings, our Directed Multi-Branch Transformer applies directional efficient
 25 attention mechanisms, effectively preserving linear crack structures and cap-
 26 turing essential spatial relationships without extensive computational costs.
 27 To seamlessly merge these multi-scale representations, our Attention-Guided
 28 Fusion module adaptively weights features from different scales, ensuring that
 29 both fine-grained details and broader contextual information are combined
 30 effectively. The introduction of cross-scale consistency losses [19] further en-
 31 hances the model’s ability to learn robust representations without manual
 32 labels. The main contributions of our work have been summarized below:

- 33 • We introduce *Crack-Segmenter*, an end-to-end self-supervised frame-
 34 work for pavement-crack segmentation that requires no pixel-level or
 35 weak annotations, thereby reducing annotation costs significantly.
- 36 • We design three modules: the *Scale-Adaptive Embedder* (SAE) for

multi-resolution feature extraction, the *Directional Attention Transformer* (DAT) to preserve elongated crack geometry, and the *Attention-Guided Fusion* (AGF) to adaptively merge scale-specific representations.

- We develop inter-scale and intra-scale consistency losses to enhance coherent feature representations, substantially improving model learning without manual supervision.
- We evaluated *Crack-Segmenter* on ten public crack datasets against 13 state-of-the-art fully supervised models. It outperformed every baseline with statistically significant gains.

2. Related Works

2.1. Fully-Supervised Crack Segmentation Methods

Fully supervised methods have utilized rich pixel-level annotations to develop specialized models for accurate pavement crack segmentation. For example, Dung et al. [20] adopted a Fully Convolutional Network (FCN) based on a VGG16 encoder for automated concrete crack detection, where high-level features were upsampled via deconvolution to produce crack masks. While this approach proved the feasibility of end-to-end crack segmentation, the basic FCN architecture struggled with limited context and spatial consistency, often missing fine crack details. Researchers subsequently improved crack segmentation by integrating multi-scale feature fusion and boosting techniques. Yang et al. [21] introduced a Feature Pyramid and Hierarchical Boosting Network (FPHBN) that fuses semantic information from deep and shallow layers and utilizes a feature pyramid to enhance detection of cracks at different scales. This method improved the generalization to various crack widths and backgrounds by considering multi-level features simultaneously. Similarly, Liu et al. [22] proposed a deeply supervised encoder-decoder network (DeepCrack) which aggregates multi-scale features from multiple network layers, capturing both fine and coarse crack patterns for more robust segmentation. These multi-scale approaches demonstrated higher accuracy than plain FCNs, but the heavy pooling in their backbones could still cause some loss of fine spatial information. To better preserve crack locality, later studies adopted encoder-decoder architectures like U-Net. Jenkins et al. [23] showed that a vanilla U-Net can effectively segment road cracks at the

1 pixel level, and further improvements added attention gating to suppress
 2 irrelevant background features. Building on this idea, Pan et al. [24] devel-
 3 oped an attention-enhanced U-Net variant called SCHNet, which incorpo-
 4 rates parallel spatial, channel, and feature-pyramid attention modules into
 5 a VGG19-based network. Transformer-based architectures have also been
 6 explored: Guo et al. [25] proposed a Crack Transformer (CT) model using
 7 a Swin Transformer encoder and all-MLP decoder, showing robust perfor-
 8 mance in detecting long, complex cracks even under noisy conditions. Simi-
 9 larly, [26] embedded a Transformer encoder within a U-shaped CNN in their
 10 CrackFormer network, substantially improving segmentation continuity and
 11 accuracy for thin cracks. However, these supervised approaches are heavily
 12 dependent on pixel-level annotations, which are time-consuming, costly, and
 13 practically difficult to obtain at a large scale [27].

14 *2.2. Weakly-Supervised and Semi-Supervised Crack Segmentation*

15 To alleviate the burden of dense labeling, researchers have developed
 16 weakly supervised frameworks that learn from coarse annotations (e.g. image-
 17 level labels). Al-Huda et al. [28] proposed a two-stage weakly supervised
 18 method based on class activation mapping and iterative refinement: a crack
 19 classification network first produces initial pixel indications of cracks at mul-
 20 tiple scales, then a U-Net with an attention mechanism is trained on these
 21 noisy masks and incrementally fine-tunes the predictions. He et al. [29]
 22 similarly employed image-level labels to drive crack segmentation by using
 23 a generative adversarial localization strategy. They trained a U-GAT-IT
 24 model to generate class activation maps of cracks, iteratively erased detected
 25 regions to discover new crack areas, and then converted these refined maps
 26 into pseudo-labels to train a segmentation network. Semi-supervised meth-
 27 ods have emerged to exploit unlabeled roadway images alongside limited
 28 labeled data, further reducing the need for annotations. Shim et al. [30] pio-
 29 neered an adversarial learning-based semi-supervised segmentation approach
 30 for concrete crack detection. Their model employed multiscale feature extrac-
 31 tors and a generative adversarial network to train on labeled and unlabeled
 32 data simultaneously. Building on consistency-driven learning, Shi et al. [31]
 33 developed a crack segmentation model that enforces mutual consistency con-
 34 straints between dual network predictions and incorporates a boundary-aware
 35 loss. Another innovative strategy is the two-stage CrackDiffusion framework
 36 proposed by Han et al. [32], which combines an unsupervised anomaly detec-
 37 tion stage with a supervised refinement stage. In the first stage, a diffusion-

1 based inpainting model removes cracks from images to generate crack-free
 2 counterparts and uses the differences (via structural similarity measures) to
 3 localize cracks without labels. In the second stage, those initial crack maps
 4 inform a U-Net segmentation model that learns to produce precise crack
 5 masks. Despite these advancements, weakly and semi-supervised segmenta-
 6 tion methods still depend on pseudo-labels or sparse annotations. This re-
 7 quirement makes it impractical to annotate large-scale, high-resolution pave-
 8 ment crack segmentation datasets which has completely no pixel-level labels.
 9 Because of this limitation, researchers have turned their attention toward
 10 fully self-supervised segmentation methods. These fully self-supervised ap-
 11 proaches aim to eliminate the need for ground-truth pixel-level annotations
 12 entirely during the training process. However, research in the area of fully
 13 self-supervised pavement distress segmentation remains limited and underex-
 14 plored. Moreover, the few existing pavement crack segmentation techniques
 15 that claim to be fully self-supervised still suffer from significant limitations
 16 and drawbacks.

17 *2.3. Limitations*

18 Self-supervised learning (SSL) is a representation learning paradigm in
 19 which the supervisory signal is derived automatically from the structure of
 20 the unlabelled data itself, rather than from externally provided human anno-
 21 tations [33]. In the context of semantic segmentation, self-supervised segmen-
 22 tation extends this idea from learning image-level representations to learning
 23 pixel-level assignments without any human-drawn masks. The model first
 24 designs an intrinsic task such as grouping pixels with similar colour statis-
 25 tics, enforcing consistency between differently augmented views, and then
 26 treats the resulting pseudo-labels as ground truth for training [34, 35]. Be-
 27 cause these labels arise automatically from the data, the pipeline needs no
 28 external annotations at any stage. SSL has been applied to pavement and
 29 other civil-infrastructure images through tasks such as self-training for cracks
 30 [36, 37], pavement-surface anomaly detection [38, 39], sidewalk-quality clas-
 31 sification [40]. Some studies have attempted to apply SSL in segmentation
 32 of pavement cracks. For example, [17] propose SS-YOLO, a YOLOv8-based
 33 crack segmentation model that fuses CBAM and Gaussian multi-head self-
 34 attention with curriculum learning-driven pseudo-labeling. However, SS-
 35 YOLO still starts from a fully supervised YOLOv8 backbone trained on real
 36 crack masks before it ever “self-labels” unannotated data, so it isn’t end-to-
 37 end self-supervised. Similarly, Zhang et al. [18] proposed a dual cycle-GAN

1 that learns to translate crack image patches into GT-like structure patterns
 2 (and back) using an unpaired “structure library” of binary skeletons. Con-
 3 versely, because that library is built from pixel-precise, human-annotated
 4 curves (e.g., VOC object-boundary masks, Berkeley contour annotations,
 5 and public crack GTs), the method still depends on existing pixel-level labels
 6 and isn’t truly end-to-end self-supervised. Song et al. [41] also proposed a
 7 two-stage pavement-crack framework: an improved U-Net (U-Net augmented
 8 with residual blocks and attention gates) is first contrastively pre-trained on
 9 unlabeled crack (background) patches, then fully fine-tuned with pixel-level
 10 ground-truth masks during the second stage. Because the network does not
 11 directly learn the image-to-mask mapping until this second stage where ev-
 12 ery gradient is computed against human-annotated labels, the approach is
 13 merely semi-supervised, not an end-to-end fully self-supervised segmentation
 14 method. Ma et al. [42] introduced UP-CrackNet, pavement-crack detector
 15 that trains a conditional GAN to inpaint randomly masked regions of crack-
 16 free road images. At inference, cracks are segmented by thresholding the
 17 pixel-wise reconstruction residuals, allowing annotation-free training. One
 18 problem with UP-CrackNet is that since it never sees real crack patterns
 19 during optimization, it must treat every unfamiliar texture as a defect. This
 20 makes the network yield a high error map and treats labels artifacts present
 21 in the crack image (e.g. leaf, tyre mark) as a crack, inflating false positives
 22 that a crack-aware model could reject.

23 3. Methodology

24 3.1. Problem Structure and Overview

25 Pavement binary crack segmentation involves accurately classifying ev-
 26 ery pixel in an image as either crack or non-crack. Currently, this task relies
 27 heavily on extensive ground truth pixel-level annotations (masks). That is,
 28 given a batch of input pavement images $I \in \mathbb{R}^{B \times H \times W \times C}$ where B , H , W , and
 29 C denote the batch size, height, width, and number of channels respectively,
 30 the goal of this task is to predict their corresponding binary segmentation
 31 maps $S \in \{0, 1\}^{B \times H \times W \times 1}$, where each pixel is assigned a value of 1 if it
 32 belongs to a crack and 0 otherwise. This supervised learning approach has
 33 achieved significant success in binary crack segmentation. However, obtain-
 34 ing the precise ground truth annotations for supervised learning is both costly
 35 and impractical for large-scale crack datasets.

To address this limitation, we propose a fully self-supervised framework that completely eliminates the need for ground truth labels or masks by learning robust feature representations directly from the input images. Our framework utilizes multi-scale feature extraction, directional attention mechanisms, and adaptive scale fusion to identify and segment pavement cracks accurately. The challenge here is formulating a learning paradigm that can effectively differentiate crack pixels from non-crack pixels without explicit supervision, thus making it scalable and cost-efficient.

3.2. Overall Framework

The proposed self-supervised segmentation framework comprises of three primary modules designed to collectively address the aforementioned challenge: the **Scale Adaptive Embedder** (Φ_{SAE}), **Directional Attention Transformer** (f_{DAT}), and **Attention-Guided Fusion** (f_{AGF}). The overall framework for our proposed architecture has been shown in Figure 1.

The proposed framework begins with the Scale Adaptive Embedder module, which processes a batch of input images I simultaneously at multiple scales: fine, small, and large. Specifically, given a batch of input image I , the module produces three distinct feature embeddings with embedding dimension D :

$$\Phi_{\text{SAE}}(I) = \{F_f, F_s, F_l\}, \quad (1)$$

where $F_f, F_s \in \mathbb{R}^{B \times D \times H \times W}$ and $F_l \in \mathbb{R}^{B \times D \times \frac{H}{2} \times \frac{W}{2}}$. F_f, F_s, F_l correspond to the fine, small, and large scale feature embeddings respectively. This multi-scale representation ensures comprehensive capture of cracks of varying widths and complexities, enhancing sensitivity to fine details and broad spatial context simultaneously.

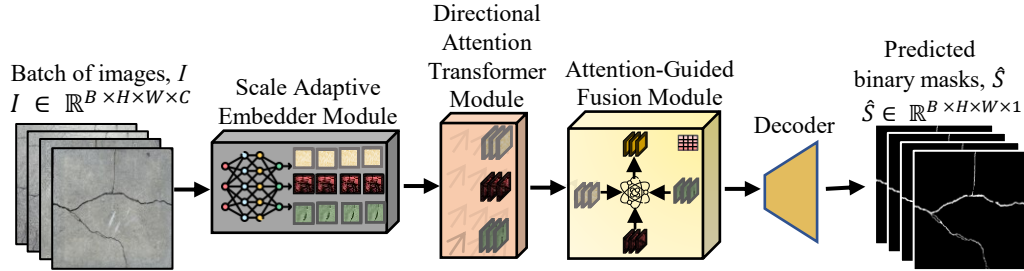


Figure 1: Overall framework for our proposed architecture.

Next, each scale-specific feature embedding from the Scale Adaptive Embedder undergoes further refinement through the Directional Attention Transformer module. This transformer applies efficient attention with directed convolutions to maintain and emphasize linear crack structures within feature maps, essential for accurate pavement crack identification. Formally, this transformation is represented as:

$$F'_f = f_{\text{DAT}}(F_f) \quad (2)$$

$$F'_s = f_{\text{DAT}}(F_s) \quad (3)$$

$$F'_l = f_{\text{DAT}}(F_l) \quad (4)$$

where $F'_f, F'_s \in \mathbb{R}^{B \times D \times H \times W}$ and $F'_l \in \mathbb{R}^{B \times D \times \frac{H}{2} \times \frac{W}{2}}$. F'_f, F'_s, F'_l represent the refined feature embeddings for the different scales after passing through the Directional Attention Transformer module. This module enhances local contextual consistency within each scale, preserving crucial spatial relationships pertinent to cracks.

Finally, the Attention-Guided Fusion module integrates these refined scale-specific features from the Directional Attention Transformer module using attention-based adaptive weighting to form a unified, robust representation. Specifically, the large-scale feature map F'_l is upsampled and projected to match the spatial dimensions of the other scales. Subsequently, the fusion module computes attention weights and merges the scales:

$$F_{\text{fused}} = f_{\text{AGF}}(F'_f, F'_s, \text{upsample}(F'_l)), \quad (5)$$

where $F_{\text{fused}} \in \mathbb{R}^{D \times H \times W}$ and *upsample* represents the upsampling operation. This adaptive fusion ensures optimal integration of both detailed crack structures and broader contextual information.

Finally, the fused feature representation F_{fused} is processed through a linear decoding layer to produce the final segmentation map prediction \hat{S} :

$$\hat{S} = f_{\text{decode}}(F_{\text{fused}}), \quad (6)$$

where f_{decode} is the decoding layer and $\hat{S} \in \mathbb{R}^{B \times H \times W \times 1}$ represents the predicted segmentation maps.

Through cross-scale consistency losses, specifically inter-scale and intra-scale self-supervised losses, our model learns to produce consistent and accurate segmentation predictions without any ground truth annotations (masks).

1 Thus, the proposed framework efficiently addresses pavement crack segmen-
 2 tation challenges, significantly reducing annotation costs while maintaining
 3 high segmentation performance. The next section goes into details on the
 4 modules used.

5 3.3. Scale Adaptive Embedder

6 The **Scale Adaptive Embedder (SAE)** module is designed to effec-
 7 tively capture multi-scale spatial information from pavement images, crucial
 8 for accurately identifying cracks of varying sizes and complexities. This mod-
 9 ule embeds the input image into feature representations at fine, small, and
 10 large spatial scales, enhancing the model’s ability to detect both detailed
 11 and broad crack structures simultaneously. Figure 2 shows the architecture
 12 of the Scale Adaptive Embedder.

13 Given a batch of input pavement images $I \in \mathbb{R}^{B \times C \times H \times W}$, the SAE module
 14 first applies convolutional projection operations at three distinct scales to
 15 produce feature maps tailored for fine, small, and large-scale analyses.

16 For the fine-scale embedding, the convolution operation is mathematically
 17 defined as:

$$18 \quad F_f = \sigma(W_f * I + b_f), \quad F_f \in \mathbb{R}^{B \times D \times H \times W}, \quad (7)$$

19 where $W_f \in \mathbb{R}^{D \times C \times 1 \times 1}$ represents the convolutional kernel weights, $b_f \in$
 20 \mathbb{R}^D are biases, $*$ denotes the convolution operation, and σ is a nonlinear
 21 activation function. This fine-scale operation ensures the capture of detailed
 22 and fine-grained crack patterns.

23 For the small-scale embedding, we similarly define the convolution oper-
 24 ation as:

$$25 \quad F_s = \sigma(W_s * I + b_s), \quad F_s \in \mathbb{R}^{B \times D \times H \times W}, \quad (8)$$

26 where $W_s \in \mathbb{R}^{D \times C \times 3 \times 3}$ with appropriate padding and stride set to 1 to main-
 27 tain spatial dimensions. This operation captures crack structures at inter-
 28 mediate spatial resolutions, preserving local contextual relationships.

29 The large-scale embedding convolution operation is expressed as:

$$30 \quad F_l = \sigma(W_l * I + b_l), \quad F_l \in \mathbb{R}^{B \times D \times \frac{H}{2} \times \frac{W}{2}}, \quad (9)$$

31 where $W_l \in \mathbb{R}^{D \times C \times 3 \times 3}$ with stride set to 2 and appropriate padding to reduce
 32 spatial dimensions. This large-scale embedding helps in capturing broader
 33 spatial contexts and large-scale pavement defects.

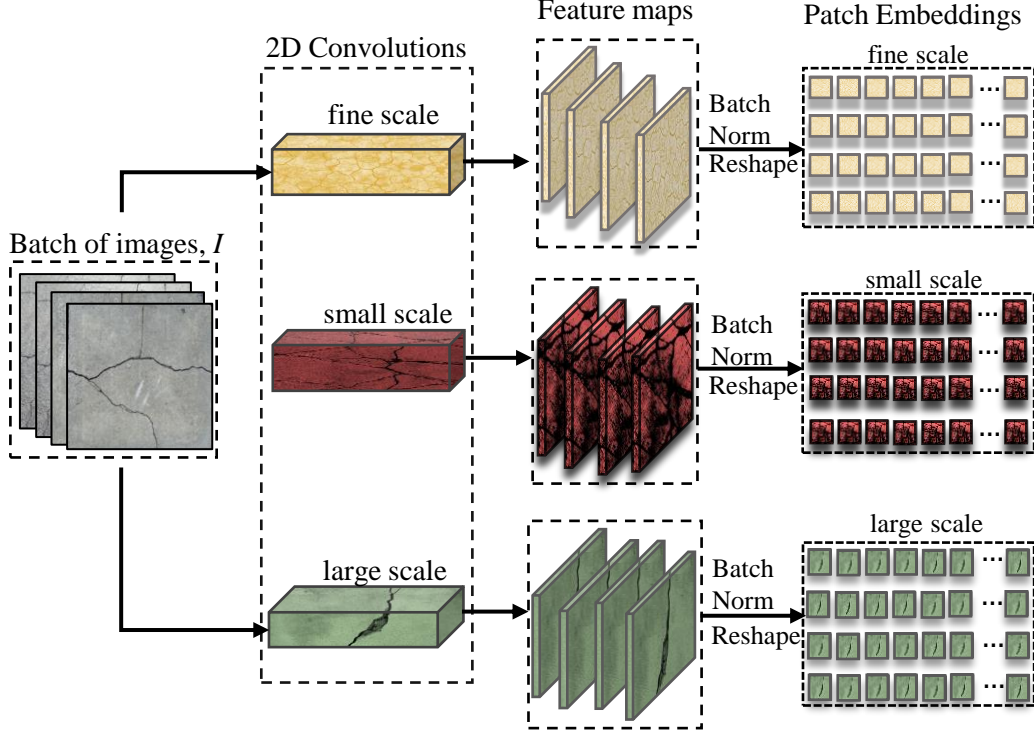


Figure 2: Scale Adaptive Embedder Module.

After convolutional projections, each feature map undergoes batch normalization to stabilize learning and enhance convergence:

$$F'_f = \text{BatchNorm}(F_f) \quad (10)$$

$$F'_s = \text{BatchNorm}(F_s) \quad (11)$$

$$F'_l = \text{BatchNorm}(F_l) \quad (12)$$

where the Batch Normalization operation $\text{BatchNorm}(\cdot)$ standardizes feature maps across each batch, thus improving training stability.

Finally, these normalized feature maps are reshaped and transposed to form sequences of patch embeddings compatible as input to the Directional Attention Transformer. Mathematically, the reshaping and transposition operation is defined as:

$$F''_{\text{scale}} \in \mathbb{R}^{B \times (H_{\text{scale}} W_{\text{scale}}) \times D}, \quad (13)$$

1 where scale denotes the specific scale (fine, small, or large), and $H_{\text{scale}}, W_{\text{scale}}$
2 represent the corresponding spatial dimensions of each scale-specific embed-
3 ding.

4 This multi-scale adaptive embedding strategy addresses the segmentation
5 challenge of capturing both narrow, hairline cracks and wider crack forma-
6 tions within pavement images. By employing distinct yet complementary
7 scale-specific convolutions, the SAE module ensures robust feature extrac-
8 tion at multiple resolutions, effectively supporting downstream modules for
9 precise segmentation without manual annotations. In the next section, we
10 will talk about the Directional Attention Module.

11 3.4. Directional Attention Transformer

12 The **Directional Attention Transformer (DAT)** module is designed
13 to explicitly model directional spatial relationships, significantly enhancing
14 the detection and segmentation of elongated, linear crack structures in pave-
15 ment images. This module integrates multi-scale embeddings from the Scale
16 Adaptive Embedder (SAE). It refines feature representations using spatially-
17 directed attention mechanisms that effectively distinguish critical crack fea-
18 tures from background noise. The architecture for the DAT module has been
19 shown in Figure 3

20 The DAT module takes as input the multi-scale embeddings from the
21 SAE module, denoted as

$$22 \quad F''_{\text{scale}} \in \mathbb{R}^{B \times (H_{\text{scale}} W_{\text{scale}}) \times D},$$

23 where $\text{scale} \in \{f, s, l\}$ corresponds to fine, small, and large scales respectively,
24 B is the batch size, and D is the embedding dimension. These embeddings are
25 first normalized via Layer Normalization to stabilize and improve training:

$$26 \quad \hat{F}_{\text{scale}} = \text{LayerNorm}(F''_{\text{scale}}), \quad \hat{F}_{\text{scale}} \in \mathbb{R}^{B \times (H_{\text{scale}} W_{\text{scale}}) \times D}. \quad (14)$$

27 Subsequently, the normalized embeddings are reshaped back to their spa-
28 tial dimensions:

$$29 \quad \hat{F}_{\text{scale}}^{\text{spatial}} \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}. \quad (15)$$

30 Directional convolutions are then applied to capture elongated structural
31 patterns essential for crack segmentation. We define these convolutions as a
32 function g , parameterized by kernels W_k and biases b_k . Specifically, for each
33 direction k (e.g., horizontal (1, 3), vertical (3, 1)), we compute:

$$Q_k = g(\hat{F}_{\text{scale}}^{\text{spatial}}; W_{Q_k}, b_{Q_k}), \quad Q_k \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}, \quad (16)$$

$$K_k = g(\hat{F}_{\text{scale}}^{\text{spatial}}; W_{K_k}, b_{K_k}), \quad K_k \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}, \quad (17)$$

where

$$g(X; W, b) = W * X + b, \quad (18)$$

with $*$ denoting convolution. Here $W_{Q_k}, W_{K_k} \in \mathbb{R}^{D \times D \times k_h \times k_w}$ are directional kernels and b_{Q_k}, b_{K_k} are biases.

Values are obtained via point-wise convolution:

$$V = W_V * \hat{F}_{\text{scale}}^{\text{spatial}} + b_V, \quad V \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}. \quad (19)$$

Attention maps are computed by softmax-normalizing the element-wise similarity of queries and keys:

$$A_k = \text{softmax}\left(\frac{Q_k \odot K_k}{\sqrt{D}}\right), \quad A_k \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}, \quad (20)$$

where \odot is element-wise multiplication. Directional context features follow:

$$C_k = A_k \odot V, \quad C_k \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}. \quad (21)$$

All C_k are concatenated and reprojected:

$$F_{\text{scale}}^{\text{attn}} = W_O(C_1 \oplus C_2 \oplus \dots \oplus C_K) + b_O, \quad F_{\text{scale}}^{\text{attn}} \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}, \quad (22)$$

where \oplus denotes channel-wise concatenation.

Finally, these features undergo Layer Normalization and a depth-wise convolutional feed-forward network:

$$\tilde{F}_{\text{scale}} = \text{LayerNorm}(F_{\text{scale}}^{\text{attn}}), \quad (23)$$

$$F_{\text{scale}}^{\text{FFN}} = \text{FFN}_{dw}(\tilde{F}_{\text{scale}}) + F_{\text{scale}}^{\text{attn}}, \quad F_{\text{scale}}^{\text{FFN}} \in \mathbb{R}^{B \times D \times H_{\text{scale}} \times W_{\text{scale}}}. \quad (24)$$

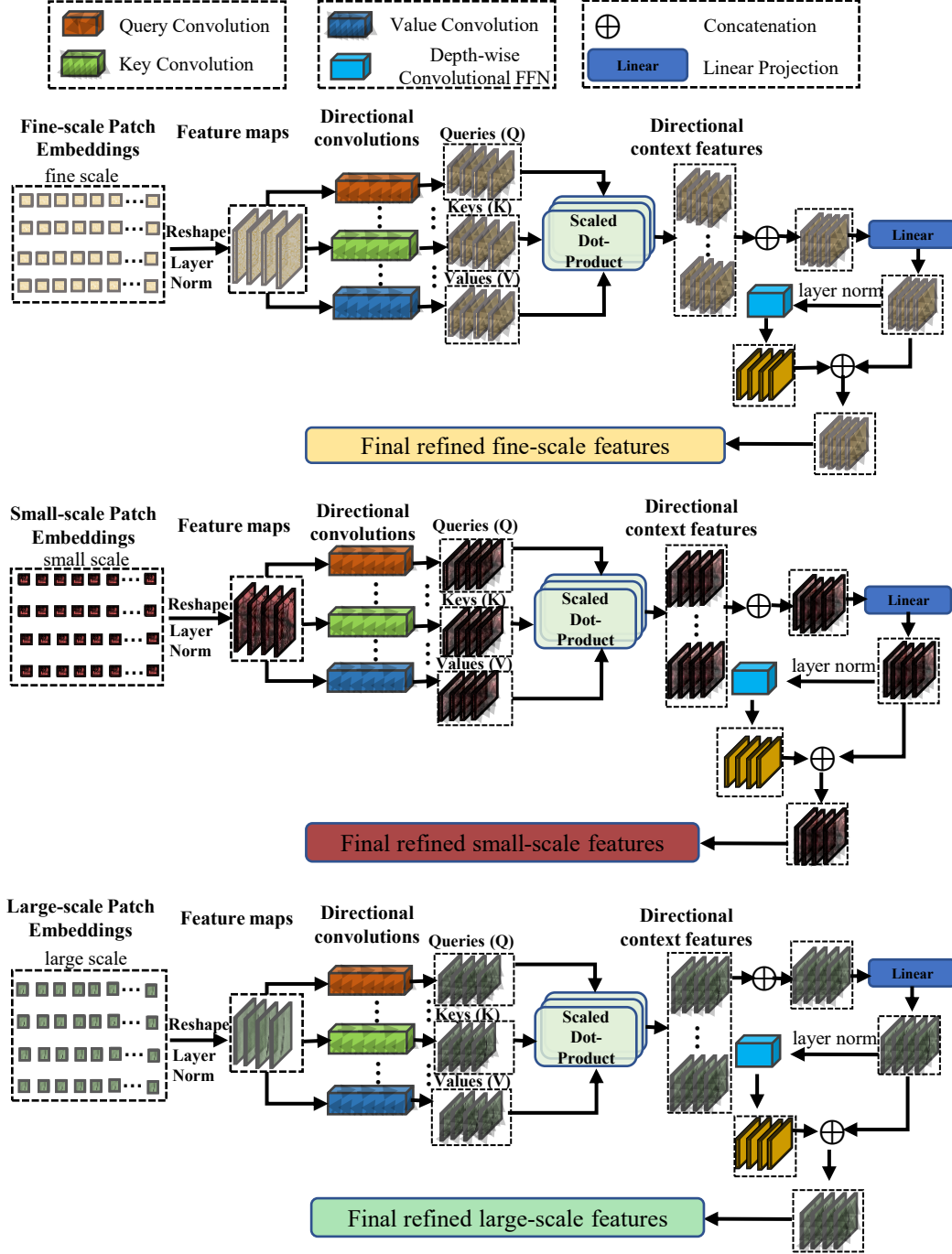


Figure 3: Directional Attention Transformer Module.

The output $F_{\text{scale}}^{\text{FFN}}$ thus encodes directional spatial relationships, enhancing multi-scale crack segmentation. These refined features feed into the Attention-Guided Fusion module.

3.5. Attention-Guided Fusion Module

The **Attention-Guided Fusion (AGF)** module intelligently integrates multi-scale feature representations into a unified, robust feature map. Following the Scale Adaptive Embedder and Directional Attention Transformer modules, this module uses adaptive attention mechanisms to determine the optimal combination of features from the fine, small, and large scales. The AGF module applies dynamic weights to each scale-specific feature according to its contextual relevance. This ensures that the final representation captures essential pavement crack details across multiple resolutions. Figure 4 shows the full architecture for the AGF module.

Let $F_f^{\text{FFN}} \in \mathbb{R}^{B \times D \times H \times W}$, $F_s^{\text{FFN}} \in \mathbb{R}^{B \times D \times H \times W}$, and $F_l^{\text{FFN}} \in \mathbb{R}^{B \times D \times \frac{H}{2} \times \frac{W}{2}}$ denote refined features from the fine, small, and large scales, respectively, obtained from the Directional Attention Transformer. To align spatial dimensions across scales, the large-scale feature map F_l^{FFN} is first upsampled and projected:

$$F_l^{\text{proj}} = \text{Conv}_{1 \times 1}(\text{Upsample}(F_l^{\text{FFN}})), \quad F_l^{\text{proj}} \in \mathbb{R}^{B \times D \times H \times W}, \quad (25)$$

where $\text{Conv}_{1 \times 1}$ is a 1×1 convolution to reduce dimensionality and match spatial dimensions, and Upsample denotes bilinear interpolation.

Next, these scale-specific feature maps are concatenated along the channel dimension:

$$F_{\text{cat}} = [F_l^{\text{proj}}; F_s^{\text{FFN}}; F_f^{\text{FFN}}], \quad F_{\text{cat}} \in \mathbb{R}^{B \times 3D \times H \times W}. \quad (26)$$

The composite map F_{cat} undergoes an attention-based weighting mechanism:

$$A = \sigma(W_A * F_{\text{cat}} + b_A), \quad A \in \mathbb{R}^{B \times 3 \times H \times W}, \quad (27)$$

where $W_A \in \mathbb{R}^{3 \times 3D \times 1 \times 1}$, $b_A \in \mathbb{R}^3$, and σ is the sigmoid activation. A provides scale-specific attention weights.

The concatenated feature map (F_{cat}) is then split back into its scale-specific components (F_l^{split} , F_s^{split} , F_f^{split}):

$$F_l^{\text{split}} \in \mathbb{R}^{B \times D \times H \times W}, \quad (28)$$

$$F_s^{\text{split}} \in \mathbb{R}^{B \times D \times H \times W}, \quad (29)$$

$$F_f^{\text{split}} \in \mathbb{R}^{B \times D \times H \times W}. \quad (30)$$

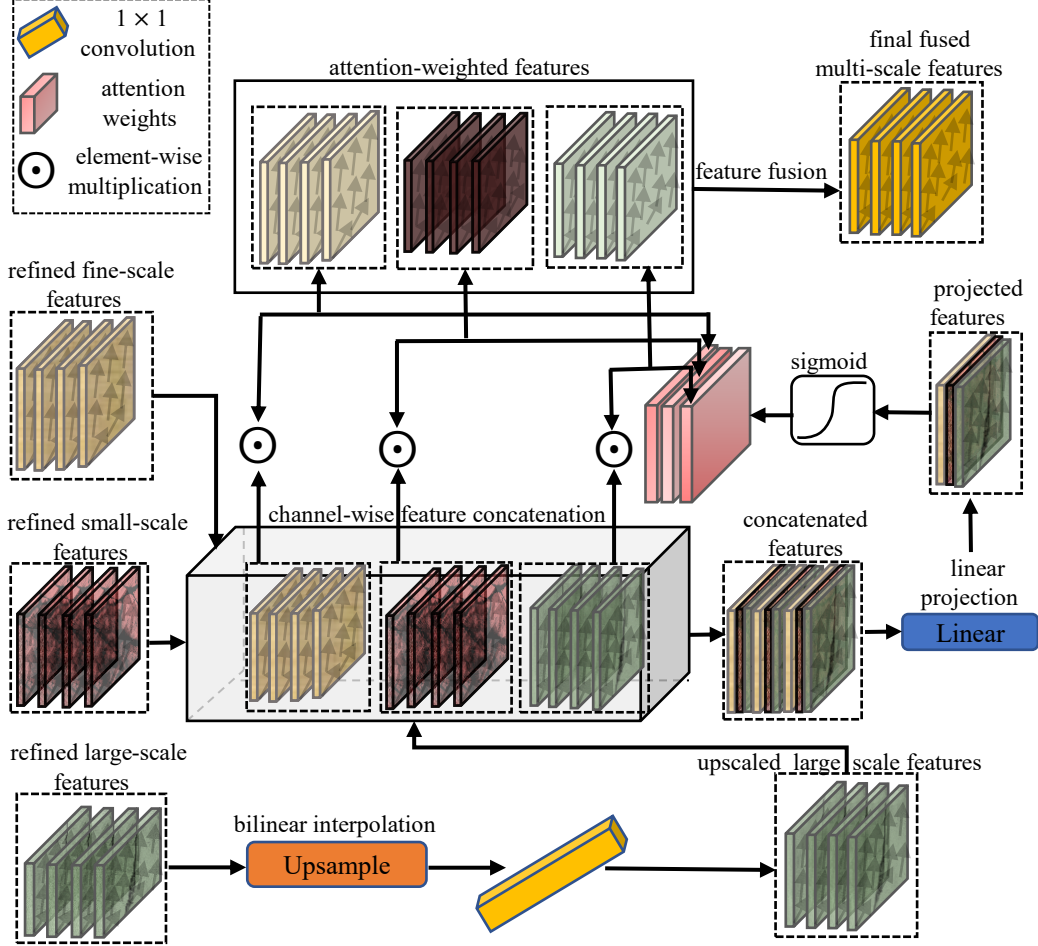


Figure 4: Attention-Guided Fusion Module.

Each component is weighted by its corresponding attention slice:

$$F_l^{\text{weighted}} = F_l^{\text{split}} \odot A[:, 0 : 1, :, :], \quad (31)$$

$$F_s^{\text{weighted}} = F_s^{\text{split}} \odot A[:, 1 : 2, :, :], \quad (32)$$

$$F_f^{\text{weighted}} = F_f^{\text{split}} \odot A[:, 2 : 3, :, :], \quad (33)$$

where \odot denotes element-wise multiplication.

Finally, the attention-weighted features are summed:

$$F_{\text{fused}} = F_l^{\text{weighted}} + F_s^{\text{weighted}} + F_f^{\text{weighted}}, \quad F_{\text{fused}} \in \mathbb{R}^{B \times D \times H \times W}. \quad (34)$$

This adaptive fusion strategy ensures multi-scale information is optimally integrated for accurate pavement crack segmentation.

3.6. Cross-scale Consistency Loss

To effectively address the challenge of self-supervised pavement crack segmentation, our framework incorporates a cross-scale consistency loss [43], composed of two complementary components: Inter-scale Consistency Loss and Intra-scale Consistency Loss. These losses are specifically designed to promote consistency and coherence in the learned feature representations across and within scales, significantly enhancing the robustness and accuracy of the model without reliance on manual annotations.

3.6.1. Inter-scale Consistency Loss

The Inter-scale Consistency Loss enforces similarity by minimizing the cosine distance between feature representations at different scales (fine, small, large), encouraging the model to learn scale-invariant crack features. Let $\mathbf{G}^f, \mathbf{G}^s, \mathbf{G}^l \in \mathbb{R}^d$ be the contextual feature vectors from fine, small, and large scales, respectively. The cosine similarity between two vectors is defined as:

$$\cos(\mathbf{G}^x, \mathbf{G}^y) = \frac{\mathbf{G}^x \cdot \mathbf{G}^y}{\|\mathbf{G}^x\| \|\mathbf{G}^y\|}, \quad x, y \in \{f, s, l\}. \quad (35)$$

The Inter-scale Consistency Loss $\mathcal{L}_{\text{inter}}(\mathbf{G}^x, \mathbf{G}^y)$ minimizes the cosine dissimilarity:

$$\mathcal{L}_{\text{inter}}(\mathbf{G}^x, \mathbf{G}^y) = 1 - \cos(\mathbf{G}^x, \mathbf{G}^y). \quad (36)$$

In practice, this loss is computed pairwise between fine–small and small–large scales, with a weighting factor:

$$\mathcal{L}_{\text{inter}} = \lambda_1 \left[\mathcal{L}_{\text{inter}}(\mathbf{G}^f, \mathbf{G}^s) + \mathcal{L}_{\text{inter}}(\mathbf{G}^s, \mathbf{G}^l) \right], \quad (37)$$

where λ_1 is the weighting factor for the inter-scale consistency loss.

3.6.2. Intra-scale Consistency Loss

The Intra-scale Consistency Loss improves internal consistency within each scale-specific feature representation. Given an attention map $\mathbf{A} \in \mathbb{R}^{L \times L}$ and the identity matrix $\mathbf{I} \in \mathbb{R}^{L \times L}$, we enforce \mathbf{A} to resemble \mathbf{I} via an L_1 loss:

$$\mathcal{L}_{\text{intra}}(\mathbf{A}) = \frac{1}{L^2} \sum_{i=1}^L \sum_{j=1}^L |A_{ij} - I_{ij}|. \quad (38)$$

1 We weight this loss as:

$$2 \quad \mathcal{L}_{\text{intra}} = \lambda_2 \mathcal{L}_{\text{intra}}(\mathbf{A}), \quad (39)$$

3 where λ_2 is the weighting factor for the intra-scale consistency loss.

4 3.6.3. Cross-Entropy Loss

5 To further support self-supervision, a pseudo-target was derived from the
6 model’s output itself $\mathbf{T} = [T_n] \in \{0, 1\}^N$ by taking the highest predicted
7 crack probability for each pixel:

$$8 \quad T_n = \begin{cases} 1, & \text{if } O_n \geq 0.5, \\ 0, & \text{otherwise,} \end{cases} \quad n = 1, \dots, N, \quad (40)$$

9 where $\mathbf{O} = [O_n] \in [0, 1]^N$ is the model’s output probability map and N is
10 the total number of pixels.

11 The binary Cross-Entropy Loss between \mathbf{O} and \mathbf{T} is then defined as

$$12 \quad \mathcal{L}_{\text{CE}}(\mathbf{O}, \mathbf{T}) = -\frac{1}{N} \sum_{n=1}^N \left[T_n \log(O_n) + (1 - T_n) \log(1 - O_n) \right], \quad (41)$$

13 where \mathcal{L}_{CE} denotes the cross-entropy loss, O_n is the predicted probability that
14 pixel n belongs to the crack class, and T_n is the corresponding pseudo-target
15 label.

16 3.6.4. Total Consistency Loss

17 Finally, the overall cross-scale consistency loss integrates all components:

$$18 \quad \mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}}(\mathbf{O}, \mathbf{T}) + \frac{1}{B} \sum_{b=1}^B \left(\mathcal{L}_{\text{inter}}^{(b)} + \mathcal{L}_{\text{intra}}^{(b)} \right), \quad (42)$$

19 where B is the batch size.

20 These self-supervisory losses ensure multi-scale representations remain co-
21 herent, internally consistent, and aligned with the model’s predictions, signif-
22 icantly enhancing segmentation accuracy without ground truth annotations
23 (masks).

1 4. Experimental details

2 In this section, the ten public datasets and the preprocessing steps used
3 to prepare data for training and testing are described. Each dataset in-
4 cludes unique crack patterns and imaging conditions, which enables testing
5 the model across varied scenarios. The routine adopted for training the pro-
6 posed model including the training settings and evaluation metrics is then
7 explained. Our study is complemented with ablation studies that helps iso-
8 late the contribution of each module in our proposed model.

9 4.1. Datasets and Preprocessing

10 Experiments were conducted on ten publicly available datasets to as-
11 sess the effectiveness of the proposed self supervised segmentation approach.
12 Each dataset contains images capturing different crack patterns across var-
13 ious surfaces, lighting conditions, and environmental scenarios. Utilizing
14 these diverse datasets allows us to effectively assess the generalization capa-
15 bility and robustness of our segmentation model under different real-world
16 conditions. Specifically, we employed the CFD dataset [44], Crack500 [45],
17 CrackTree200 [46], DeepCrack [22], Eugen Miller [47], Forest [48], GAPs384
18 [49], Rissbilder [50], Sylvie [51], and Volker [50]. Each dataset is characterized
19 by distinct types of cracks, surface materials, and varying illumination condi-
20 tions, making them well-suited for training and validating our self-supervised
21 model’s performance. Table 1 provides a concise overview of each dataset’s
22 key characteristics, highlighting the variations in crack patterns and imaging
23 environments.

24 Furthermore, each dataset was split into training and testing subsets,
25 maintaining an 80:20 split to ensure consistency in our evaluation protocol.
26 This division provides a fair comparison of our model’s predictive perfor-
27 mance on previously unseen data. To enhance the generalization and robust-
28 ness of our model, several data augmentation strategies were applied during
29 training. These augmentations included random horizontal and vertical flips,
30 rotation transformations, and scaling variations. These augmentation tech-
31 niques simulate the variability encountered in practical pavement inspection
32 scenarios, enabling the model to better generalize and accurately segment
33 cracks under diverse conditions.

Dataset	Crack Types	Surface Material	Lighting Conditions
CFD	Thin linear cracks	Asphalt pavement	Outdoor daylight, shadows, oil stains
Crack500	Hairline, wide cracks	Asphalt road surfaces	Mixed outdoor, varied weather
CrackTree200	Linear, alligator cracks	Asphalt pavement	Low contrast, uneven lighting
DeepCrack	Pavement, stone cracks	Asphalt concrete; stone	Daylight, some laser-lit
Eugen Miller	Random cracks	Tunnel concrete	Tunnel lighting
Forest	Thin linear cracks	Asphalt pavement	Outdoor daylight, shadows
GAPs384	Longitudinal, transverse, block	Asphalt roads	Dry daylight
Rissbilder	Architectural cracks	Concrete, masonry	Varied lighting
Sylvie	Linear, network cracks	Asphalt pavement	Outdoor varied lighting
Volker	Structural cracks	Concrete facades	Field conditions, well-lit

Table 1: Comparison of crack detection datasets

1 4.2. Implementation details

2 Our proposed self-supervised segmentation model was developed using
3 the PyTorch deep learning library. The AdamW optimizer was used with a
4 weight decay of 1×10^{-5} to mitigate potential overfitting. The initial learning
5 rate was established at 1×10^{-4} . To enhance training efficiency, a learning
6 rate scheduler that reduced the learning rate by a factor of 0.5 whenever
7 validation performance stagnated for five consecutive epochs was used.

8 Only the original images from the datasets were used during training,
9 explicitly excluding any segmentation masks to ensure a genuinely self su-
10 pervised learning scenario.. For all experiments, the datasets were randomly
11 divided into training and validation subsets, ensuring consistency and fair-
12 ness across evaluations. For model evaluation and performance benchmark-

ing, predicted segmentation masks were directly compared against the ground truth masks available within the validation datasets. Training was performed with a batch size of 8 for up to 500 epochs, incorporating early stopping with a patience of 100 epochs to prevent unnecessary computations and overfitting.

The performance of our model was compared with 14 state-of-the-art fully supervised segmentation models, including FCN [52], U-Net [53], U-Net++ [54], PSPNet [55], PAN [56], MAnet [57], LinkNet [58], FPN [59], DeepLabV3 [60], DeepLabV3+ [61], UPerNet [62], Segformer [63], and CrackFormer [26]. In total, all 14 segmentation models (the 13 fully-supervised baselines and our proposed self-supervised model) were each trained on all the 10 datasets, yielding a total of 140 experimental runs. Each baseline model was trained using the same experimental settings as our proposed model for consistency.

All computations were performed using the hardware and software configurations detailed in Table 2.

Component	Details
GPU	NVIDIA A40 (48 GB)
Framework	PyTorch 2.7
Programming Language	Python 3.9.12
CUDA Version	11.8
Optimizer	AdamW
Learning Rate Scheduler	Adaptive learning rate annealing

Table 2: Summary of the hardware and software environment used in experiments.

4.3. Evaluation metrics

Evaluation of the models were conducted using several standard metrics, including mean Intersection over Union (mIoU) and Dice coefficient. While these metrics assess overall overlap and similarity, additional metrics such as the XOR metric and Hammoud Distance (HM) were incorporated to further capture spatial disagreement and misalignment, providing deeper insights into the segmentation quality.

Mean Intersection over Union (mIoU). The mIoU metric calculates the average overlap between predicted masks and ground truth masks across different classes. For a particular class, IoU is computed by dividing the intersection of the prediction and ground truth by their union:

$$\text{IoU}_c = \frac{|P_c \cap G_c|}{|P_c \cup G_c|},$$

where P_c and G_c represent the predicted and ground truth pixel sets for class c . The mIoU then averages the IoU scores across all C classes:

$$\text{mIoU} = \frac{1}{C} \sum_{c=1}^C \text{IoU}_c.$$

Dice Coefficient. The Dice coefficient measures how closely the predicted segmentation aligns with the ground truth. It emphasizes regions with smaller or finer details, making it particularly useful for pavement cracks. Dice is computed as:

$$\text{Dice} = \frac{2 \times |P \cap G|}{|P| + |G|},$$

where P and G denote the predicted and actual sets of crack pixels, respectively.

XOR. The XOR metric quantifies discrepancies between the predicted and ground truth masks. It highlights areas exclusively classified as crack or non-crack in one mask but not in the other, thereby capturing mismatches effectively:

$$\text{XOR} = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W (P_{ij} \oplus G_{ij}),$$

where P_{ij} and G_{ij} are the predicted and ground truth binary labels at pixel (i, j) , and \oplus denotes the logical XOR operation. A lower XOR value indicates better segmentation performance.

Hammoud Distance (HM). The Hammoud distance evaluates the spatial dissimilarity between predicted and ground truth masks. It quantifies the extent of spatial misalignment or inconsistency, defined mathematically as:

$$\text{HM}(P, G) = \max \left\{ \max_{p \in P} \min_{g \in G} \|p - g\|_2, \max_{g \in G} \min_{p \in P} \|g - p\|_2 \right\},$$

where P and G are the sets of crack-pixel coordinates in the predicted and ground truth masks, respectively, and $\|\cdot\|_2$ denotes the Euclidean distance.

1 Lower HM values indicate higher segmentation accuracy, emphasizing better
2 spatial alignment between prediction and ground truth.

3 5. Results and Discussion

4 5.1. Quantitative results

5 Table 3 provides a comprehensive performance comparison between our
6 proposed self-supervised segmentation model and various state-of-the-art fully
7 supervised methods across ten pavement crack datasets. The assessment uti-
8 lized four metrics: mean mIoU, Dice score, XOR, and Hammoud distance.

9 Our proposed model achieved significantly superior performance com-
10 pared to existing methods across nearly all datasets. On the CFD dataset,
11 our model recorded a remarkable mIoU of 0.8875 and a Dice score of 0.9340,
12 substantially surpassing U-Net++ (the best competitor among supervised
13 methods) with an mIoU of 0.5257 and Dice of 0.6869. Moreover, our method
14 exhibited notably lower XOR (0.6138) and HM (0.6138) values, demonstrat-
15 ing minimal spatial discrepancies and better alignment of predicted cracks.

16 Similarly, on the CRACK500 dataset, our self-supervised model achieved
17 outstanding results, with an mIoU of 0.9332 and Dice of 0.9647, greatly
18 exceeding the best-performing supervised method, Linknet, which attained
19 mIoU and Dice scores of 0.6449 and 0.7838, respectively. The XOR and HM
20 values of 0.1957 and 0.1629 further confirmed our model’s exceptional spatial
21 accuracy and lower false-positive rates.

22 The DeepCrack and Forest datasets also highlighted our model’s robust-
23 ness, achieving mIoU values of 0.8217 and 0.8167, respectively. These values
24 substantially outperformed the next best methods, U-Net++ (0.7166 mIoU
25 on DeepCrack) and U-Net (0.5392 mIoU on Forest). The high Dice scores
26 of 0.8952 (DeepCrack) and 0.8896 (Forest) reinforced the accuracy improve-
27 ments enabled by our framework’s modules.

28 The CrackTree200 dataset, known for challenging crack patterns, demon-
29 strated a significant performance leap with our model attaining an mIoU of
30 0.8670 and a Dice score of 0.9223, far superior to the second-best, U-Net, with
31 0.4861 mIoU and 0.6498 Dice scores. Similarly, our approach excelled on the
32 GAPs dataset, delivering an mIoU of 0.8096 and Dice of 0.8854, surpassing
33 Segformer’s mIoU of 0.4016 and Dice of 0.5714 by a large margin.

34 Performance on Eugen Miller was also notably superior, with our model
35 achieving an mIoU of 0.8451 and Dice of 0.9071, outperforming U-Net++

Table 3: Validation Results of our proposed and other models across all datasets.

Model	CFD				CRACK500				DeepCrack				Forest			
	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow
FCN	0.3509	0.5155	1.7027	0.6491	0.5750	0.7088	0.8035	0.4250	0.6123	0.7448	0.7623	0.3877	0.3796	0.5460	1.5054	0.6204
Segformer	0.4734	0.6398	1.6703	0.6421	0.6381	0.7776	0.6747	0.4249	0.6556	0.7910	0.6538	0.3694	0.4743	0.6406	1.4148	0.6110
PSPNet	0.3552	0.5191	1.2426	0.6695	0.6207	0.7653	0.6716	0.4197	0.6576	0.7925	1.0336	0.4658	0.3840	0.5487	1.2346	0.6523
Linknet	0.4648	0.6232	1.7669	0.6835	0.6449	0.7838	0.6987	0.4263	0.7047	0.8243	1.3199	0.5122	0.4883	0.6412	1.5864	0.6554
FPN	0.4212	0.5833	1.3548	0.6373	0.6316	0.7740	0.6503	0.4348	0.6783	0.8078	0.9757	0.4544	0.4343	0.5976	1.2351	0.6181
Unet	0.5123	0.6743	2.2378	0.7018	0.6398	0.7796	0.7171	0.4174	0.7059	0.8255	1.2395	0.4970	0.5392	0.6974	1.5924	0.6441
PAN	0.3654	0.5166	78.3467	0.9853	0.6231	0.7674	0.6937	0.4240	0.6535	0.7891	0.8343	0.4317	0.4580	0.6261	1.2560	0.6160
DeepLabV3	0.3167	0.4744	1.9940	0.7353	0.6420	0.7816	0.6699	0.4132	0.6807	0.8092	1.0298	0.4622	0.4821	0.6477	1.3812	0.6352
DeepLabV3Plus	0.3863	0.5526	1.3633	0.6674	0.6376	0.7784	0.7092	0.4166	0.6641	0.7961	0.9354	0.4437	0.4610	0.6273	1.2988	0.6079
CrackFormer	0.4494	0.6178	1.4242	0.6186	0.6306	0.7702	0.8168	0.4435	0.6179	0.7592	0.6708	0.3755	0.4450	0.6139	1.3694	0.6201
UPerNet	0.4674	0.6352	1.5864	0.6286	0.6432	0.7813	0.7292	0.4255	0.6581	0.7929	0.6514	0.3644	0.4724	0.6381	1.4423	0.6179
MAnet	0.5166	0.6761	1.6948	0.6589	0.6341	0.7754	0.7453	0.4284	0.7004	0.8216	1.3568	0.5158	0.5174	0.6752	1.6298	0.6495
UnetPlusPlus	0.5257	0.6869	1.7524	0.6635	0.6443	0.7834	0.6703	0.4202	0.7166	0.8344	1.2447	0.4983	0.5457	0.7044	1.5658	0.6388
Crack-Segmenter (ours)	0.8875	0.9340	0.6138	0.6138	0.9332	0.9647	0.1957	0.1629	0.8217	0.8952	0.3685	0.2405	0.8167	0.8896	0.5836	0.4003
Crack-Segmenter-v0 (ours)	0.5435	0.7025	0.8597	0.4565	0.8572	0.9210	0.1782	0.1428	0.6955	0.7898	0.5003	0.3045	0.5723	0.7257	0.7695	0.4277
Crack-Segmenter-v1 (ours)	0.5432	0.7023	0.8603	0.4568	0.5360	0.6965	0.8708	0.4640	0.5698	0.7026	0.7064	0.4303	0.5554	0.7102	0.7972	0.4446
Crack-Segmenter-v2 (ours)	0.5360	0.6965	0.8708	0.4640	0.5711	0.7247	0.7699	0.4289	0.6836	0.8013	0.5388	0.3164	0.5711	0.7247	0.7699	0.4289
Crack-Segmenter-v3 (ours)	0.5779	0.7302	0.7300	0.4221	0.7457	0.8407	0.3033	0.2543	0.6797	0.7996	0.5590	0.3203	0.5153	0.6617	0.7999	0.4847

Model	CrackTree200				GAPs				Eugen Miller			
	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow
FCN	0.065	0.1215	13.8985	0.9353	0.3648	0.5128	1.4705	0.6352	0.6215	0.7660	0.5168	0.3785
Segformer	0.2916	0.4509	12.3826	0.9299	0.4016	0.5714	1.3561	0.6268	0.5973	0.7402	0.6956	0.4543
PSPNet	0.2930	0.4477	6.7798	0.9178	0.2907	0.4432	1.2396	0.7015	0.6062	0.7479	0.5292	0.3905
Linknet	0.4804	0.6411	12.2496	0.9355	0.3996	0.5659	1.7396	0.6894	0.6021	0.7414	0.5740	0.4178
FPN	0.2962	0.4469	7.6561	0.9197	0.3068	0.4624	1.3438	0.7369	0.4682	0.6161	0.6412	0.4997
Unet	0.4861	0.6498	11.4862	0.9301	0.3832	0.5508	1.5567	0.6753	0.6315	0.7655	0.4998	0.3771
PAN	0.2755	0.4284	6.7856	0.9185	0.2761	0.4259	1.2340	0.7191	0.5796	0.7265	0.5651	0.4189
DeepLabV3	0.3178	0.4708	8.4792	0.9221	0.3685	0.5343	1.5615	0.6844	0.6291	0.7633	0.5928	0.4039
DeepLabV3Plus	0.2788	0.4243	9.0466	0.9205	0.3168	0.4705	1.5988	0.7019	0.6047	0.7457	0.5501	0.3999
CrackFormer	0.2656	0.4189	12.3523	0.9348	0.3470	0.5112	1.6896	0.6595	0.4756	0.6352	0.8786	0.6001
UPerNet	0.2569	0.4049	12.9757	0.9328	0.4029	0.5724	1.3435	0.6121	0.5835	0.7295	0.6135	0.4350
MAnet	0.4198	0.5773	11.9556	0.9330	0.3793	0.5414	1.6618	0.6844	0.5892	0.7322	0.5547	0.4110
UnetPlusPlus	0.4255	0.5927	8.9129	0.9183	0.3926	0.5561	1.6479	0.6858	0.7072	0.8214	0.4798	0.3505
Crack-Segmenter (ours)	0.8670	0.9223	5.0154	0.8409	0.8096	0.8854	0.6542	0.4497	0.8451	0.9071	0.3952	0.3669
Crack-Segmenter-v0 (ours)	0.1572	0.2716	5.1871	0.8428	0.6353	0.7711	0.6264	0.3647	0.8597	0.9245	0.1635	0.1403
Crack-Segmenter-v1 (ours)	0.1474	0.2568	5.8197	0.8526	0.6312	0.7682	0.6293	0.3688	0.8016	0.8860	0.2237	0.1984
Crack-Segmenter-v2 (ours)	0.1406	0.2453	5.1946	0.8594	0.6353	0.7711	0.6264	0.3647	0.6860	0.8080	0.3556	0.3140
Crack-Segmenter-v3 (ours)	0.1385	0.2430	5.0391	0.8615	0.4953	0.6450	0.7485	0.5047	0.7230	0.8369	0.3078	0.2770

Model	Rissbilder				Sylvie				Volker			
	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow
FCN	0.5098	0.6707	0.8074	0.4902	0.7498	0.8439	0.3016	0.2502	0.6658	0.7981	0.4685	0.3342
Segformer	0.5169	0.6810	0.7938	0.4912	0.6600	0.7899	0.3960	0.2939	0.6632	0.7973	0.4732	0.3407
PSPNet	0.5042	0.6685	0.7935	0.5339	0.5548	0.6998	0.4102	0.3592	0.6865	0.8132	0.4625	0.3401
Linknet	0.6084	0.7497	0.8332	0.4941	0.6430	0.7750	0.5250	0.3740	0.7235	0.8344	0.4871	0.3439
FPN	0.5822	0.7348	0.7652	0.4932	0.5993	0.7456	0.4893	0.3898	0.7026	0.8242	0.4505	0.3341
Unet	0.6449	0.7836	0.8240	0.4836	0.6558	0.7891	0.4501	0.3269	0.7454	0.8532	0.4788	0.3349
PAN	0.5415	0.7015	0.7529	0.5039	0.6073	0.7504	0.4030	0.3438	0.6790	0.8076	0.4574	0.3405
DeepLabV3	0.5970	0.7474	0.7754	0.4857	0.6639	0.7968	0.3633	0.2972	0.7213	0.8378	0.4448	0.3280
DeepLabV3Plus	0.5650	0.7210	0.7781	0.4965	0.6508	0.7826	0.3683	0.3038	0.6871	0.8139	0.4748	0.3467
CrackFormer	0.4638	0.6313	0.8590	0.5126	0.4829	0.6402	0.6936	0.5363	0.6182	0.7631	0.5042	0.3614
UPerNet	0.5150	0.6794	0.8045	0.4872	0.6663	0.7946	0.3450	0.2680	0.6682	0.8009	0.4628	0.3331
MAnet	0.6370	0.7771	0.8190	0.4811	0.6254	0.7642	0.4140	0.3246	0.7027	0.8245	0.5236	0.3598
UnetPlusPlus	0.6564	0.7920	0.7964	0.4744	0.6719	0.8021	0.5109	0.3890	0.7641	0.8659	0.4655	0.3274
Crack-Segmenter (ours)	0.7998	0.8785	0.5047	0.4344	0.8891	0.9352	0.3341	0.3122	0.6707	0.7955	0.6434	0.6240
Crack-Segmenter-v0 (ours)	0.8075	0.8932	0.2398	0.1925	0.9063	0.9503	0.1060	0.0937	0.8746	0.9330	0.1437	0.1254
Crack-Segmenter-v1 (ours)	0.8079	0.8934	0.2394	0.1921	0.9065	0.9504	0.1057	0.0935	0.8761	0.9338	0.1421	0.1239
Crack-Segmenter-v2 (ours)	0.7894	0.8810	0.2611	0.2106	0.7873	0.8745	0.2326	0.2127	0.8637	0.9265	0.1560	0.1363
Crack-Segmenter-v3 (ours)	0.6740	0.7979	0.3898	0.3260	0.8802	0.9354	0.1344	0.1198	0.5804	0.7303	0.4637	0.4196

1 with 0.7072 mIoU and 0.8214 Dice. These results indicate our model’s effec-
 2 tiveness even in diverse structural conditions.

3 However, the performance on the Volker dataset was slightly lower, with
 4 an mIoU of 0.6707 and Dice of 0.7955, compared to U-Net++, which achieved
 5 an mIoU of 0.7641 and Dice of 0.8659. This suggests room for further re-
 6 finement of our multi-scale attention mechanisms to handle highly variable
 7 crack structures effectively. Figure 5 shows mIoU and Dice scores of Crack-
 8 Segmenter and all the baseline models across the different datasets sum-
 9 marised in a radar plot.

10 The consistently superior quantitative results clearly demonstrate the ef-
 11 fectiveness of the integrated modules within our self-supervised segmentation
 12 framework. Specifically, the Scale-Adaptive Embedder efficiently captures
 13 comprehensive multi-scale feature details; the Directional Attention Trans-
 14 former emphasizes linear crack structures, improving the detection accuracy;
 15 and the Attention-Guided Fusion optimally merges multi-scale features, en-
 16 hancing overall segmentation performance. Collectively, these modules fa-
 17 cilitate accurate, annotation-free crack segmentation, significantly advancing
 18 the state-of-the-art in pavement distress assessment.

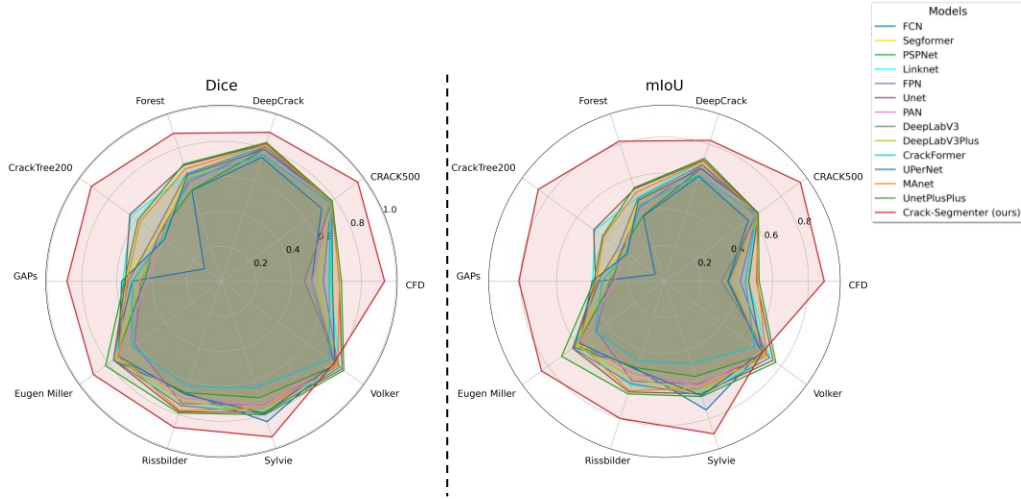


Figure 5: Radar plots for validation Dice and mIoU scores of all the baseline models and Crack-Segmenter across all the 10 datasets.

1 5.2. Ablation studies

2 To evaluate the individual contributions of each proposed module within
3 our self-supervised segmentation framework, we conducted comprehensive
4 ablation experiments on the DeepCrack dataset. This systematic analysis
5 isolates the effectiveness of each component: the Scale-Adaptive Embedder
6 (SAE), Directional Attention Transformer (DAT), and Attention-Guided Fu-
7 sion (AGF), and examines their interactions when combined in various con-
8 figurations.

9 5.2.1. Experimental Setup

10 We evaluated six distinct architectural variants using mIoU and Dice
11 score as primary metrics. Each variant represents a specific combination of
12 the proposed modules, allowing us to quantify both individual contributions
13 to segmentation performance. The following variants were systematically
14 evaluated:

15 **Baseline:** Standard U-Net architecture with ResNet-18 encoder, repre-
16 senting conventional supervised segmentation without self-supervised com-
17 ponents.

18 **Crack-Segmenter:** Complete architecture incorporating all three pro-
19 posed modules.

20 **Crack-Segmenter-v0:** Integration of SAE module only, focusing on
21 multi-scale representation learning capabilities.

22 **Crack-Segmenter-v1:** Combination of SAE and DAT modules, empha-
23 sizing directional attention alongside multi-scale features.

24 **Crack-Segmenter-v2:** Integration of SAE and AGF modules, combin-
25 ing multi-scale embeddings with attention-guided fusion.

26 **Crack-Segmenter-v3:** Combination of DAT and AGF modules without
27 multi-scale embeddings.

28 5.2.2. Analysis

29 To assess the individual contributions of each proposed modules, we con-
30 ducted systematic ablation experiments on the DeepCrack dataset. These
31 studies are crucial for understanding the role each component plays in the
32 overall performance of our self-supervised segmentation framework. All ab-
33 lation variants were evaluated using the four key metrics: mIoU, Dice score,
34 XOR Score, and HM Score.

35 We begin by defining the **Baseline** model, which consists of a standard
36 U-Net architecture with a ResNet-18 encoder and none of the proposed self-

1 supervised modules. This configuration serves as a conventional supervised
2 benchmark. As shown in Table 4, it achieved an mIoU of 0.5866 and a Dice
3 score of 0.7258. These results establish a foundational point of comparison
4 for the proposed architectural variations.

5 When only the SAE module was included in the architecture (**Crack-**
6 **Segmenter-v0**), we observed a substantial improvement in both metrics,
7 reaching an mIoU of 0.6955 and a Dice score of 0.7898. This confirms the
8 strong impact of multi-scale representations in capturing crack structures of
9 varying widths. SAE provides the model with diverse spatial resolutions,
10 enabling it to better recognize both fine-grained and coarse crack patterns,
11 which are typically missed in single-scale encoders.

12 In the **Crack-Segmenter-v1** variant, we added the DAT module on top
13 of SAE while leaving out AGF. Interestingly, this combination led to a drop
14 in performance (mIoU of 0.5698 and Dice of 0.7026), falling even below the
15 baseline. This suggests that applying directional attention without a proper
16 fusion mechanism may not be sufficient for effective feature integration. DAT
17 focuses on enhancing linear crack continuity through spatially aware convo-
18 lutions, but without adaptive fusion to resolve scale-level redundancies, it
19 may introduce conflicting or misaligned representations.

20 The combination of SAE and AGF in **Crack-Segmenter-v2** demon-
21 strated stronger performance, with 0.6836 mIoU and 0.8013 Dice. These
22 results reinforce the importance of adaptive feature fusion when dealing
23 with multi-scale representations. AGF dynamically learns to assign relevance
24 weights across different spatial resolutions, ensuring optimal use of both local
25 and contextual information during segmentation. The effective integration
26 of scale-specific details allows the model to better adapt to irregular crack
27 geometries and surrounding textures.

28 We then evaluated **Crack-Segmenter-v3**, which includes only DAT and
29 AGF, excluding SAE. This variant resulted in an mIoU of 0.6797 and Dice
30 score of 0.7996. Although this is an improvement over the baseline, the ab-
31 sence of SAE resulted in weaker feature diversity, reducing the ability of the
32 fusion and attention mechanisms to operate on rich, scale-aware represen-
33 tations. This again highlights the importance of SAE as a foundation for
34 multi-scale learning.

35 Finally, the full model (**Crack-Segmenter**) that integrates all three
36 modules: SAE, DAT, and AGF achieved the highest performance across
37 all metrics: 0.8217 mIoU, 0.8952 Dice, 0.3685 XOR, and 0.2405 HM. These
38 results confirm that the complete architecture benefits from the combined

1 strengths of its components. SAE enables comprehensive multi-scale feature
 2 extraction, DAT enhances directional continuity, and AGF fuses these fea-
 3 tures adaptively. Together, this integration result in a robust and coherent
 4 segmentation map that generalizes well across complex crack structures and
 5 challenging visual conditions.

Table 4: Ablation study results on the DeepCrack dataset. Each variant systematically evaluates different module combinations to assess individual contributions.

Model Variant	Metric				Module Components		
	mIoU \uparrow	Dice \uparrow	XOR \downarrow	HM \downarrow	SAE	DAT	AGF
Baseline	0.5866	0.7258	0.7675	0.4134	\times	\times	\times
Crack-Segmenter	0.8217	0.8952	0.3685	0.2405	\checkmark	\checkmark	\checkmark
Crack-Segmenter-v0	0.6955	0.7898	0.5003	0.3405	\checkmark	\times	\times
Crack-Segmenter-v1	0.5698	0.7026	0.7064	0.4303	\checkmark	\checkmark	\times
Crack-Segmenter-v2	0.6836	0.8013	0.5388	0.3164	\checkmark	\times	\checkmark
Crack-Segmenter-v3	0.6797	0.7996	0.5590	0.3203	\times	\checkmark	\checkmark

6 These ablation results validate the design of our self-supervised architec-
 7 ture and demonstrate how each module contributes to enhancing segmenta-
 8 tion quality. The clear performance gains of the full model emphasize the
 9 necessity of combining multi-scale embeddings, directional attention, and
 10 guided fusion in an integrated manner for accurate pavement crack segmen-
 11 tation.

12 5.3. Model Explainability

13 Interpreting model behavior is essential for validating its reliability, es-
 14 pecially in safety-critical applications such as pavement crack detection. To
 15 assess whether the proposed *Crack-Segmenter* attends to relevant spatial
 16 structures, attention map visualizations were employed to analyze its focus
 17 during inference.

18 Attention maps were extracted from the final transformer block at each
 19 of the three spatial scales (fine, small, and large) in the architecture. These
 20 layers capture the model’s most refined representations across different reso-
 21 lutions. For blocks with multiple attention heads, attention weights were av-
 22 eraged across heads to produce a single map per scale, ensuring clarity while
 23 preserving dominant spatial patterns. All attention maps were then normal-
 24 ized to the range $[0, 1]$ and overlaid as heatmaps on the original DeepCrack

images, where color intensity reflects attention strength. Regions receiving strong attention appear warmer (e.g., red or yellow), while areas with low attention are cooler (e.g., blue), enabling clear visual identification of the model’s focus during segmentation.

Figure 6 presents example attention visualizations for sample images in the DeepCrack dataset. The model consistently concentrates attention along actual crack regions, while suppressing distractors such as shadows, surface texture, and background noise. This focused behavior highlights the effectiveness of the Directional Attention Transformer and Attention-Guided Fusion modules in guiding the network toward semantically meaningful features. These visualizations confirm that the model captures relevant structural cues without explicit supervision, reinforcing both the design rationale and the effectiveness of the proposed framework.

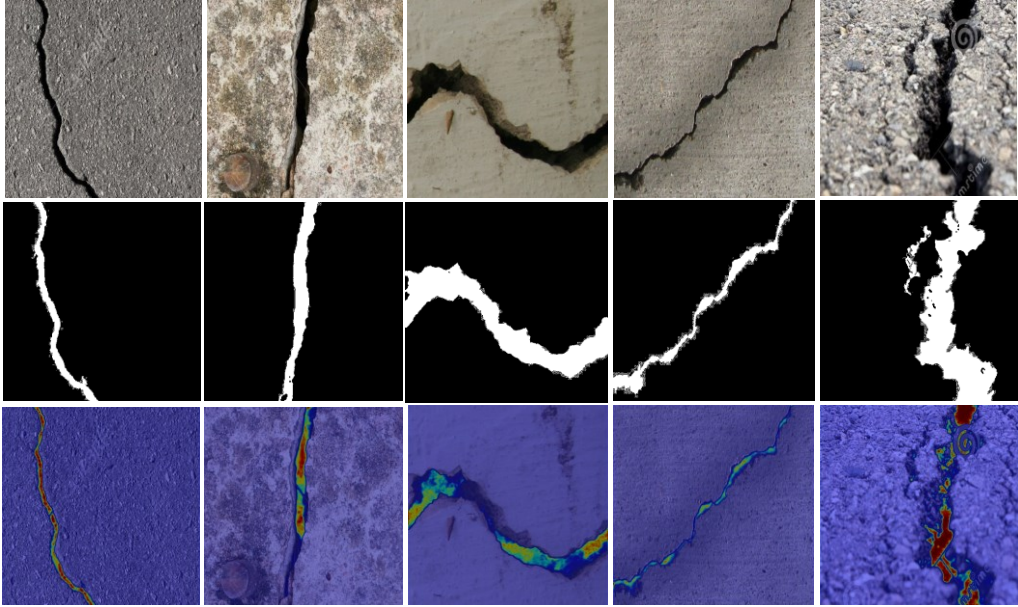


Figure 6: Attention map visualization of *Crack-Segmenter* on samples images from DeepCrack dataset. Top row: input images. Middle row: predicted masks. Bottom row: attention map overlays highlighting regions the model focused on during crack segmentation.

5.4. Statistical Analysis

To rigorously assess the statistical significance of our proposed *Crack-Segmenter*, we conducted detailed statistical analysis comparing its perfor-

1 mance against state-of-the-art segmentation methods. These analysis focused
2 on evaluating consistency and superiority across multiple segmentation met-
3 rics, including mIoU, Dice score, XOR, and HM. Table 5 summarizes the
4 mean performance and standard deviations of each method averaged over all
5 datasets.

6 From Table 5, it is evident that the proposed *Crack-Segmenter* achieved
7 consistently superior results across all evaluation metrics. Specifically, our
8 model attained the highest mIoU (0.8340 ± 0.0714) and Dice score (0.9008
9 ± 0.0458), indicating significantly more accurate segmentation performance
10 compared to other baseline models. Additionally, it showed the lowest XOR
11 (0.9309 ± 1.4432) and HM (0.4446 ± 0.2013) values, reflecting fewer misclas-
12 sifications and better spatial alignment with ground-truth crack structures.

Table 5: Mean Performance and Standard Deviation by Model accross all datasets

Model	mIoU [†]	Dice [†]	XOR [‡]	HM [‡]
DeepLabV3	0.5419 ± 0.1569	0.6863 ± 0.1434	1.7292 ± 2.4292	0.5367 ± 0.2004
DeepLabV3Plus	0.5252 ± 0.1525	0.6712 ± 0.1434	1.7123 ± 2.6091	0.5305 ± 0.1915
FCN	0.4894 ± 0.2017	0.6228 ± 0.2123	2.2237 ± 4.1299	0.5106 ± 0.2018
FPN	0.5121 ± 0.1475	0.6593 ± 0.1378	1.5562 ± 2.1700	0.5518 ± 0.1777
Linknet	0.5760 ± 0.1104	0.7180 ± 0.0930	2.1780 ± 3.5753	0.5532 ± 0.1849
MAnet	0.5722 ± 0.1111	0.7165 ± 0.0976	2.1355 ± 3.4871	0.5447 ± 0.1866
PAN	0.5059 ± 0.1529	0.6540 ± 0.1469	9.1329 ± 24.3942	0.5702 ± 0.2333
PSPNet	0.4953 ± 0.1525	0.6446 ± 0.1428	1.4397 ± 1.9048	0.5450 ± 0.1865
Segformer	0.5372 ± 0.1273	0.6880 ± 0.1146	2.0511 ± 3.6555	0.5184 ± 0.1894
UPerNet	0.5334 ± 0.1366	0.6829 ± 0.1266	2.0954 ± 3.8467	0.5105 ± 0.1941
Unet	0.5944 ± 0.1111	0.7369 ± 0.0923	2.1082 ± 3.3477	0.5388 ± 0.1948
UnetPlusPlus	0.6050 ± 0.1265	0.7439 ± 0.1050	1.8047 ± 2.5484	0.5366 ± 0.1867
CrackFormer	0.4796 ± 0.1187	0.6361 ± 0.1120	2.1258 ± 3.6134	0.5662 ± 0.1662
Crack-Segmenter (Ours)	0.8340 ± 0.0714	0.9008 ± 0.0458	0.9309 ± 1.4432	0.4446 ± 0.2013

[†]Higher values indicate better performance (mIoU, Dice). [‡]Lower values indicate better performance (XOR, HM). **Bold** values indicate best performance for each metric.

13 To determine whether these performance improvements were statistically
14 meaningful, paired *t*-tests were conducted between our model and each base-
15 line method, as shown in Table 6. Our model exhibited statistically signif-
16 icant improvements in both mIoU and Dice score compared to all baseline
17 methods, as indicated by positive mean differences and very low *p*-values
18 ($p < 0.01$ and mostly $p < 0.001$). For instance, when compared with
19 commonly used segmentation models such as U-Net and DeepLabV3, the
20 proposed method demonstrated substantial mean improvements in mIoU
21 ($+0.2396$ and $+0.2921$, respectively) and Dice score ($+0.1639$ and $+0.2144$,
22 respectively), all statistically significant at $p < 0.01$.

Furthermore, the highest statistical significance was noted against the CrackFormer model, with a mean difference of +0.3544 in mIoU and +0.2647 in Dice score (both $p < 0.001$). This strong statistical evidence confirms the effectiveness of integrating multi-scale embeddings, directional attention mechanisms, and adaptive fusion within our architecture. These improvements likely stem from the enhanced capability of our model to accurately capture diverse crack patterns and textures, as demonstrated consistently across datasets. Figure 7 shows distribution of the dice and mIoU scores of Crack-Segmenter and the baseline models.

Table 6: Statistical significance tests comparing Crack-Segmenter against baseline methods. Mean differences, t-statistics, p-values, and significance levels are reported for each evaluation metric. Positive mean differences for mIoU and Dice indicate superior performance.

Metric	Statistic	Models					
		FCN	SegFormer	PSPNet	LinkNet	FPN	U-Net
mIoU	mean diff.	+0.3446	+0.2968	+0.3387	+0.2581	+0.3220	+0.2396
	t-statistic	4.829	6.091	5.840	5.560	5.657	5.089
	p-value	0.000935	0.000181	0.000247	0.000352	0.000311	0.000655
	significance	***	***	***	***	***	***
Dice	mean Diff.	+0.2779	+0.2128	+0.2562	+0.1828	+0.2415	+0.1639
	t-statistic	3.886	5.204	5.084	5.125	4.927	4.586
	p-value	0.003696	0.000561	0.000659	0.000624	0.000816	0.001317
	significance	**	***	***	***	***	**

Metric	Statistic	PAN	DeepLabV3	DeepLabV3+	CrackFormer	UPerNet	MANet	U-Net++
mIoU	mean diff.	+0.3281	+0.2921	+0.3088	+0.3544	+0.3006	+0.2618	+0.2290
	t-statistic	5.682	4.852	5.394	7.523	5.837	5.649	4.449
	p-value	0.000301	0.000906	0.000437	0.000036	0.000248	0.000314	0.001603
	significance	***	***	***	***	***	***	**
Dice	mean Diff.	+0.2468	+0.2144	+0.2295	+0.2647	+0.2178	+0.1843	+0.1568
	t-statistic	4.799	4.180	4.573	6.516	4.887	4.993	3.968
	p-value	0.000975	0.002377	0.001341	0.000109	0.000863	0.000746	0.003264
	significance	***	**	*	***	***	***	*

Note: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$. Paired t-tests assume the same dataset was used across all models.

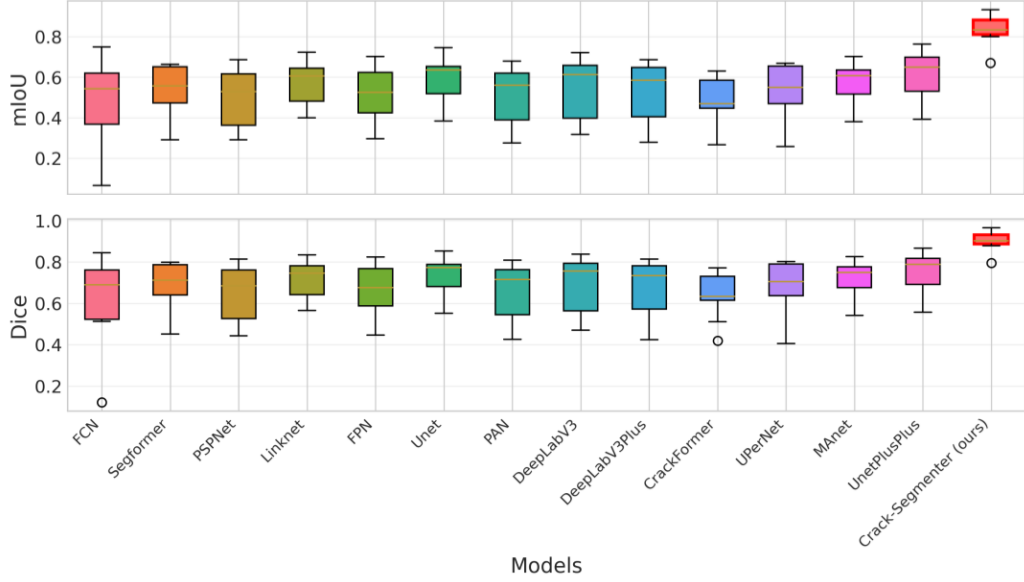


Figure 7: Box plot showing the mIoU and Dice Scores of Crack-Segmenter against other model baselines.

5.5. Qualitative Results

Figure 8 presents visual comparisons between our self-supervised method and supervised baselines across nine datasets. Each row displays a different dataset, enabling comprehensive assessment of segmentation performance across diverse crack patterns and imaging conditions.

On the CFD dataset with thin linear cracks, our method produces clean, continuous segmentations closely matching ground truth. While baseline methods capture general crack structure, they show inconsistent crack width. Our multi-scale embeddings effectively capture fine details at appropriate resolutions, resulting in well-defined boundaries without excessive dilation.

For CRACK500’s vertical cracks with varying widths, our approach maintains consistent segmentation quality throughout, accurately preserving narrow and wider sections. The Directional Attention Transformer emphasizes vertical continuity, leading to better boundary precision and endpoint detection compared to baselines.

The DeepCrack dataset contains horizontal cracks with irregular edges. All methods perform competently, but ours shows superior noise suppression while maintaining sharp boundaries. The Attention-Guided Fusion module balances detail preservation with contextual consistency, reducing back-

1 ground artifacts present in baseline predictions. On Forest’s multiple branch-
2 ing cracks, our method excels at maintaining connectivity at junctions, pre-
3 serving complete network structure.

4 CrackTree200 presents dense interconnected crack networks. Our ap-
5 proach captures intricate patterns while maintaining clear boundaries be-
6 tween closely spaced cracks. Multi-scale processing enables simultaneous
7 detection of major branches and fine subsidiary cracks, whereas baselines ei-
8 ther over-segment (merging adjacent cracks) or under-segment (missing finer
9 branches).

10 For GAPS’ diagonal cracks with variable contrast, our method demon-
11 strates robust performance even in low-contrast regions, maintaining consis-
12 tent quality along entire crack lengths. The self-supervised learning identifies
13 patterns based on structural characteristics rather than intensity contrasts
14 alone.

15 The Eugen Miller dataset contains curved tunnel cracks under varying
16 lighting. Our approach provides smoother curve representations with fewer
17 discretization artifacts than baselines. The directional attention adapts to
18 changing crack orientations effectively.

19 On Rissbilder’s intersecting cracks, our method accurately segments both
20 primary cracks and intersections without creating artificial connections or
21 gaps. This demonstrates effective handling of complex spatial relationships
22 through attention-guided fusion.

23 For Sylvie’s prominent horizontal cracks, all methods achieve satisfactory
24 results, but ours produces the most consistent width representation along
25 the entire length. This consistency stems from balanced feature extraction
26 across scales.

27 These qualitative results confirm that our self-supervised approach achieves
28 segmentation quality matching or exceeding supervised methods while elim-
29 inating manual annotation requirements. The visual evidence demonstrates
30 superior crack continuity preservation, consistent boundary maintenance,
31 and adaptation to diverse crack morphologies across different pavement con-
32 ditions.

Dataset	Input image	Ground truth mask	Crack segmenter (ours)	CrackFormer	Segformer	Unet	UnetPlusPlus
CFD							
CRACK 500							
DeepCrack							
Forest							
Crack Tree200							
GAPs							
Eugen Miller							
Rissbilder							
Sylvie							

Figure 8: Qualitative comparison of predicted segmentation masks across multiple datasets. Each row shows the input image, ground truth, and predictions from the chosen baseline models (CrackFormer, Segformer, Unet, UnetPlusPlus) and the proposed *Crack-Segmenter*.

1 6. Practical Applications

2 The proposed *Crack-Segmenter* framework offers several practical bene-
3 fits in construction and infrastructure management. Efficient pavement crack
4 detection is vital for ensuring road safety, enhancing infrastructure durability,
5 and optimizing asset management. Eliminating manual annotation costs en-
6 ables transportation agencies to conduct regular, comprehensive inspections
7 more affordably.

8 A primary application is preventive pavement maintenance planning. The
9 accurate crack segmentation maps from our model support early identifica-
10 tion and timely repair of pavement defects, enabling efficient budget alloca-
11 tion and reducing long-term maintenance expenses.

12 Additionally, our method seamlessly integrates with automated inspec-
13 tion technologies, such as unmanned aerial vehicles (UAVs) and vehicle-
14 mounted systems, enabling rapid, continuous monitoring of large road net-
15 works with minimal human involvement. This automation increases inspec-
16 tion efficiency and consistency across expansive areas.

17 Construction and engineering firms can apply this approach for proactive
18 monitoring of new or rehabilitated pavement surfaces, swiftly identifying
19 potential structural weaknesses. This early detection allows for timely cor-
20 rective measures, extending pavement lifespan and enhancing public safety.

21 Finally, the annotation-free nature of our approach facilitates adoption
22 in resource-constrained regions, democratizing access to advanced pavement
23 assessment tools and supporting equitable infrastructure management prac-
24 tices globally.

25 7. Conclusion

26 This study presented *Crack-Segmenter*, a fully self-supervised segmenta-
27 tion framework developed specifically for pavement crack detection. Our ap-
28 proach successfully eliminates the dependence on costly and labor-intensive
29 pixel-level annotations, addressing a significant limitation of existing seg-
30 mentation methods. The proposed model integrates three innovative mod-
31 ules: the Scale-Adaptive Embedder (SAE), Directional Attention Trans-
32 former (DAT), and Attention-Guided Fusion (AGF). Each module targets
33 specific challenges in crack segmentation, collectively enhancing the model’s
34 capability to accurately detect diverse crack patterns.

35 Experimental evaluations conducted on ten publicly available datasets
36 demonstrated the effectiveness of the proposed framework. *Crack-Segmenter*

significantly outperformed state-of-the-art fully supervised methods across multiple metrics, including mIoU, Dice score, XOR, and Hammoud Distance (HM). Comprehensive statistical analysis confirmed the statistical significance of these improvements, with notably strong performance gains observed over prominent baseline methods.

Ablation studies further underscored the importance of each proposed module. Specifically, SAE facilitated effective multi-scale feature extraction, capturing cracks across varied widths and complexities. DAT enhanced spatial coherence and continuity, crucial for linear crack structures. Finally, AGF intelligently fused these features, emphasizing contextually relevant information at each scale. Collectively, these modules delivered superior segmentation accuracy, confirming their individual and combined effectiveness.

Attention map visualizations provided interpretability, confirming that *Crack-Segmenter* correctly focused on meaningful crack regions, reducing background distractions. Such explainability strengthens confidence in deploying this method in practical, real-world pavement monitoring tasks.

In future work, extending this framework to handle other pavement distress types and further optimizing computational efficiency could expand its utility. Additionally, exploring methods to integrate temporal information from sequential pavement inspections may improve detection robustness over time, enhancing preventive maintenance practices.

References

- [1] M. Barman, J. Munch, U. M. Arepalli, Cost/benefit analysis of the effectiveness of crack sealing techniques, Research Report MN/RC 2019-26, University of Minnesota, Duluth. Department of Civil Engineering, prepared for the Local Road Research Board (June 2019).
URL <https://rosap.ntl.bts.gov/view/dot/61830>
- [2] B. A. Kyem, E. K. O. Denteh, J. K. Asamoah, A. Aboah, Pavecap: The first multimodal framework for comprehensive pavement condition assessment with dense captioning and pci estimation, ArXiv abs/2408.04110 (2024).
URL <https://api.semanticscholar.org/CorpusID:271768959>
- [3] S. L. H. Lau, E. Chong, X. Yang, X. Wang, Automated pavement crack segmentation using u-net-based convolutional neural network, IEEE Access 8 (2020) 114892–114899. doi:10.1109/ACCESS.2020.3003638.

- 1 [4] B. Agyei Kyem, J. K. Asamoah, A. Aboah, Context-cracknet: A
2 context-aware framework for precise segmentation of tiny cracks in
3 pavement images, *Construction and Building Materials* 484 (2025)
4 141583. doi:<https://doi.org/10.1016/j.conbuildmat.2025.141583>.
5 URL <https://www.sciencedirect.com/science/article/pii/S0950061825017337>
- 6 [5] T. Wen, S. Ding, H. Lang, J. J. Lu, Y. Yuan, Y. Peng,
7 J. Chen, A. Wang, Automated pavement distress segmenta-
8 tion on asphalt surfaces using a deep learning network, *Inter-
9 national Journal of Pavement Engineering* 24 (2) (2023)
10 2027414. arXiv:<https://doi.org/10.1080/10298436.2022.2027414>,
11 doi:10.1080/10298436.2022.2027414.
12 URL <https://doi.org/10.1080/10298436.2022.2027414>
- 13 [6] C. Xiang, V. J. Gan, L. Deng, J. Guo, S. Xu, Unified weakly and
14 semi-supervised crack segmentation framework using limited coarse
15 labels, *Engineering Applications of Artificial Intelligence* 133 (2024)
16 108497. doi:<https://doi.org/10.1016/j.engappai.2024.108497>.
17 URL <https://www.sciencedirect.com/science/article/pii/S0952197624006559>
- 18 [7] J. Li, C. Yuan, X. Wang, G. Chen, G. Ma, Semi-supervised
19 crack detection using segment anything model and deep trans-
20 fer learning, *Automation in Construction* 170 (2025) 105899.
21 doi:<https://doi.org/10.1016/j.autcon.2024.105899>.
22 URL <https://www.sciencedirect.com/science/article/pii/S0926580524006356>
- 23 [8] A. A. Neema Jakisa Owor, Yaw Adu-Gyamfi, M. Amo-
24 Boateng, Pavesam – segment anything for pavement dis-
25 tress, *Road Materials and Pavement Design* 0 (0) (2024)
26 1–25. arXiv:<https://doi.org/10.1080/14680629.2024.2374863>,
27 doi:10.1080/14680629.2024.2374863.
28 URL <https://doi.org/10.1080/14680629.2024.2374863>
- 29 [9] G. Li, J. Wan, S. He, Q. Liu, B. Ma, Semi-supervised semantic segmen-
30 tation using adversarial learning for pavement crack detection, *IEEE
31 Access* 8 (2020) 51446–51459. doi:10.1109/ACCESS.2020.2980086.
- 32 [10] S. Shim, J. Kim, G.-C. Cho, S.-W. Lee, Multiscale and adversar-
33 ial learning-based semi-supervised semantic segmentation approach for

- 1 crack detection in concrete structures, *IEEE Access* 8 (2020) 170939–
2 170950. doi:10.1109/ACCESS.2020.3022786.
- 3 [11] B. A. Kyem, J. K. Asamoah, Y. Huang, A. Aboah, Weather-
4 adaptive synthetic data generation for enhanced power line in-
5 spection using stargan, *IEEE Access* 12 (2024) 193882–193901.
6 doi:10.1109/ACCESS.2024.3520120.
- 7 [12] W. Wang, C. Su, Semi-supervised semantic segmentation network
8 for surface crack detection, *Automation in Construction* (2021).
9 doi:10.1016/J.AUTCON.2021.103786.
- 10 [13] M. Mohammed, Z. Han, Y. Li, Z. Al-Huda, W.-D. Wang, En-
11 hanced pavement crack segmentation with minimal labeled data:
12 a triplet attention teacher-student framework, *International Jour-
13 nal of Pavement Engineering* 2024, VOL. 25 (2024) 2400562.
14 doi:10.1080/10298436.2024.2400562.
- 15 [14] Z. Jian, J. Liu, Cross teacher pseudo supervision: Enhancing semi-
16 supervised crack segmentation with consistency learning, *Adv. Eng. In-
17 form.* 59 (C) (Jan. 2024). doi:10.1016/j.aei.2023.102279.
18 URL <https://doi.org/10.1016/j.aei.2023.102279>
- 19 [15] H. Feng, W. Li, Z. Luo, Y. Chen, S. Fatholahi, M. Cheng, C. Wang, J. M.
20 Junior, J. Li, Gcn-based pavement crack detection using mobile lidar
21 point clouds, *IEEE Transactions on Intelligent Transportation Systems*
22 23 (2021) 11052–11061. doi:10.1109/tits.2021.3099023.
- 23 [16] H. Feng, L. Ma, Y. Yu, Y. Chen, J. Li, Scl-gcn: Stratified con-
24 trastive learning graph convolution network for pavement crack
25 detection from mobile lidar point clouds, *International Journal of
26 Applied Earth Observation and Geoinformation* 118 (2023) 103248.
27 doi:<https://doi.org/10.1016/j.jag.2023.103248>.
28 URL <https://www.sciencedirect.com/science/article/pii/S1569843223000705>
- 29 [17] M. Sohaib, M. J. Hasan, M. Shah, Z. Zheng, A robust self-supervised ap-
30 proach for fine-grained crack detection in concrete structures, *Scientific
31 Reports* 14 (06 2024). doi:10.1038/s41598-024-63575-x.

- 1 [18] K. Zhang, Y. Zhang, H. Cheng, Self-supervised structure learning for
2 crack detection based on cycle-consistent generative adversarial net-
3 works, *Journal of Computing in Civil Engineering* 34 (2020) 04020004.
4 doi:10.1061/(ASCE)CP.1943-5487.0000883.
- 5 [19] S. Karimijafarbigloo, R. Azad, A. Kazerouni, D. Merhof, Ms-former:
6 Multi-scale self-guided transformer for medical image segmentation,
7 in: I. Oguz, J. Noble, X. Li, M. Styner, C. Baumgartner, M. Rusu,
8 T. Heinmann, D. Kontos, B. Landman, B. Dawant (Eds.), *Medical*
9 *Imaging with Deep Learning*, Vol. 227 of *Proceedings of Machine*
10 *Learning Research*, PMLR, 2024, pp. 680–694.
11 URL <https://proceedings.mlr.press/v227/karimijafarbigloo24a.html>
- 12 [20] C. V. Dung, L. D. Anh, Autonomous concrete crack detection using
13 deep fully convolutional neural network, *Automation in Construction*
14 99 (2019) 52–58. doi:<https://doi.org/10.1016/j.autcon.2018.11.028>.
15 URL <https://www.sciencedirect.com/science/article/pii/S0926580518306745>
- 16 [21] F. Yang, L. Zhang, S. Yu, D. Prokhorov, X. Mei, H. Ling, Feature pyra-
17 mid and hierarchical boosting network for pavement crack detection,
18 *IEEE Transactions on Intelligent Transportation Systems* 21 (4) (2020)
19 1525–1535. doi:10.1109/TITS.2019.2910595.
- 20 [22] Y. Liu, J. Yao, X. Lu, R. Xie, L. Li, Deepcrack: A deep hierarchical
21 feature learning architecture for crack segmentation, *Neurocomputing*
22 338 (2019) 139–153. doi:<https://doi.org/10.1016/j.neucom.2019.01.036>.
23 URL <https://www.sciencedirect.com/science/article/pii/S0925231219300566>
- 24 [23] M. D. Jenkins, T. A. Carr, M. I. Iglesias, T. Buggy, G. Morison, A
25 deep convolutional neural network for semantic pixel-wise segmentation
26 of road and pavement surface cracks, in: *2018 26th European signal*
27 *processing conference (EUSIPCO)*, IEEE, 2018, pp. 2120–2124.
- 28 [24] Y. Pan, G. Zhang, L. Zhang, A spatial-channel hierarchi-
29 cal deep learning network for pixel-level automated crack
30 detection, *Automation in Construction* 119 (2020) 103357.
31 doi:<https://doi.org/10.1016/j.autcon.2020.103357>.
32 URL <https://www.sciencedirect.com/science/article/pii/S0926580520309377>

- 1 [25] F. Guo, Y. Qian, J. Liu, H. Yu, Pavement crack detection based on
2 transformer network, *Automation in Construction* 145 (2023) 104646.
3 doi:<https://doi.org/10.1016/j.autcon.2022.104646>.
4 URL <https://www.sciencedirect.com/science/article/pii/S0926580522005167>
- 5 [26] H. Liu, J. Yang, X. Miao, C. Mertz, H. Kong, Crackformer
6 network for pavement crack segmentation, *IEEE Transactions*
7 *on Intelligent Transportation Systems* 24 (9) (2023) 9240–9252.
8 doi:10.1109/TITS.2023.3266776.
- 9 [27] B. A. Kyem, E. K. O. Denteh, J. K. Asamoah, K. A. Tutu, A. Aboah,
10 Advancing pavement distress detection in developing countries: A
11 novel deep learning approach with locally-collected datasets, *ArXiv*
12 *abs/2408.05649* (2024).
13 URL <https://api.semanticscholar.org/CorpusID:271854773>
- 14 [28] Z. Al-Huda, B. Peng, R. N. A. Algburi, S. Alfasly, T. Li, Weakly super-
15 vised pavement crack semantic segmentation based on multi-scale object
16 localization and incremental annotation refinement, *Applied Intelligence*
17 53 (11) (2023) 14527–14546. doi:10.1007/s10489-022-04212-w.
18 URL <https://doi.org/10.1007/s10489-022-04212-w>
- 19 [29] T. He, H. Li, Z. Qian, C. Niu, R. Huang, Research on
20 weakly supervised pavement crack segmentation based on de-
21 fect location by generative adversarial network and target re-
22 optimization, *Construction and Building Materials* 411 (2024) 134668.
23 doi:<https://doi.org/10.1016/j.conbuildmat.2023.134668>.
24 URL <https://www.sciencedirect.com/science/article/pii/S0950061823043891>
- 25 [30] S. Shim, J. Kim, G.-C. Cho, S.-W. Lee, Multiscale and adversar-
26 ial learning-based semi-supervised semantic segmentation approach for
27 crack detection in concrete structures, *IEEE Access* 8 (2020) 170939–
28 170950. doi:10.1109/ACCESS.2020.3022786.
- 29 [31] T. Shi, Y. Wang, Y. Fang, Y. Zhang, Semi-supervised segmentation
30 model for crack detection based on mutual consistency constraint
31 and boundary loss, *Engineering Applications of Artificial Intelligence*
32 139 (2025) 109683. doi:<https://doi.org/10.1016/j.engappai.2024.109683>.
33 URL <https://www.sciencedirect.com/science/article/pii/S0952197624018414>

- 1 [32] C. Han, H. Yang, T. Ma, S. Wang, C. Zhao, Y. Yang, Crackdiffusion: A
2 two-stage semantic segmentation framework for pavement crack combin-
3 ing unsupervised and supervised processes, *Automation in Construction*
4 160 (2024) 105332. doi:<https://doi.org/10.1016/j.autcon.2024.105332>.
5 URL <https://www.sciencedirect.com/science/article/pii/S0926580524000682>
- 6 [33] T. Uelwer, J. Robine, S. S. Wagner, M. Höftmann, E. Upschulte,
7 S. Konietzny, M. Behrendt, S. Harmeling, A survey on self-supervised
8 representation learning (2023). arXiv:2308.11455.
9 URL <https://arxiv.org/abs/2308.11455>
- 10 [34] M. Caron, P. Bojanowski, A. Joulin, M. Douze, Deep clustering for
11 unsupervised learning of visual features (2019). arXiv:1807.05520.
12 URL <https://arxiv.org/abs/1807.05520>
- 13 [35] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Gold-
14 stein, F. Bordes, A. Bardes, G. Mialon, Y. Tian, A. Schwarzschild, A. G.
15 Wilson, J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsiavash,
16 Y. LeCun, M. Goldblum, A cookbook of self-supervised learning (2023).
17 arXiv:2304.12210.
18 URL <https://arxiv.org/abs/2304.12210>
- 19 [36] D. A. Beyene, D. Q. Tran, M. B. Maru, T. Kim, S. Park, S. Park,
20 Unsupervised domain adaptation-based crack segmentation using
21 transformer network, *Journal of Building Engineering* 80 (2023) 107889.
22 doi:<https://doi.org/10.1016/j.jobbe.2023.107889>.
23 URL <https://www.sciencedirect.com/science/article/pii/S2352710223020697>
- 24 [37] Q. D. Nguyen, H.-T. Thai, S. D. Nguyen, Self-training method for
25 structural crack detection using image blending-based domain mixing
26 and mutual learning, *Automation in Construction* 170 (2025) 105892.
27 doi:<https://doi.org/10.1016/j.autcon.2024.105892>.
28 URL <https://www.sciencedirect.com/science/article/pii/S0926580524006289>
- 29 [38] Z. Lin, H. Wang, S. Li, Pavement anomaly detection based on trans-
30 former and self-supervised learning, *Automation in Construction* 143
31 (2022) 104544. doi:<https://doi.org/10.1016/j.autcon.2022.104544>.
32 URL <https://www.sciencedirect.com/science/article/pii/S0926580522004150>

- 1 [39] A Self-Supervised Learning Technique for Road Defects Detec-
2 tion Based on Monocular Three-Dimensional Reconstruction,
3 Vol. Volume 3: 21st International Conference on Advanced Ve-
4 hicle Technologies; 16th International Conference on Design
5 Education of International Design Engineering Technical Con-
6 ferences and Computers and Information in Engineering Con-
7 ference. arXiv:[https://asmedigitalcollection.asme.org/IDETC-](https://asmedigitalcollection.asme.org/IDETC-CIE/proceedings-pdf/IDETC-CIE2019/59216/V003T01A021/6453236/v003t01a021-detc2019-98135.pdf)
8 CIE/proceedings-pdf/IDETC-CIE2019/59216/V003T01A021/6453236/v003t01a021-
9 detc2019-98135.pdf, doi:10.1115/DETC2019-98135.
10 URL <https://doi.org/10.1115/DETC2019-98135>
- 11 [40] H.-Y. Yoon, J.-H. Kim, J.-W. Jeong, Classification of the sidewalk condi-
12 tion using self-supervised transfer learning for wheelchair safety driving,
13 Sensors 22 (1) (2022). doi:10.3390/s22010380.
14 URL <https://www.mdpi.com/1424-8220/22/1/380>
- 15 [41] Q. Song, W. Yao, H. Tian, Y. Guo, R. C. Muniyandi, Y. An, Two-stage
16 framework with improved u-net based on self-supervised contrastive
17 learning for pavement crack segmentation, Expert Syst. Appl. 238 (Part
18 F) (2024) 122406.
19 URL <https://doi.org/10.1016/j.eswa.2023.122406>
- 20 [42] N. Ma, R. Fan, L. Xie, Up-cracknet: Unsupervised pixel-wise road crack
21 detection via adversarial image restoration, IEEE Transactions on In-
22 telligent Transportation Systems 25 (2024) 13926–13936.
23 URL <https://api.semanticscholar.org/CorpusID:267311537>
- 24 [43] S. Karimijafarbigloo, R. Azad, A. Kazerouni, D. Merhof, Ms-former:
25 Multi-scale self-guided transformer for medical image segmentation,
26 in: I. Oguz, J. Noble, X. Li, M. Styner, C. Baumgartner, M. Rusu,
27 T. Heinmann, D. Kontos, B. Landman, B. Dawant (Eds.), Medical
28 Imaging with Deep Learning, Vol. 227 of Proceedings of Machine
29 Learning Research, PMLR, 2024, pp. 680–694.
30 URL <https://proceedings.mlr.press/v227/karimijafarbigloo24a.html>
- 31 [44] Y. Shi, L. Cui, Z. Qi, F. Meng, Z. Chen, Automatic road crack detec-
32 tion using random structured forests, IEEE Transactions on Intelligent
33 Transportation Systems 17 (12) (2016) 3434–3445.

- 1 [45] L. Zhang, F. Yang, Y. D. Zhang, Y. J. Zhu, Road crack detection using
2 deep convolutional neural network, in: Image Processing (ICIP), 2016
3 IEEE International Conference on, IEEE, 2016, pp. 3708–3712.
- 4 [46] Q. Zou, Y. Cao, Q. Li, Q. Mao, S. Wang, Cracktree: Automatic crack
5 detection from pavement images, Pattern Recognition Letters 33 (3)
6 (2012) 227–238.
- 7 [47] S. Ham, S. Bae, H. Kim, I. Lee, G.-P. Lee, D. Kim, Training a semantic
8 segmentation model for cracks in the concrete lining of tunnel, Journal
9 of Korean Tunnelling and Underground Space Association (2021) 549–
10 558doi:10.9711/KTAJ.2021.23.6.549.
- 11 [48] Y. Shi, L. Cui, Z. Qi, F. Meng, Z. Chen, Automatic road crack detec-
12 tion using random structured forests, IEEE Transactions on Intelligent
13 Transportation Systems 17 (12) (2016) 3434–3445.
- 14 [49] M. Eisenbach, R. Stricker, D. Seichter, K. Amende, K. Debes, M. Sessel-
15 mann, D. Ebersbach, U. Stoeckert, H.-M. Gross, How to get pavement
16 distress detection ready for deep learning? a systematic approach., in:
17 International Joint Conference on Neural Networks (IJCNN), 2017, pp.
18 2039–2047.
- 19 [50] M. Pak, S. Kim, Crack detection using fully convolutional network in
20 wall-climbing robot, in: J. J. Park, S. J. Fong, Y. Pan, Y. Sung (Eds.),
21 Advances in Computer Science and Ubiquitous Computing, Springer
22 Singapore, Singapore, 2021, pp. 267–272.
- 23 [51] R. Amhaz, S. Chambon, J. Idier, V. Baltazart, Automatic crack detec-
24 tion on two-dimensional pavement images: An algorithm based on min-
25 imal path selection, IEEE Transactions on Intelligent Transportation
26 Systems 17 (10) (2016) 2718–2729. doi:10.1109/TITS.2015.2477675.
- 27 [52] E. Shelhamer, J. Long, T. Darrell, Fully convolutional networks for se-
28 mantic segmentation, 2015 IEEE Conference on Computer Vision and
29 Pattern Recognition (CVPR) (2014) 3431–3440.
30 URL <https://api.semanticscholar.org/CorpusID:1629541>
- 31 [53] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks
32 for biomedical image segmentation, in: N. Navab, J. Hornegger, W. M.

- 1 Wells, A. F. Frangi (Eds.), Medical Image Computing and Computer-
2 Assisted Intervention – MICCAI 2015, Springer International Publish-
3 ing, Cham, 2015, pp. 234–241.
- 4 [54] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, J. Liang, Unet++: Re-
5 designing skip connections to exploit multiscale features in image seg-
6 mentation, IEEE Transactions on Medical Imaging 39 (6) (2020) 1856–
7 1867. doi:10.1109/TMI.2019.2959609.
- 8 [55] H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid scene parsing network,
9 in: 2017 IEEE Conference on Computer Vision and Pattern Recognition
10 (CVPR), 2017, pp. 6230–6239. doi:10.1109/CVPR.2017.660.
- 11 [56] H. Li, P. Xiong, J. An, L. Wang, Pyramid attention network for semantic
12 segmentation, in: Proceedings of the British Machine Vision Conference
13 (BMVC), 2018, p. 285. doi:10.48550/arXiv.1805.10180.
- 14 [57] R. LI, S. Zheng, C. Zhang, C. Duan, J. Su, L. Wang, P. Atkinson, Mul-
15 tiattention network for semantic segmentation of fine-resolution remote
16 sensing images, IEEE Transactions on Geoscience and Remote Sensing
17 PP (2021) 1–13. doi:10.1109/TGRS.2021.3093977.
- 18 [58] A. Chaurasia, E. Culurciello, Linknet: Exploiting encoder repre-
19 sentations for efficient semantic segmentation, in: 2017 IEEE Vi-
20 sual Communications and Image Processing (VCIP), 2017, pp. 1–4.
21 doi:10.1109/VCIP.2017.8305148.
- 22 [59] A. Kirillov, R. Girshick, K. He, P. Dollar, Panoptic feature pyra-
23 mid networks, in: Proceedings of the IEEE/CVF Conference on Com-
24 puter Vision and Pattern Recognition (CVPR), 2019, pp. 6392–6401.
25 doi:10.1109/CVPR.2019.00656.
- 26 [60] L.-C. Chen, G. Papandreou, F. Schroff, H. Adam, Rethinking atrous
27 convolution for semantic image segmentation, ArXiv abs/1706.05587 (06
28 2017). doi:10.48550/arXiv.1706.05587.
29 URL <https://api.semanticscholar.org/CorpusID:22655199>
- 30 [61] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-
31 decoder with atrous separable convolution for semantic image segmenta-
32 tion, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), Com-

- puter Vision – ECCV 2018, Springer International Publishing, Cham, 2018, pp. 833–851.
- [62] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, J. Sun, Unified perceptual parsing for scene understanding, in: European Conference on Computer Vision, 2018, pp. 432–448.
URL <https://api.semanticscholar.org/CorpusID:50781105>
- [63] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Álvarez, P. Luo, Segformer: Simple and efficient design for semantic segmentation with transformers, in: Neural Information Processing Systems, 2021, pp. 12077–12090.
URL <https://api.semanticscholar.org/CorpusID:235254713>