



UNIVERSITY OF LEEDS

This is a repository copy of *GLoU-MiT: Lightweight Global-Local Mamba-Guided U-mix transformer for UAV-based pavement crack segmentation*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/227785/>

Version: Accepted Version

Article:

Shan, J., Huang, Y. orcid.org/0000-0002-1220-6896, Jiang, W. et al. (2 more authors) (2025) GLoU-MiT: Lightweight Global-Local Mamba-Guided U-mix transformer for UAV-based pavement crack segmentation. *Advanced Engineering Informatics*, 65 (Part D). 103384. ISSN 1474-0346

<https://doi.org/10.1016/j.aei.2025.103384>

This is an author produced version of an article published in *Advanced Engineering Informatics*, made available under the terms of the Creative Commons Attribution License (CC-BY), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited.

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 **GLoU-MiT: Lightweight Global-Local Mamba-Guided U-**
2 **Mix Transformer for UAV-based Pavement Crack**
3 **Segmentation**

4 Jinhuan Shan^{a, b}, Yue Huang^c, Wei Jiang^{a, b, *}, Dongdong Yuan^{a, b}, Feiyang Guo^{a, b}

5

6 a. Key Laboratory for Special Area Highway Engineering of Ministry of Education, Chang'an

7 University, Xi'an 710064, China

8 b. School of Highway, Chang'an University, Xi'an 710064, China

9 c. Institute for Transport Studies, University of Leeds, Leeds LS2 9JT, UK

10

11 Email: jhshan@chd.edu.cn (J. Shan), y.huang1@leeds.ac.uk (Y. Huang), jiangwei@chd.edu.cn (W.

12 Jiang), ddy@chd.edu.cn (D. Yuan), 2024021063@chd.edu.cn (F. Guo)

13

14 *Corresponding author: Prof. Wei Jiang, jiangwei@chd.edu.cn

15

16

17

18

19 **Abstract:** The utility of Unmanned Aerial Vehicles (UAVs) for routine pavement
20 distresses inspection has been increasingly recognized due to their efficiency, flexibility,
21 safety, and low-cost automation. However, UAV-acquired high-altitude images present
22 unique challenges for deep learning-based semantic segmentation models, such as
23 minute crack details, blurred boundaries, and high levels of environmental noise. We
24 propose GLoU-MiT, a lightweight segmentation model designed to address the
25 difficulties of UAV-based pavement crack segmentation. Our model integrates a U-
26 shaped Mix Transformer architecture for efficient hierarchical feature extraction, a
27 Global-Local Mamba-Guided Skip Connection for improved feature alignment and
28 computational efficiency, and a Boundary/Semantic Deep Supervision Refinement
29 Module to enhance segmentation precision in complex scenarios. Extensive
30 experiments on UAV-Crack500, CrackSC and Crack500 datasets demonstrate that
31 GLoU-MiT effectively improves segmentation accuracy, particularly in low-contrast
32 and complex background environments, making it a robust solution for UAV-based
33 pavement crack inspection tasks. Furthermore, inference speed and energy
34 consumption evaluations conducted on the Jetson Orin Nano (8GB) show that our
35 model achieves an excellent balance between accuracy, energy efficiency, and speed.
36 The code will be released at: <https://github.com/SHAN-JH/GLoU-MiT>.

37

38 **Keywords:** Pavement crack, Vision mamba, Vison Transformer, Semantic
39 segmentation, Skip connection

40

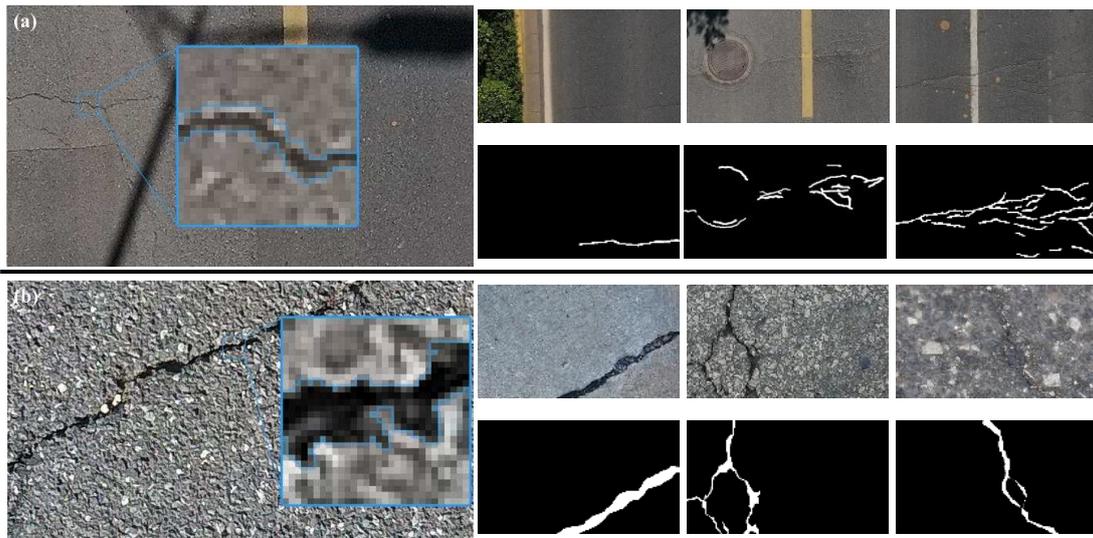
41 **1. Introduction**

42 During the operation of pavement, early micro-cracks often emerge in the
43 pavement structure under the dual influence of vehicular loads and climatic variations,
44 as well as the underlying geological conditions [1–4]. Although these initial cracks may
45 not significantly impact the usability of the pavement, their progression, especially
46 under the combined effects of rainwater penetration and recurrent vehicular pressures,
47 can rapidly evolve into various forms of distresses such as potholes, subsidence, and
48 scouring, severely undermining the overall performance of asphalt pavements.
49 Therefore, early detection and timely intervention of such pavement distresses are
50 crucial, not only reducing the maintenance costs but also effectively extending the
51 lifespan of the pavement [5]. However, the tasks of regular road inspection and accurate
52 diagnosis of pavement distresses demand substantial manpower and financial resources
53 from road maintenance departments.

54 With significant breakthroughs in computer science in recent years, especially the
55 efficiency and efficacy demonstrated by artificial intelligence in handling labor-
56 intensive and repetitive tasks, there are new avenues for achieving automated, high-
57 precision, and cost-effective pavement disease detection [6–10]. This includes, but is
58 not limited to, automated inspections using drones [11–13] or autonomous vehicles
59 [14–16], and distresses identification using advanced algorithms like object detection
60 and semantic segmentation [17–19].

61 Accurate maintenance decision-making relies heavily on high-quality disease data,
62 which in turn depends on efficient automated detection technologies and precision
63 image capturing strategies. For instance, drones, known for their flexibility and high
64 level of automation, can work in conjunction with automated charging stations to
65 facilitate continuous inspection operations [20]. However, there are three main
66 challenges in semantic segmentation of pavement distresses like cracks using drone-
67 captured images: class imbalance, irregularity of edges, and scene noise. The constraint
68 of safe flying altitude results in a lower proportion of cracks in the captured images,

69 thereby limiting the performance of detection algorithms. Additionally, pavement crack
70 areas, characterized by random, sparse, and diverse pixel compositions along with
71 uncommon textures and edges, further increase the difficulty of accurate crack
72 segmentation. Most current detection algorithms and datasets are designed for close-
73 range photographed images (Fig. 1 (b)), hence applying them directly to high-altitude
74 drone images yields suboptimal results. Inappropriate choices of loss functions or
75 network architecture might prevent the model from effectively capturing minute cracks,
76 and in extreme cases, the model might predict all pixels as background, resulting in
77 entirely black images. Moreover, the limitations in flying altitude mean that the
78 captured images contain considerable noise and are heavily influenced by
79 environmental factors such as changes in lighting, which makes it challenging to
80 achieve effective segmentation results with limited training datasets (Fig. 1 (a)).



81 **Fig. 1 (a) UAV-captured image of UAV-Crack500 dataset, and (b) Phone-captured**
82 **image of Crack500 dataset**

83 The introduction of Vision Transformers significantly enhances the model's
84 capability to capture global contextual information, allowing it to transcend the
85 limitations of a local perspective and greatly improve its ability to detect cracks in
86 complex scenarios [21–24]. However, the quadratic computational complexity of
87 Vision Transformers significantly restricts their deployment and application on edge
88 devices. Mamba [25], as an efficient sequence modeling architecture, has attracted

89 considerable attention after being introduced into the visual domain. With its linear
90 computational complexity and robust long-range dependency modeling, Mamba has
91 demonstrated outstanding performance in various visual tasks. Despite these
92 advantages, Vision Mamba's application in crack semantic segmentation remains
93 limited. Recent studies have explored different Mamba-based architectures for crack
94 detection, aiming to balance efficiency and segmentation accuracy. ULNet [26]
95 introduced a cross-visual Mamba feature extraction module and a frequency-domain
96 feature extraction branch, enhancing fine crack detection while maintaining low
97 computational costs. CrackMamba [27] leveraged VMambaV2 as the encoder and
98 introduced a Snake Scan module, which reshapes crack feature sequences based on
99 their natural development patterns, improving feature extraction for complex crack
100 structures. MambaCrackNet [28] integrated Vision Mamba with depthwise separable
101 residual convolutions, forming a hybrid CNN-Mamba segmentation network, which
102 significantly improved crack detection accuracy while maintaining robustness to patch
103 size and training sample variations. SCSegamba [29] further optimized Mamba-based
104 segmentation by proposing a Structure-Aware Visual State Space module, which
105 combines Gated Bottleneck Convolution (GBC) and a Structure-Aware Scanning
106 Strategy (SASS) to enhance morphological feature modeling and semantic continuity
107 of cracks, achieving high segmentation accuracy with only 2.8M parameters and
108 demonstrating excellent real-world deployment performance.

109 Despite these advancements, existing studies indicate that Vision Mamba has not
110 yet achieved performance on par with traditional Convolutional Neural Networks
111 (CNNs) and Transformers in crack segmentation tasks [30]. One of the key reasons for
112 this limitation is the unique nature of crack segmentation. Slender cracks in complex
113 environments require both long-range dependency modeling to differentiate them from
114 background textures and local feature extraction to precisely delineate their boundaries.
115 Current Mamba-based models struggle to balance these two aspects effectively, leading
116 to performance gaps compared to CNNs and Transformers. Therefore, further

117 architectural improvements are necessary to fully unlock Mamba's potential for crack
118 segmentation. Another limitation of Mamba-based crack segmentation models is their
119 lack of evaluation on edge devices. This is primarily because Mamba relies on custom
120 CUDA kernels for its scan operations, which are difficult to track and optimize via
121 TorchScript. As a result, pure Mamba models exhibit slower inference speeds on
122 resource-constrained devices, limiting their practicality for real-time crack detection
123 and UAV-based inspections.

124 To address these challenges, we propose integrating Mamba with traditional CNN
125 and Transformer models, leveraging their complementary strengths. CNNs excel at
126 capturing local texture and edge details, while Transformers efficiently model global
127 dependencies. By incorporating Mamba into this hybrid framework, our model can
128 enhance long-range contextual understanding while preserving fine structural details,
129 achieving a better balance between accuracy and computational efficiency in complex
130 crack segmentation scenarios. We first introduce a U-shaped segmentation framework
131 based on Mix Transformer, which leverages the hierarchical structure and efficient self-
132 attention mechanisms of SegFormer for feature extraction, progressively upsampling
133 and using skip connections to restore crack details. Next, a lightweight Global-Local
134 Mamba-Guided Skip Connection based on Vision Mamba is employed to progressively
135 filter out redundant information from the encoder, reducing the dimensionality of
136 feature maps through direct addition operations, thereby lowering computational
137 complexity in the decoder. Finally, a Boundary/Semantic Deep Supervision
138 Refinement Module is integrated to refine crack boundaries and semantic information,
139 enhancing the model's performance on UAV images. We conduct a comparative
140 analysis of the commonly used close-range pavement defect datasets (Crack500 [29]
141 and CrackSC [21]) and our collected UAV-based pavement defect dataset (UAV-
142 Crack500 [30]). This comparison helps to understand the distributional differences
143 between datasets and provides a theoretical foundation for model improvements.

144 Our contributions can be summarized as follows:

- 145 (1) We propose GLoU-MiT, a lightweight and efficient UAV-based pavement
146 crack segmentation model. Built on a U-shaped Mix Transformer framework, it
147 balances local and global feature extraction while reducing computational cost,
148 making it suitable for edge deployment.
- 149 (2) We introduce a Global-Local Mamba-Guided Skip Connection to enhance
150 feature alignment while reducing computational complexity. By progressively
151 filtering redundant encoder details and directly adding feature maps instead of
152 concatenation, this mechanism improves both efficiency and segmentation
153 accuracy.
- 154 (3) To refine crack segmentation, particularly in low-contrast or complex
155 backgrounds, we integrate a Boundary/Semantic Deep Supervision Refinement
156 Module. This module enhances fine-grained crack boundary detection and
157 semantic consistency, leading to improved F_1 -score and Crack IoU, especially
158 for thin and indistinct cracks.
- 159 (4) Extensive experiments on UAV-Crack500, CrackSC, and Crack500 datasets
160 demonstrate moderate improvements in F_1 -score and Crack IoU while
161 maintaining high efficiency. Furthermore, inference speed and energy
162 consumption evaluations on Jetson Orin Nano (8GB) confirm the model's
163 practical deployment feasibility.

164 **2. Related Works**

165 To address the challenges of minute crack details, blurred boundaries, and high
166 levels of environmental noise in crack segmentation tasks, deep learning model design
167 has primarily focused on three key enhancements. First, multi-scale feature extraction
168 and fusion techniques are employed to effectively capture cracks of varying sizes and
169 intricate patterns[31–34]. Second, advanced attention mechanisms are integrated to
170 highlight critical crack regions and suppress irrelevant background noise[35–38]. Third,
171 boundary refinement strategies are developed to improve the precision of segmentation
172 along crack edges, ensuring accurate delineation even in complex scenarios[39,40].

173 These design considerations have become essential for advancing the performance of
174 deep learning models in crack segmentation. This section summarizes these
175 foundational improvements while introducing the design philosophy of the proposed
176 model.

177 **2.1 Multi-scale Feature Extraction and Fusion**

178 In Convolutional Neural Networks (CNNs), feature extraction transforms raw
179 input data into higher-level representations, aiding subsequent higher-level tasks. This
180 extraction is typically achieved through the convolution, pooling, and other basic
181 modules of the backbone network, enabling understanding of higher-level semantics.
182 Efficient backbone models are crucial for effective feature extraction and high-level
183 semantic representation. Downsampling in backbones, though reducing parameter
184 count and boosting robustness, also diminishes feature map dimensions. For semantic
185 segmentation tasks, mapping high-level segmentation back to original image sizes
186 without losing edge and detail information is essential. FCN [41] first combined
187 features from different stages with transposed convolutional upsampling for end-to-end
188 per-pixel semantic segmentation. U-Net [42], with skip connections, conjoins encoder
189 and decoder stage feature maps, enabling the decoder to relearn details lost during
190 encoding. DeepLabv3+ [43] utilized Atrous Spatial Pyramid Pooling to grasp multi-
191 scale contextual information, merging high-level with low-level feature maps for better
192 edge and detail detection. The skip connections in U-Net not only serve to obtain richer
193 feature information during the decoding stage but also help alleviate the vanishing
194 gradient problem during training, accelerate network convergence, and preserve
195 detailed information in images. These functions of skip connections are what enable U-
196 Net to perform exceptionally well in tasks such as medical image segmentation.

197 In this paper, we enhance the skip connections by integrating Local and Global
198 Mamba into the skip connections at different layers, allowing the model to fully
199 leverage the information from the encoder. This approach facilitates the extraction of

200 effective local and global information, thereby improving the accuracy of crack
201 segmentation in complex backgrounds.

202 **2.2 Enhancements in Attention Mechanisms and Long-Range Dependency**

203 In Convolutional Neural Networks (CNNs), a sequence of convolutional and
204 nonlinear layers is employed for feature extraction from a global receptive field.
205 However, this method traditionally treats all input regions uniformly, potentially
206 leading to substantial noise in complex backgrounds and thereby impairing network
207 efficacy. Deep learning has integrated attention mechanisms, inspired by the neural
208 attention processes of the human brain, to address this issue. These mechanisms allow
209 for differential weighting of features, focusing the network on more effective feature
210 representations while inhibiting less discriminative ones. This approach effectively
211 minimizes distractions from background noise and irrelevant areas, consequently
212 augmenting model performance.

213 Attention mechanisms can be divided into categories like channel attention, spatial
214 attention, and self-attention, depending on their area of focus. Squeeze-and-Excitation
215 Networks (SENet) [44] apply global average pooling and fully connected layers to
216 discern channel-wise feature dependencies. The Bottleneck Attention Module (BAM)
217 [45] combines channel and spatial attentions, effectively enhancing feature extraction
218 without augmenting network depth. The Convolutional Block Attention Module
219 (CBAM) [46] integrates and decouples spatial and channel attentions, improving
220 computational efficiency. Originating in natural language processing, self-attention
221 mechanisms have been adeptly transposed to the realm of computer vision. These
222 mechanisms, through queries, keys, and values, assign varying weights based on inter-
223 feature relationships across different positions, thus aiding models in more effectively
224 capturing contextual data within sequences. The Swin Transformer [47] incorporates a
225 window-based self-attention mechanism, blending traditional CNN structures with
226 attention strategies. This model leverages hierarchical attention mechanisms to
227 assimilate both global and local image information, markedly enhancing its

228 performance. Vision Mamba [48–50] is built upon a state space model (SSM)
229 architecture, which has been adapted for vision tasks. This model incorporates a form
230 of linear attention, making it efficient for processing high-resolution images and
231 handling long-range dependencies.

232 In this paper, we introduce Global-Local Mamba-Guided Skip Connection. This
233 approach leverages the long-range dependency capabilities and linear complexity of the
234 Mamba model, thereby enhancing the model's spatial information perception efficiently
235 without significantly increasing computational overhead. This strategy ensures that the
236 model can effectively capture both local details and global context, improving overall
237 performance in handling complex scenarios.

238 **2.3 Refinement of Segmentation Edges**

239 In the field of deep learning, particularly for semantic segmentation tasks,
240 downsampling is crucial for extracting high-level semantic features. However, this
241 process typically results in the loss of edge detail information. High-level semantic
242 features, while effective for classifying categories, suffer from low resolution. In
243 contrast, lower-level features, although higher in resolution and capable of generating
244 sharp, detailed boundaries, also contain significant background noise. Thus, bridging
245 the information flow between high and low-level features is essential for achieving
246 precise edge segmentation. To effectively combine low and high-level information,
247 focusing on edge information through specific convolutional architectures and loss
248 functions is also crucial for improving boundary segmentation precision. Gated-SCNN
249 [51] employs a dual-branch structure to separately process semantic and edge
250 information, incorporating boundary loss to enhance edge definition in the prediction
251 output. The Boundary-Aware Segmentation Network (BASNet) [52] combines a
252 prediction network with a refinement module, using a hybrid loss function to capture
253 predictions at various resolutions and post-refinement losses. Given that edge points
254 often exhibit uncertainty in segmentation predictions (with confidence levels around
255 0.5), PointRend [53] identifies these uncertain points in coarse segmentation maps. It

256 then re-predicts these points using a Multi-Layer Perceptron (MLP) that integrates both
 257 coarse and detailed features, thereby refining the edges. Deformable Convolution [54]
 258 introduces a learnable offset into its receptive field, enhancing the flexibility of
 259 convolution to align with actual boundary shapes more closely, thus improving feature
 260 extraction and boundary prediction capabilities. Additionally, the boundary loss [55]
 261 proposed by Hoel Kervadec et al. focuses on minimizing the interface area between
 262 segmentation boundaries and ground truth, thereby enhancing the network's capability
 263 in contour space prediction.

264 Inspired by deep supervision and refinement techniques, we effectively
 265 amalgamate boundary and semantic details from various refinement layers using a
 266 boundary/semantic fusion head, thereby substantially improving the performance of
 267 model to discern boundary and semantic nuances at different scales. This advancement
 268 not only bolsters crack segmentation accuracy but also systematically mitigates the
 269 interference caused by environmental noise pixels.

270 **3. Proposed Architecture**

271 **3.1 Preliminaries**

272 **State Space Model (SSM):** The state space model is a mathematical framework
 273 used to describe the representation of the current state and predict future states based
 274 on given inputs. Specifically, the model derives a predicted output function $y(t) \in \mathbf{R}$
 275 from the continuous input function $x(t) \in \mathbf{R}$ and the hidden state representation
 276 $h(t) \in \mathbf{R}^N$, as shown in Equation (1).

$$277 \quad \begin{cases} h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) = \mathbf{C}h(t) + \mathbf{D}x(t) \end{cases} \quad (1)$$

278 where $h'(t)$ is the time derivative of the state $h(t)$, indicating how the state evolves
 279 over time; $\mathbf{A} \in \mathbf{R}^{N \times N}$ is the state transition matrix, determines how the hidden state
 280 updates over time.; $\mathbf{B} \in \mathbf{R}^{N \times 1}$ is the input control matrix, defining how the input

281 influences the state; $\mathbf{C} \in \mathbf{R}^{1 \times N}$ is the output matrix, which translates the state to the
 282 output; $\mathbf{D} \in \mathbf{R}$ is the direct transmission matrix, representing the direct influence of the
 283 input on the output; N represents the latent state dimension. The former part of
 284 Equation (1) is referred to as the **state equation**, while the latter part is called the **output**
 285 **equation**. $\mathbf{D}x(t)$ directly influences the output $y(t)$ by bypassing the state variable
 286 $h(t)$, in a manner similar to a shortcut connection. Consequently, SSM further
 287 simplifies Equation (1) and omits \mathbf{D} (or equivalently, sets $\mathbf{D}=\mathbf{0}$), as shown in Equation
 288 (2).

$$289 \quad \begin{cases} h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t) \\ y(t) = \mathbf{C}h(t) \end{cases} \quad (2)$$

290 **Discretization:** In Equation (2), the input is a continuous time-based signal.
 291 However, since real-world data is typically discrete, it is necessary to derive an
 292 equivalent equation in the discrete-time domain, as shown in Equation (3).

$$293 \quad \begin{cases} h_k = \bar{\mathbf{A}}h_{k-1} + \bar{\mathbf{B}}x_k \\ y_k = \bar{\mathbf{C}}h_k \end{cases} \quad (3)$$

294 where h_k is the hidden state, representing the system's state at time step k ; x_k is the
 295 input signal, representing the input provided to the system at time step k ; y_k is the
 296 output signal, representing the output computed from the hidden state h_k ; $\bar{\mathbf{A}}$ is the
 297 discrete state transition matrix; $\bar{\mathbf{B}}$ is the discrete input control matrix; $\bar{\mathbf{C}}$ is the discrete
 298 output matrix.

299 Mamba (S6 model) adopted the Zero-Order Hold (ZOH) discretization method to
 300 convert the continuous-time state equations into discrete form. The ZOH method
 301 assumes that the input remains constant within each sampling interval, effectively
 302 transforming the continuous-time system into a discrete-time system by holding the last
 303 known input value constant until the next sampling point.

304 (1) Discrete state transition matrix $\bar{\mathbf{A}}$

305 Assuming that the input signal remains constant within a time step Δ (ZOH
 306 method), i.e., $x(t) = x_k$ ($t \in [k\Delta, (k+1)\Delta]$). The solution for $h(t)$ of Equation (2) can be
 307 obtained using the matrix exponential method:

$$308 \quad h(t) = e^{\mathbf{A}(t-k\Delta)}h_k + \int_0^{t-k\Delta} e^{\mathbf{A}\tau}\mathbf{B}x_k d\tau \quad (4)$$

309 where τ is integration variable representing continuous time within a single time step.

310 At $t = (k+1)\Delta$, i.e., after discretization to the next step:

$$311 \quad h_{k+1} = e^{\mathbf{A}\Delta}h_k + \left(\int_0^{\Delta} e^{\mathbf{A}\tau} d\tau\right)\mathbf{B}x_k \quad (5)$$

312 Thus, according to Equation (3), the discrete transition matrix is:

$$313 \quad \bar{\mathbf{A}} = e^{\mathbf{A}\Delta} \quad (6)$$

314 (2) Discrete input control matrix $\bar{\mathbf{B}}$

315 According to Equation (3) and (5), the discrete input control matrix is:

$$316 \quad \bar{\mathbf{B}} = \left(\int_0^{\Delta} e^{\mathbf{A}\tau} d\tau\right)\mathbf{B} \quad (7)$$

317 In control theory, this integral has an analytical solution:

$$318 \quad \int_0^{\Delta} e^{\mathbf{A}\tau} d\tau = \mathbf{A}^{-1}(e^{\mathbf{A}\Delta} - \mathbf{I}) \quad (8)$$

319 Thus, the discrete input matrix is:

$$320 \quad \bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{\mathbf{A}\Delta} - \mathbf{I})\mathbf{B} \quad (9)$$

321 where \mathbf{I} represents the identity matrix. Using the first-order Taylor expansion for $e^{\mathbf{A}\Delta}$

322 ($e^{\mathbf{A}\Delta} \approx \mathbf{I} + \mathbf{A}\Delta$), matrix $\bar{\mathbf{B}}$ can be further simplified:

$$323 \quad \bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{\mathbf{A}\Delta} - \mathbf{I})\mathbf{B} \approx \mathbf{A}^{-1}(\mathbf{A}\Delta)\mathbf{B} = \Delta\mathbf{B} \quad (10)$$

324 (3) Discrete output matrix $\bar{\mathbf{C}}$

325 In a state space model, the temporal evolution of the system is primarily governed
 326 by the state equation, while the output equation serves as an instantaneous mapping and
 327 does not influence the state evolution. This equation merely describes how the current
 328 state $h(t)$ is mapped to the output $y(t)$, without involving differentiation with respect

329 to t or any temporal accumulation effects. Therefore, it represents a static
 330 (instantaneous) linear transformation that does not contribute to the system’s temporal
 331 dynamics. As a result, the output matrix remains unchanged after discretization:

$$332 \quad \bar{\mathbf{C}} = \mathbf{C} \quad (11)$$

333 **Initialization (A):** In Equation (2), if matrix \mathbf{A} is initialized with random values
 334 during training, the model may struggle to achieve optimal results. This is because the
 335 next state is not only influenced by the current state but also by prior states. To address
 336 this, the High-order Polynomial Projection Operators (HiPPO) method is introduced,
 337 which produces a hidden state that effectively memorizes its history. This enhances the
 338 model's ability to handle long-range dependencies, allowing it to capture recent tokens
 339 efficiently while attenuating the influence of older tokens. Such a design helps the
 340 model maintain long-term memory while focusing more on recent information. Based
 341 on this, Mamba introduces two simplified initialization methods for both the complex
 342 and real cases [56], aiming to optimize the handling of long-range dependencies in
 343 different scenarios, as shown in Equation (12).

$$344 \quad \mathbf{A}_{\text{init}} = \begin{cases} -\frac{1}{2} + ni, & \text{S4D-Lin} \\ -(n+1), & \text{S4D-Real} \end{cases} \quad (12)$$

345 where n is state dimension index of N , and i is imaginary unit ($i^2 = -1$).

346 **Selection Mechanisms (B, C and Δ):** The SSM and S4 model (Structured SSM)
 347 exhibit limitations in certain key tasks due to their fixed linear time-invariant (LTI)
 348 nature. In these models, the entire historical state is compressed into a single
 349 representation, and the state evolution matrices \mathbf{A} , \mathbf{B} , \mathbf{C} remain static, meaning they
 350 cannot dynamically adjust to different inputs. This inherent limitation prevents the
 351 model from selectively focusing on or ignoring specific inputs, reducing its adaptability
 352 in complex tasks requiring contextual awareness. To overcome this limitation, Mamba
 353 (Selective SSM, S6 model) introduces an input-dependent selection mechanism, where

354 \mathbf{B} , \mathbf{C} , Δ are dynamically adjusted based on the input sequence, while \mathbf{A} remains
 355 fixed. Specifically:

356 (1) Fixed parameters \mathbf{A} :

357 The matrix \mathbf{A} , which governs the evolution of the hidden state, is initialized
 358 (Equation (12)) and remains input-invariant throughout training. This design ensures
 359 stable memory retention and structured state dynamics, allowing the model to
 360 efficiently encode long-range dependencies while maintaining numerical stability.
 361 However, $\bar{\mathbf{A}}$ can still adapt indirectly through changes in Δ (Equation (6)).

362 (2) Input-dependent parameters \mathbf{B} , \mathbf{C} , Δ :

363 In contrast, \mathbf{B} , \mathbf{C} and Δ are all input-dependent and adapt dynamically, as shown
 364 in Equation (13). This mechanism dynamically adjusts the weights over time based on
 365 the input, improving the model's efficiency in handling long-range dependencies.

$$366 \quad \begin{cases} \mathbf{B} = S_{\mathbf{B}}(x) \\ \mathbf{C} = S_{\mathbf{C}}(x) \\ \Delta = \tau_{\Delta}(\mathbf{Parameter} + S_{\Delta}(x)) \end{cases} \quad (13)$$

367 where, $S_{\mathbf{B}}(x) = \mathbf{Linear}_N(x)$ and $S_{\mathbf{C}}(x) = \mathbf{Linear}_N(x)$ are fully connected layers that
 368 project the input to dimension N ; τ_{Δ} represents the Softplus activation function,
 369 ensuring numerical stability; $\mathbf{Parameter}$ refers to a learnable bias term that is
 370 independent of the input sequence x ; $S_{\Delta}(x) = \mathbf{Broadcast}_D(\mathbf{Linear}_1(x))$ is a linear
 371 transformation applied to x , followed by broadcasting to match the required shape D .

372 Visual State Space (VSS) Block: Originally, the Mamba model was primarily
 373 designed for processing sequential data and demonstrated remarkable capabilities in
 374 fields such as natural language processing. To extend its application to the visual
 375 domain, researchers have proposed extended models, such as Vim [57], VMamba [48],
 376 and LocalMamba [49]. These models transform two-dimensional image data into one-
 377 dimensional sequences along various directions, thereby enabling the S6 architecture
 378 to be directly applied to visual information processing. Specifically, the Vim model
 379 employs a bidirectional scanning mechanism to simultaneously extract contextual

380 information with both forward and backward directions; VMamba introduces a four-
381 directional cross-scanning strategy to comprehensively capture global features from the
382 top, bottom, left, and right; and LocalMamba is designed with a localized scanning
383 strategy, focusing on capturing fine-grained image details.

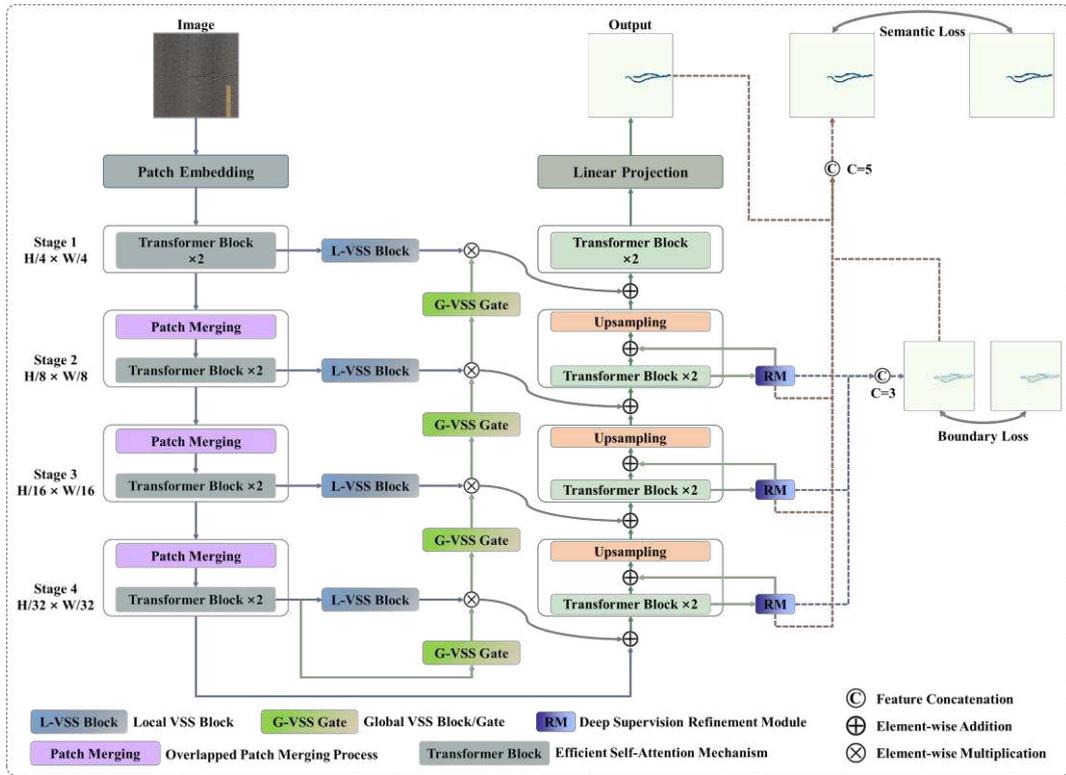
384 **3.2 Comprehensive Architecture**

385 The following section introduces the lightweight pavement crack segmentation
386 model **GLoU-MiT** we designed, as shown in Fig. 2. The model consists of three main
387 components: (1) a U-shaped segmentation framework based on Mix Transformer for
388 hierarchical feature extraction; (2) a lightweight Global-Local Mamba-Guided Skip
389 Connection based on Vision Mamba for enhanced feature alignment and computational
390 efficiency; and (3) a Boundary/Semantic Deep Supervision Refinement Module to
391 improve segmentation precision.

392 The U-shaped segmentation framework builds upon SegFormer’s Mix
393 Transformer, leveraging its hierarchical structure and efficient self-attention
394 mechanisms to extract multi-scale features. Inspired by U-Net, the model progressively
395 upsamples feature maps in the decoder and incorporates skip connections to restore
396 crack details while maintaining computational efficiency. The detailed implementation
397 of this module is described in Section 3.3.

398 Instead of direct concatenation, we introduce a Global-Local Mamba-Guided Skip
399 Connection to enhance feature alignment and computational efficiency. This
400 mechanism employs Local Mamba operations to refine fine-grained details and Global
401 Mamba operations to capture long-range dependencies. Additionally, a Cascaded
402 Gating Mechanism selectively filters redundant information while preserving key
403 semantic features, ensuring an efficient fusion of encoder and decoder features. To
404 further reduce computational cost, the skip connection applies channel reduction before
405 feature processing and expands channels back after computation. The detailed
406 implementation of this module is described in Section 3.4.

407 To enhance segmentation precision, particularly for narrow and low-contrast
 408 cracks, we incorporate a Boundary/Semantic Deep Supervision Refinement Module
 409 into the decoder. This module improves boundary detection and semantic consistency
 410 by leveraging multi-scale deep supervision, deformable convolutions, and attention
 411 mechanisms. By integrating both boundary-aware and semantic feature learning, the
 412 module enhances fine-grained segmentation. The detailed implementation of this
 413 module is described in Section 3.5.



414
415 **Fig. 2 Comprehensive Architecture**

416 **3.3 U-Mix Transformer**

417 SegFormer introduces the Mix Transformer encoder, which effectively leverages
 418 Overlapped Patch Merging Process and Efficient Self-Attention Mechanism,
 419 combining the strengths of both CNNs and Transformers to capture local and global
 420 information. This significantly improves segmentation accuracy. The Efficient Self-
 421 Attention mechanism reduces the dimensionality of the Key and Query using
 422 convolution, which enhances computational efficiency. However, SegFormer's decoder
 423 directly concatenates feature maps from different hierarchical levels and decodes them

424 through an All-MLP architecture. Although this approach reduces computational
 425 complexity, the results from concatenating feature maps of different levels tend to be
 426 suboptimal. This is because feature maps from different layers in the decoder inherently
 427 contain varied information, and directly processing them through MLP cannot fully
 428 exploit the features extracted at different levels.

429 Inspired by U-Net, we propose a symmetric **U-Mix Transformer** model (Fig. 3).
 430 Unlike U-Net, to improve computational efficiency, we do not concatenate feature
 431 maps at the same resolution. Instead, we directly add them together, followed by a Mix
 432 Transformer operation. This approach enhances efficiency while maintaining the ability
 433 to capture multi-level features effectively.

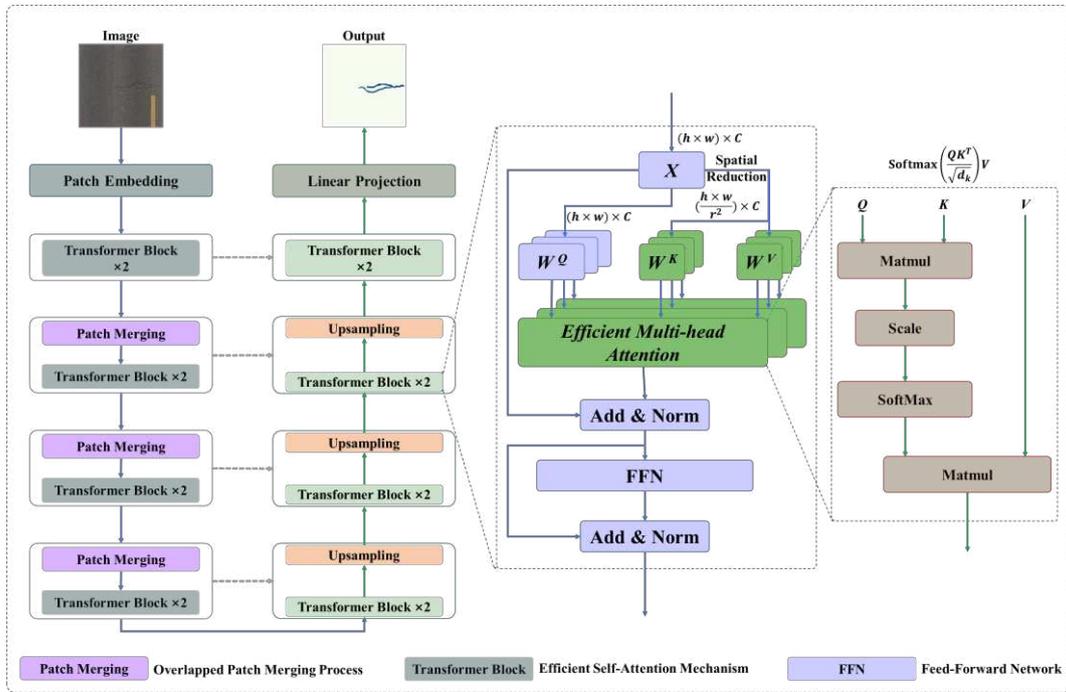
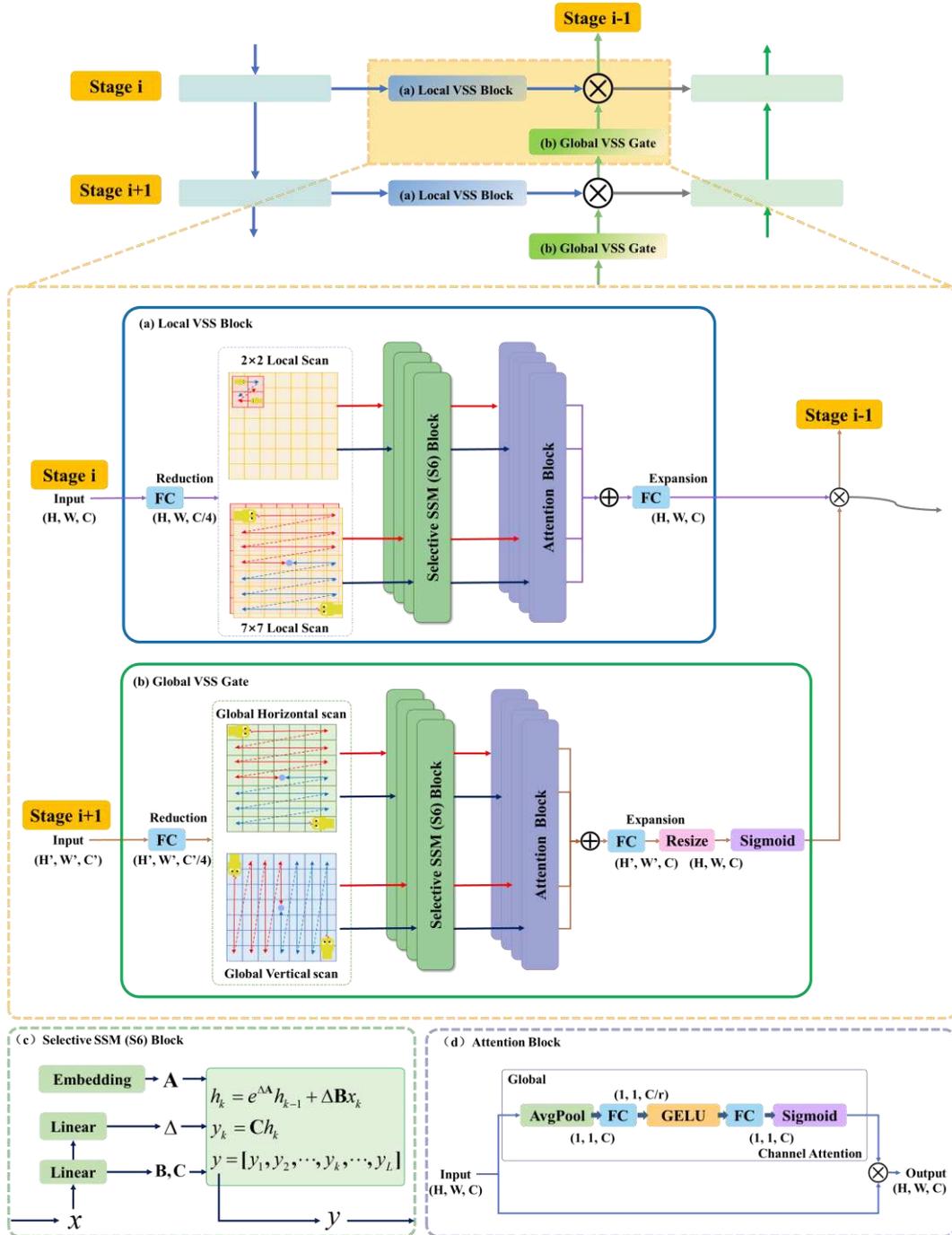


Fig. 3 U-Mix Transformer

436 3.4 Global-Local Mamba-Guided Skip Connection Module

437 Although directly adding feature maps of the same resolution but different
 438 levels—i.e., semantic-rich and detail-heavy maps—can reduce computational load, the
 439 differences in information, such as varying semantic and detail levels, may lead to
 440 suboptimal segmentation results. This is because redundant information from shallow
 441 layers can interfere with the deep semantic information. Like U-Net, the skip

442 connection concatenates the encoder's feature maps with those of the decoder, but this
 443 increases the feature map dimensionality, leading to a quadratic increase in computation
 444 for Transformer modules. To address this issue, we propose the **Lightweight Global-**
 445 **Local Mamba-Guided Skip Connection Module** (Fig. 4).



446
 447 **Fig. 4 Lightweight Global-Local Mamba-Guided Skip Connection Module**

448 First, each feature map from different layers of the encoder is processed through
 449 local Mamba operations to extract detailed information from different levels.

450 Specifically, we introduced local horizontal scans with window sizes of 2 and 7 in Stage
451 1 and Stage 2 to construct Local VSS (L-VSS) Blocks (Fig. 4(a)). Since the window
452 size of 2 inherently includes vertical scanning, we avoided adding additional vertical
453 scans for the size 7 window, reducing computational cost while still capturing vertical
454 dependencies.

455 To better fuse detailed information with deep semantic information, we propose a
456 **Mamba-based Cascaded Gate Generator**. This module starts from the deep feature
457 maps and progressively applies Global VSS (G-VSS) Blocks to extract global
458 information, incorporating scans in four directions: horizontal (left to right and right to
459 left) and vertical (top to bottom and bottom to top). A Sigmoid Gate is employed as a
460 gating mechanism to control the flow of information at each layer. By dynamically
461 adjusting the propagation of global features through the Sigmoid function, the model
462 multiplies these with the extracted local features, ensuring that important global
463 information is selectively retained while noise and irrelevant parts are suppressed.

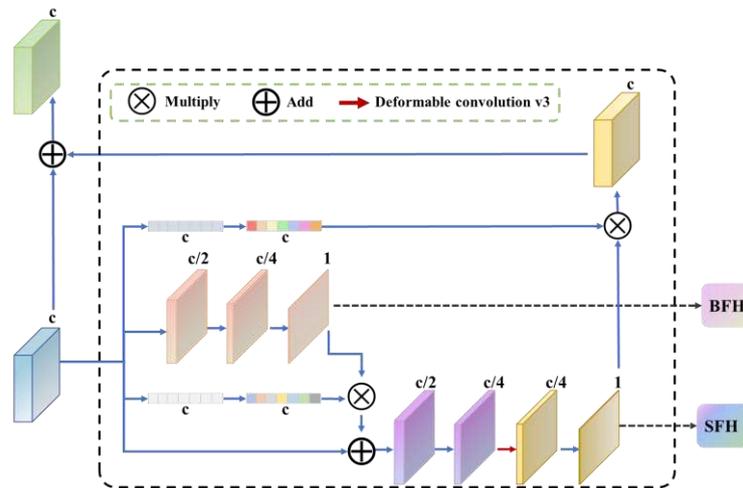
464 Through this gating mechanism, the model effectively suppresses redundant
465 information while maintaining global feature propagation. By cascading this process,
466 the feature maps obtained from the encoder are better aligned with the information
467 hierarchy in the decoder. As a result, the addition of these maps not only reduces
468 computational load but also enhances the model's ability to capture fine details.

469 To manage computational costs, VSS Blocks are not directly inserted into the skip
470 connections. Instead, we first apply channel reduction, decreasing the channels to 1/4
471 of their original number. The globally or locally scanned data is then passed into a
472 Selective Scan State Space Models (S6) [25] for computation (Fig. 4(c)), facilitating
473 effective global or local visual representation learning. Finally, the channels are
474 expanded back to their original size. To further enhance feature representation and
475 eliminate irrelevant information, we implement channel attention mechanism
476 operations in LocalMamba [49]. This allows for weighted extraction of critical feature

477 channels by dynamically adjusting their importance based on the input, ensuring that
 478 the model focuses on the most informative features while suppressing less relevant ones.

479 3.5 Deep Supervision Refinement Module

480 In our preliminary experiments, we observed that although the aforementioned
 481 structure improved segmentation accuracy for wider and more distinct cracks, its
 482 performance remained suboptimal when dealing with narrow cracks that resemble
 483 background pixels or are located in complex backgrounds. To address this issue, we
 484 propose the **Deep Supervision Refinement Module** (Fig. 5), which can be integrated
 485 into any network that requires semantic and boundary supervision. In this study, the
 486 module is inserted before the upsampling operation at each layer.



487
 488 **Fig. 5 Deep supervision refinement module**

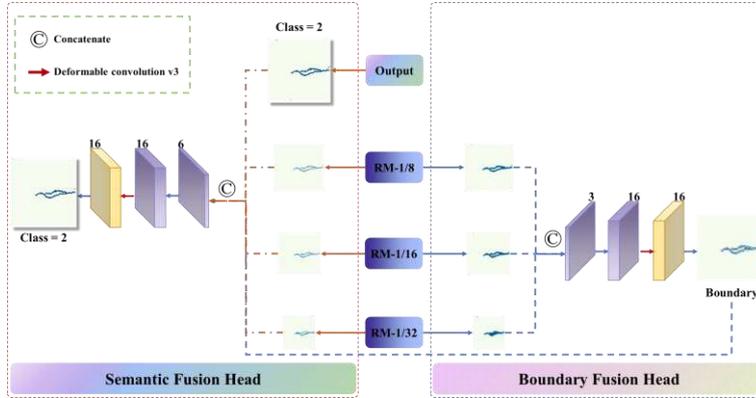
489 This module generates three outputs: a single-channel boundary feature map, a
 490 single-channel semantic feature map, and an upsampled feature map with the same
 491 number of channels as the input. First, the input feature map undergoes three
 492 convolutions to produce a single-channel boundary feature map, which is then sent
 493 directly to the boundary fusion head for further boundary refinement and boundary loss
 494 supervision. Additionally, channel boundary attention is computed on the input feature
 495 map to determine the boundary attention weights that need to be applied to each channel,
 496 and this result is added to the input feature map.

497 Next, the boundary-enhanced feature map is sent to the semantic calculation
498 module, where two standard convolutions are applied to extract semantic features and
499 reduce the channel dimensions. This is followed by a deformable convolution
500 (Deformable Convolution v3) [42] to capture the non-uniform shapes of crack
501 semantics. The output is a single-channel semantic feature map, which is passed to the
502 semantic fusion head for further crack semantic feature extraction and semantic loss
503 computation. Subsequently, semantic attention is calculated for the input feature map,
504 weighted using the semantic feature map, and then added to the input feature map.
505 Finally, the feature map is upsampled to seamlessly connect with the next layer's
506 module.

507 The main idea behind this module is that feature maps after convolution and
508 downsampling contain rich semantic information, but because their dimensions are
509 smaller than the original image, directly upsampling them may result in inaccurate
510 boundary information. By performing boundary and semantic supervision based on
511 high-level semantic information, this module refines the boundary and semantic details
512 at each layer. Additionally, a shortcut branch adds the supervised output back to the
513 original feature map, preserving the original information while enhancing boundary and
514 semantic channel supervision. This module can be integrated into any intermediate
515 feature map calculation to supervise both boundary and semantic information. To
516 capture boundary and semantic information at different scales, it is recommended to
517 apply the module before each upsampling operation.

518 Furthermore, the boundary and semantic feature maps extracted at each layer are
519 fed into the boundary and semantic fusion head for further refinement (Fig. 6). First,
520 the boundary prediction result is upsampled to the original image size and concatenated,
521 followed by a standard convolution layer and a deformable convolution layer, which
522 refine the boundary segmentation information. The boundary-refined result is then
523 concatenated with the semantic prediction result, passed through another standard
524 convolution layer and deformable convolution layer, and the boundary information is

525 used to constrain the semantic information, resulting in the final refined crack
 526 segmentation output.



527

528

Fig. 6 Boundary and Semantic Fusion Head

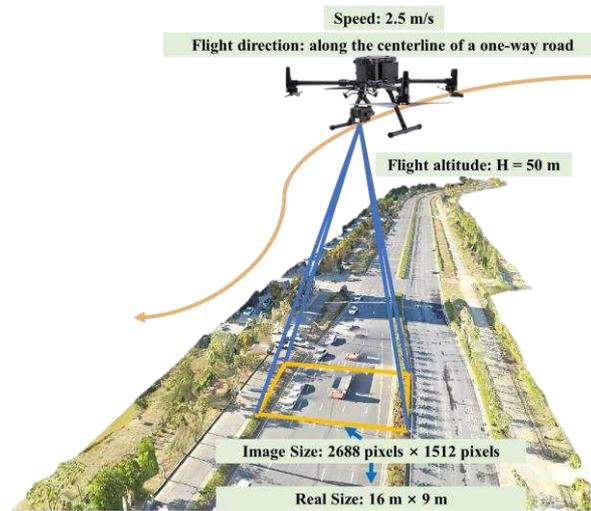
529 4. Experimental Details

530 4.1 Datasets

531 To evaluate the versatility of data acquisition methods for pavement inspection
 532 and to validate the robustness and adaptability of the proposed algorithm, this study
 533 used three distinct datasets: our UAV-Crack500 [58], which is based on long-distance
 534 pavement distresses images captured by drones, the Crack500 [59], which is based on
 535 close-distance taken with cell phones, and the CrackSC [23], which features pavement
 536 cracks captured by cell phones in the presence of complex background noise.

537 For the UAV-Crack500 dataset, given urban scenario flight altitude restrictions,
 538 the drone flying altitude was set to 50 meters. To ensure high precision and efficient
 539 image data collection covering at least a width of three lanes in one direction, a camera
 540 with 4× zoom capabilities was used at a flying speed of 2.5 m/s. The collected images
 541 have a resolution of 2688×1512 pixels, corresponding to a ground coverage area of 16
 542 m × 9 m, which equates to an actual size of 6 mm × 6 mm per pixel (Fig. 7). Due to the
 543 minor proportion that cracks occupy in drone images, directly training on full-
 544 resolution images results in the model being biased toward background predictions. To
 545 alleviate this issue and improve the proportion of crack pixels within each sample, we
 546 adopted a uniform, non-overlapping cropping strategy that divides each image into 16

547 equal-sized blocks of 672×378 pixels. These patches are then filtered to exclude those
548 without visible cracks. A total of 500 image patches with representative and complex
549 crack scenarios were meticulously selected and annotated. The selection process
550 emphasized diversity and real-world complexity, including the presence of road
551 markings, shadows, curbstones, trees, manhole covers, and road dividers (Fig. 1(a)).
552 The dataset was randomly split into a training set (250 images), a validation set (50
553 images), and a test set (200 images). As depicted, the dataset presents elongated, fine-
554 grained crack structures, where crack widths are often only a few pixels wide,
555 embedded within complex urban environments.



556

557

Fig. 7 Flight parameters for long-distance pavement monitoring

558 The Crack500 dataset comprises 500 high-resolution images of pavement captured
559 at close range using a cell phone camera, each with dimensions of 2000×1500 pixels.
560 For practicality in model training and evaluation, these images have been subdivided
561 into 16 non-overlapping regions. The dataset is stratified into different subsets for the
562 purposes of model development and performance assessment: 1896 regions are
563 allocated for training, 348 for validation, and 1124 for testing. Visual inspection of the
564 dataset reveals that the images feature relatively wide cracks, which occupy a more
565 significant proportion of the image area compared to those in UAV-based datasets.
566 Additionally, the images in the Crack500 dataset exhibit minimal variation in lighting
567 conditions and are less affected by environmental noise (Fig. 1(b)).

568 The CrackSC dataset is a newly introduced pavement crack image dataset
569 designed to address the challenges of detecting cracks in local roads with heavy
570 shadows and dense crack formations, commonly found in low-maintenance areas. The
571 dataset consists of 197 images of pavement surfaces collected using an iPhone 8 around
572 Enoree Ave, Columbia, SC. Since the dataset is provided as a whole, we randomly split
573 it into training, validation, and test sets in a 5:1:4 ratio to facilitate model evaluation.

574 4.2 Implementation Details

575 Our models were developed based on the MMSegmentation v1.2.0 framework
576 [60], which provides a standardized and modular implementation for semantic
577 segmentation. To ensure a fair and reliable comparison, all models, including our
578 proposed model and the comparison models, were trained and evaluated under the same
579 hardware and experimental settings. Specifically, all models were trained and inferred
580 on an NVIDIA Tesla T4 GPU (16GB), with only the backbone and head components
581 modified, while all other experimental settings remained identical. The backbone
582 parameters were initialized using pre-trained weights from the official repository,
583 ensuring that each model benefited from a strong initialization aligned with its
584 architecture.

585 During training, we adopted a batch size of 16 and a learning rate of $6e-5$, training
586 each model for 30,000 iterations. To improve the generalization capability of the models,
587 we applied data augmentation techniques, including flipping, rotation, color jittering,
588 and random size cropping. The cropped images were then resized to 256×256 pixels
589 before being used for training. Additionally, all models shared the same data
590 preprocessing pipeline, training strategy, optimizer settings, loss function, learning rate
591 schedule, and evaluation metrics, ensuring that the observed performance differences
592 were solely attributed to the architectural variations rather than training discrepancies.
593 To facilitate efficient model training, the study adopts AdamW ($\beta_1 = 0.9$, $\beta_2 = 0.999$,
594 *weight decay* = 0.01) as the optimizer of choice due to its effectiveness in handling
595 sparse gradients and its adaptive learning rate capabilities, which are conducive for

596 faster convergence. Additionally, a two-stage learning rate scheduling strategy was
 597 implemented. In the initial 1,500 iterations, a linear warm-up was applied, gradually
 598 increasing the learning rate from 1e-6 to the base learning rate 6e-5. Subsequently, a
 599 polynomial learning rate decay (power = 1.0) was adopted, ensuring a linear reduction
 600 in the learning rate from iteration 1,500 to 30,000, eventually reaching 0.0 at the end of
 601 training.

602 4.3 Loss Functions

603 The prevalent class imbalance, marked by a minor fraction of crack pixels
 604 compared to the overall image, severely hinders the capacity of models to discern crack
 605 features using a conventional Binary Cross-Entropy (BCE) loss function. This
 606 challenge is further pronounced in images obtained via UAVs, where the proportion of
 607 crack pixels is notably diminutive, leading to predictions with excessively fine or
 608 interrupted cracks, or in extreme cases, a complete bias towards background
 609 classification. To counteract this, Weighted Binary Cross-Entropy (w BCE) and Dice
 610 Loss have been advocated as effective loss functions to tackle the class imbalance issue.

611 The w BCE approach involves differential weighting for the positive (crack) and
 612 negative (background) classes, incentivizing the model to focus more on the sparsely
 613 represented crack class (Equation (14)). Although w BCE can mitigate the class
 614 imbalance problem, its effectiveness heavily depends on the weight parameter
 615 adjustment, often requiring extensive experimentation to determine the optimal setting
 616 for ensuring model stability and generalization.

$$617 \quad \ell_{wBCE} = -\frac{1}{N} \sum_{i=0}^N (w_1 \square y_i \square \log(y_i) + w_0 \square (1 - y_i) \square \log(1 - y_i)) \quad (14)$$

618 where N represents the number of image pixels; w_1 represents the weight for the
 619 positive samples; w_0 represents the weight for the negative samples; y_i represents the
 620 actual probability of the positive samples; y_i represents the predicted probability of the

621 positive samples. When $w_0 = w_1 = 1$, ℓ_{wBCE} reduces to the standard Binary Cross-
 622 Entropy (BCE) loss ℓ_{BCE} .

623 The Dice Loss effectively measures the degree of overlap between the predicted
 624 segmentation mask and the ground truth. It is inherently robust to class imbalance
 625 (Equation (15)). However, it is prone to gradient instability during training, especially
 626 when the intersection of the segmentation masks is minimal, potentially leading to
 627 oscillatory behaviors or convergence issues in the learning process.

$$628 \quad \ell_{Dice} = 1 - \frac{2 \sum_{i=0}^N y_i y_i}{\sum_{i=0}^N y_i^2 + \sum_{i=0}^N y_i^2} \quad (15)$$

629 To address the respective shortcomings of traditional BCE, $wBCE$, and Dice Loss,
 630 this study follows the approach in [61–65] and adopts a combined BCE and Dice loss
 631 strategy as the overall semantic loss function ℓ_s (Equation (16)). This method has been
 632 demonstrated to achieve better segmentation performance, particularly for fine crack
 633 detection.

$$634 \quad \ell_s = \ell_{BCE} + \ell_{Dice} \quad (16)$$

635 Inspired by the CE2P model [66], we propose an approach to address boundary
 636 loss by first detecting the boundaries between different semantic regions through
 637 comparing the semantic categories of adjacent pixels in the segmentation map. These
 638 boundaries are then marked as edges. To further refine the edge information, a dilation
 639 operation is applied, which widens the edges and produces an expanded edge map (with
 640 an edge width of 4). This wider edge representation helps the model capture and learn
 641 complex semantic boundary information more effectively. Given that boundary pixels
 642 constitute a smaller proportion of the total pixels, we assign a weight of 20 ($\alpha=20$) to
 643 the boundary BCE loss ℓ_b , following the empirical setting in the PIDNet model [67],

644 before combining it with the semantic loss to form the final loss function (Equation
 645 (17)):

$$646 \quad \ell = \ell_s + \alpha \ell_b \quad (17)$$

647 **4.4 Evaluation Metrics**

648 This study employs a suite of metrics for model evaluation, comprising pixel
 649 accuracy (PA), precision (Pr), recall (Re), F₁-score (F₁), and Crack Intersection over
 650 Union (IoU). Pixel accuracy quantifies the proportion of pixels correctly classified by
 651 the model relative to the total pixel count, serving as a measure of the model's overall
 652 classification efficacy (Equation (18)). Precision is defined as the ratio of true positive
 653 predictions to the total number of positive predictions made by the model, gauging the
 654 precision with which the model discerns positive cases (Equation (19)). Recall is the
 655 ratio of true positive predictions to the actual number of positive instances, evaluating
 656 the model's proficiency in detecting all positive cases (Equation (20)). The F₁-score, a
 657 weighted harmonic mean of precision and recall, balances the two metrics for a holistic
 658 performance assessment (Equation (21)). The Crack IoU metric assesses segmentation
 659 accuracy by calculating the ratio of the intersection to the union of the predicted and
 660 actual crack regions (Equation (22)).

$$661 \quad \text{PA} = \frac{TP + TN}{TP + TN + FP + FN} \quad (18)$$

$$662 \quad \text{Pr} = \frac{TP}{TP + FP} \quad (19)$$

$$663 \quad \text{Re} = \frac{TP}{TP + FN} \quad (20)$$

$$664 \quad \text{F}_1 = 2 \times \frac{\text{Pr} \times \text{Re}}{\text{Pr} + \text{Re}} \quad (21)$$

$$665 \quad \text{IoU} = \frac{TP}{TP + FP + FN} \quad (22)$$

666 where *TP* (True Positive) refers to the case when the actual class is positive and the
 667 model predicts it as positive; *TN* (True Negative) refers to the case when the actual class
 668 is negative and the model predicts it as negative; *FP* (False Positive) refers to the case
 669 when the actual class is negative but the model incorrectly predicts it as positive; *FN*

670 (False Negative) refers to the case when the actual class is positive but the model
671 incorrectly predicts it as negative.

672 Params (the number of parameters) and FLOPs (floating point operations) are used
673 in this paper as common metrics to evaluate the parameter complexity and
674 computational complexity of the model. Params refers to the total number of trainable
675 parameters in the model, typically including weights and biases, which are used to
676 measure the model's storage requirements and memory consumption during training.
677 FLOPs represent the number of floating-point operations required during a single
678 inference, which reflects the computational resources needed by the model to process
679 data.

680 Due to the unclear boundaries of cracks, along with the subjectivity and lack of
681 repeatability in manual annotations, this study, following other research [68–70], also
682 adopts a 2-pixel tolerance. To enhance the model's ability to accurately segment cracks,
683 we use the approach where a prediction is considered positive if it falls within the 2-
684 pixel dilated region of the ground truth.

685 **4.5 Reference Evaluation on Jetson Orin Nano**

686 This study systematically evaluates the deployment performance of deep learning
687 models on the NVIDIA Jetson Orin Nano (8GB) platform (Fig. 8), which delivers 40
688 TOPS of AI computational power—an 80-fold increase compared to the NVIDIA
689 Jetson Nano. Due to the limited support for ONNX and TensorRT operators in certain
690 model implementations, we adopted two deployment strategies utilizing PyTorch Just-
691 In-Time (JIT) at FP32 precision: Python-based PyTorch inference and C++-based
692 LibTorch inference. The JIT compilation process converts dynamic models from
693 Python environments into the TorchScript format, enabling cross-platform serialization
694 and optimized execution. Each deployment approach offers unique advantages. Python-
695 based PyTorch inference features short development cycles and flexible iterations,
696 making it ideal for rapid prototyping and validation. Conversely, C++-based LibTorch
697 inference is optimized for latency-sensitive applications and resource-constrained

698 environments. Performance evaluation was conducted under a standardized protocol,
 699 which included setting the Jetson platform to maximum performance mode with the
 700 highest operating frequency. Each test consisted of a 10-sample warm-up phase
 701 followed by 50 inference samples. The evaluation metrics included Energy Per Sample
 702 (EPS) (Equation (23)), Inference Latency Per Sample (LPS) (Equation (24)),
 703 Throughput (Equation (25)), and the Energy Delay Product (EDP) (Equation (26)).

$$704 \quad \text{EPC} = \frac{\sum (V(t) \times I(t) \times \Delta t)}{N} \quad (23)$$

$$705 \quad \text{LPS} = \frac{\sum (T_{\text{inference_end}} - T_{\text{inference_start}})}{N} \quad (24)$$

$$706 \quad \text{Throughput} = \frac{N}{T_{\text{total_end}} - T_{\text{total_start}}} \quad (25)$$

$$707 \quad \text{EDP} = \text{EPS} \times \text{LPS} \quad (26)$$

708 where $V(t)$ and $I(t)$ are real-time voltage and current measurements from the Jetson
 709 INA3221 sensor; N is the total number of samples; $T_{\text{inference_end}}$ and $T_{\text{inference_start}}$ represent
 710 the end and start times of the model inference phase only; $T_{\text{total_end}}$ and $T_{\text{total_start}}$ represent
 711 the end and start times of the complete pipeline of data preprocessing, model inference,
 712 and post-processing stages. During the training process, the model employed a crop
 713 size of (256, 256). To ensure consistency, the same crop size was applied during
 714 inference on edge devices. For test images with a resolution of (512, 512), we utilized
 715 a non-overlapping sliding window approach to divide the image into four patches.
 716 These patches were simultaneously input as a single batch for model inference. Finally,
 717 the prediction results were reassembled to restore the original image size, ensuring
 718 consistency and completeness of the output.



Jetson Orin Nano (8G) Specifications

AI Performance	40TOPS
GPU	1024-core NVIDIA Ampere architecture GPU with 32 Tensor Cores
GPU Frequency	625MHz (Max)
CPU	6-core Arm® Cortex®-A78AE v8.2 64-bit CPU, 1.5MB L2 + 4MB L3
CPU Frequency	1.5GHz (Max)
Memory	8GB 128-bit LPDDR5, 68GB/s
Storage	128GB NVMe Solid State Drive
Power	7W ~ 15W

719

720

Fig. 8 NVIDIA Jetson Orin Nano (8GB) platform and specifications

721 5. Results and Discussion

722 5.1 Quantitative Evaluation

723 To rigorously evaluate the effectiveness of our model, we conducted a comparative
 724 analysis against state-of-the-art semantic segmentation models. These included CNN-
 725 based models (U-Net [42], SegNeXt [71], RHACrackNet [72], TV-net (U-NetSmall)
 726 [58], CDS-Net [73]), Transformer-based models (SegFormer [74], U-MixFormer [75]),
 727 as well as advanced mamba-based models (LocalVMamba [49], Manba-UNet [50],
 728 SCSEgamba [29]). Tables 1 to 3 summarize the performance comparison of state-of-
 729 the-art methods on UAV-Crack500, CrackSC, and Crack500 datasets. Table 4 presents
 730 a comparison of model parameters and FLOPs with a fixed input size of $3 \times 512 \times 512$
 731 during testing.

732

733

Table 1 Performance comparison with the state-of-the-art methods on UAV-Crack500

Method	PA	Pr	Re	F₁	IoU
U-Net	99.12	89.52	70.65	78.97	65.25
SegNeXt-T	<u>99.21</u>	<u>91.25</u>	75.11	82.40	70.06
RHACrackNet	99.11	87.96	72.44	79.45	65.91
TV-Net	99.12	87.42	73.23	79.70	66.25
CDS-Net	99.11	87.52	72.83	79.50	65.97
SegFormer-MiT-B0	99.21	90.51	74.57	81.77	69.16
U-MixFormer-MiT-B0	99.22	90.39	75.36	82.20	69.77
SCSEgamba	99.00	91.31	64.31	75.47	60.60
LocalVMamba-T	99.16	89.56	73.09	80.49	67.36
Mamba-UNet	98.92	87.80	61.40	72.26	56.57
GLoU-MiT (Ours)	99.19	89.32	<u>76.49</u>	<u>82.41</u>	<u>70.08</u>
GLoU-MiT-DS (Ours)	99.20	88.82	77.89	83.00	70.94

734

735

Table 2 Performance comparison with the state-of-the-art methods on CrackSC

Method	PA	Pr	Re	F ₁	IoU
U-Net	98.76	86.56	59.85	70.77	54.76
SegNeXt-T	98.63	<u>88.53</u>	51.28	64.94	48.08
RHACrackNet	98.70	87.02	56.78	68.72	52.35
TV-Net	98.76	86.62	60.49	71.24	55.32
CDS-Net	98.78	88.29	59.26	70.92	54.94
SegFormer-MiT-B0	98.81	88.65	61.91	72.90	57.36
U-MixFormer-MiT-B0	98.84	88.48	64.05	<u>74.31</u>	<u>59.12</u>
SCSegamba	98.64	88.29	51.70	65.21	48.38
LocalVMamba-T	98.72	87.44	57.42	69.32	53.04
Mamba-UNet	98.34	82.41	29.04	42.94	27.34
GLoU-MiT (Ours)	98.80	87.64	<u>64.35</u>	74.21	59.00
GLoU-MiT-DS (Ours)	<u>98.82</u>	87.71	64.78	74.52	59.39

736

737

Table 3 Performance comparison with the state-of-the-art methods on Crack500

Method	PA	Pr	Re	F ₁	IoU
U-Net	97.48	83.81	72.57	77.79	63.65
SegNeXt-T	97.43	79.55	<u>78.63</u>	79.08	65.40
RHACrackNet	97.29	79.63	74.93	77.21	62.87
TV-Net	97.42	80.33	76.44	78.34	64.39
CDS-Net	97.47	<u>82.45</u>	74.28	78.15	64.14
SegFormer-MiT-B0	97.50	81.28	77.07	79.12	65.45
U-MixFormer-MiT-B0	97.56	81.21	78.46	79.81	66.40
SCSegamba	97.33	82.04	71.96	76.61	62.09
LocalVMamba-T	97.55	82.15	76.57	79.26	65.64
Mamba-UNet	97.41	81.20	75.02	77.98	63.91
GLoU-MiT (Ours)	<u>97.59</u>	80.79	80.39	80.59	67.49
GLoU-MiT-DS (Ours)	97.62	82.40	77.72	<u>79.99</u>	<u>66.66</u>

738

739

Table 4 Efficiency comparison

Method	Params (M)	Flops (G)
U-Net	29.0	203.0
SegNeXt-T	4.2	6.3
RHACrackNet	1.7	7.3
TV-Net	17.8	87.3
CDS-Net	7.2	48.5
SegFormer-MiT-B0	3.7	7.9
U-MixFormer-MiT-B0	6.4	5.2
SCSegamba	3.1	23.5
LocalVMamba-T	56.2	230.7
Mamba-UNet	19.2	1.5
GLoU-MiT (Ours)	6.6	6.5
GLoU-MiT-DS (Ours)	7.0	7.2

740

UAV-Crack500: This dataset consists of aerial images of pavement captured by

741

drones. Due to the high altitude of capture, the images have relatively low resolution,

742

with fine cracks and low contrast between the pavement and cracks, making it

743

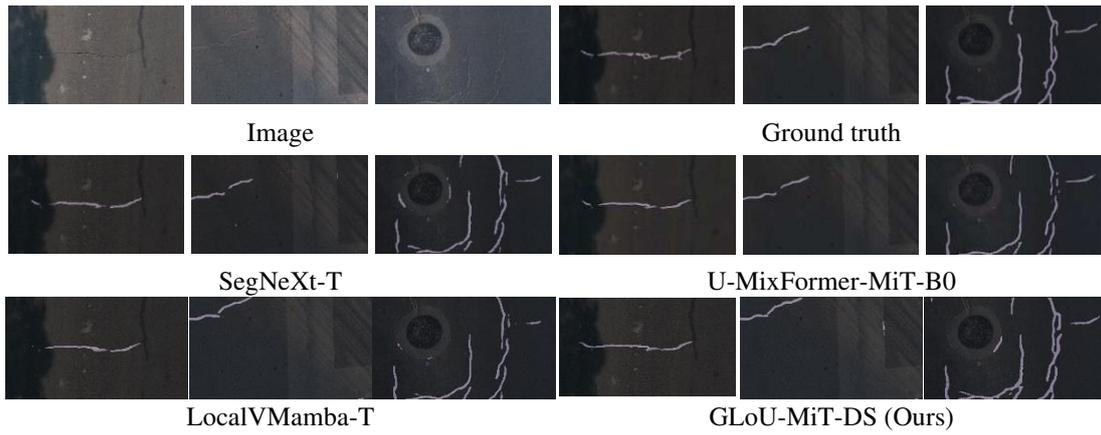
challenging to accurately segment the cracks. Although U-Net and Mamba-UNet

744 perform well in medical imaging, their segmentation performance on this dataset is
745 suboptimal. Our proposed GLoU-MiT model surpasses existing state-of-the-art models
746 based on CNNs, transformers, and Mamba architecture, outperforming the advanced
747 Vision Mamba model, LocalVMamba, with a 1.98% improvement in F_1 -score and a
748 2.72% increase in Crack IoU. After incorporating the DS module, GLoU-MiT-DS
749 achieves additional gains of 0.59% in F_1 -score and 0.86% in Crack IoU with a slight
750 increase of approximately 0.4M parameters and 0.7 GFLOPs.

751 **CrackSC:** The CrackSC dataset contains narrow cracks with significant
752 environmental interference. Strong models like SegNeXt-T and Mamba-UNet perform
753 worse than the traditional U-Net on this dataset. This is likely because U-Net applies
754 convolutions and downsampling directly at higher resolutions, which is advantageous
755 for detecting small cracks. However, U-Net’s convolution and downsampling at the
756 original resolution result in a substantial overhead in parameters and FLOPs. The
757 lightweight model U-MixFormer, which is carefully designed, achieved promising
758 results. In comparison, our GLoU-MiT model experienced a slight performance
759 decrease of approximately 0.1% in both F_1 and IoU. However, after incorporating the
760 DS supervision module, GLoU-MiT showed significant improvements, with F_1 and IoU
761 increasing by 0.21% and 0.27%, respectively.

762 **Crack500:** In Crack500, the cracks occupy a larger portion of the image and the
763 background is relatively simple, leading to smaller performance differences among
764 various SOTA models. Under these circumstances, traditional U-Net is outperformed
765 by more lightweight models. Our GLoU-MiT model achieves the best results,
766 improving upon the advanced transformer model U-MixFormer by 0.78% in F_1 and
767 1.09% in IoU. However, after adding the DS supervision module, the segmentation
768 performance decreased. This could be attributed to the edge supervision in the DS
769 module, which is generated by dilating the edges by 4 pixels. This approach may have
770 had a negative impact on the segmentation of clearly defined crack boundaries.

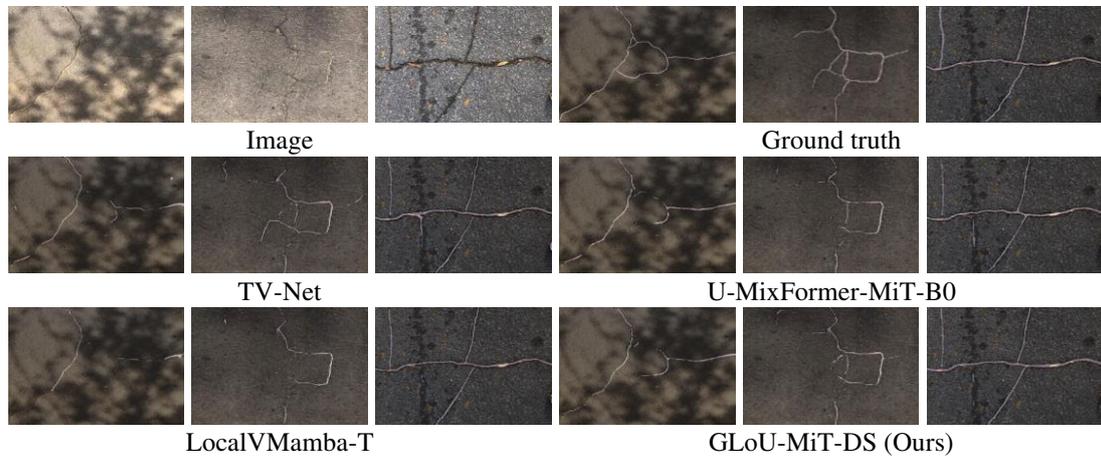
771



772

Fig. 9 Qualitative results on UAV-Crack500

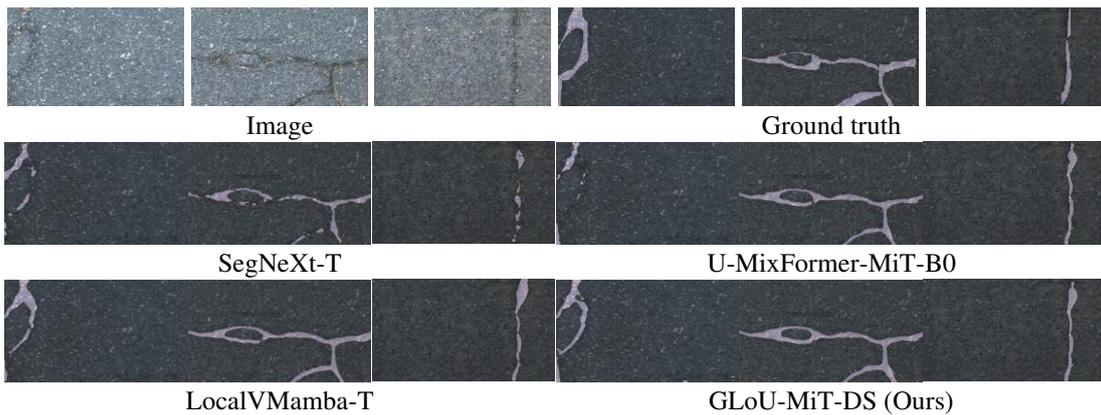
773



774

Fig. 10 Qualitative results on CrackSC

775



776

Fig. 11 Qualitative results on Crack500

777 5.2 Qualitative Evaluation

778 We selected three relatively challenging pavement images from each of the three
 779 datasets, and chose the best-performing models from the CNN, Transformer, and

780 Mamba network architectures for visualizing the segmentation results, as shown in Figs.
781 9-11. From these figures, it is evident that our model excels at distinguishing crack and
782 pavement pixels in complex scenes. For narrow and elongated cracks, segmentation
783 models often encounter discontinuities. However, our model performs better than other
784 models in producing continuous crack segmentations, significantly improving
785 segmentation accuracy.

786 **5.3 Energy-Delay Performance Evaluation**

787 As shown in Fig. 12, compiling the TorchScript model into C++ using LibTorch
788 significantly improves inference performance and efficiency by eliminating the
789 additional overhead and dynamic scheduling issues associated with the Python
790 interpreter. The results in Fig. 12(a) demonstrate that LibTorch-based inference
791 significantly reduces energy consumption per sample across most models, primarily
792 due to its static thread management strategy, which efficiently utilizes hardware
793 resources and minimizes thread switching and synchronization overhead.

794 However, for lightweight models such as SegFormer, U-MixFormer, and GLoU-
795 MiT, where computational loads are lower, the CUDA cores on the Jetson device are
796 not fully utilized. In these cases, the fixed threading model of LibTorch introduces task
797 scheduling overhead, resulting in an EDP that is not always significantly lower than
798 Python inference (Fig. 12(d)). Conversely, in heavier models like Mamba-UNet and
799 LocalVMamba, where computational resources are fully leveraged, LibTorch’s
800 performance advantage becomes more pronounced.

801 Among all models, Mamba-UNet exhibits the highest latency and energy
802 consumption, despite having lower FLOPs, due to the computationally expensive
803 Mamba operations applied on high-resolution inputs. In contrast, SegFormer achieves
804 a balance between inference speed and energy efficiency. Our proposed GLoU-MiT
805 model demonstrates competitive inference speed and energy consumption compared to
806 SegNeXt-T, but achieves higher segmentation accuracy across all datasets.

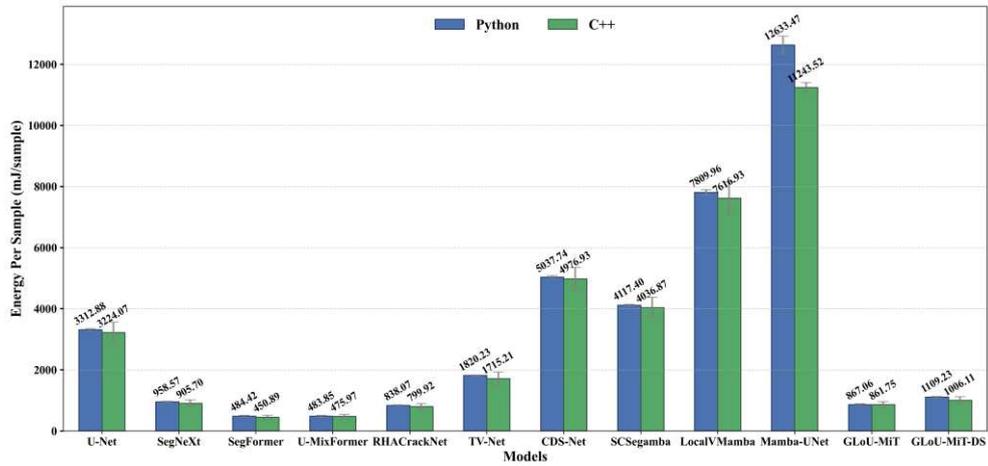
807 The Energy-Delay Product (EDP), a comprehensive measure of inference
808 efficiency, shows that after compilation into C++, GLoU-MiT achieves an EDP of
809 8.06×10^4 , which is 9.3% lower than SegNeXt-T (8.89×10^4), while improving the F₁-
810 score on the Crack500 dataset by 1.51%. Although the EDP is slightly higher compared
811 to U-MixFormer-MiT-B0 (2.56×10^4), GLoU-MiT improves the F₁-score on the
812 Crack500 dataset by 0.78%, demonstrating a favorable trade-off between accuracy and
813 efficiency.

814 With the incorporation of the Deep Supervision Refinement (DS) module, the EDP
815 of GLoU-MiT-DS increases to 10.77×10^4 , primarily due to the additional computation
816 introduced by boundary and semantic supervision. However, this results in a notable
817 improvement in segmentation performance, increasing the F₁-score on the UAV-
818 Crack500 dataset to 83% (a 0.59% improvement over GLoU-MiT) and boosting the F₁-
819 score on the CrackSC dataset from 74.21% to 74.52%.

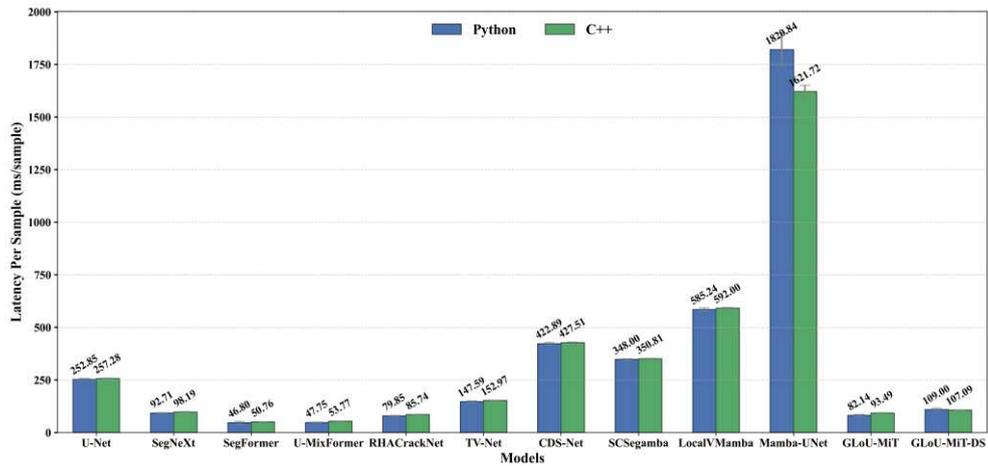
820 These findings suggest that while LibTorch inference generally reduces energy
821 consumption, its efficiency gains vary depending on model complexity and workload
822 distribution. Moreover, the DS module significantly enhances crack boundary
823 segmentation, demonstrating the trade-off between accuracy and computational
824 overhead in UAV-based crack detection scenarios.

825 Integrating the findings from Sections 5.1 to 5.3, it becomes evident that the
826 number of parameters and FLOPs is not necessarily indicative of a model’s real-world
827 inference efficiency. Instead, inference performance is jointly determined by
828 architectural design choices, the distribution and parallelizability of computation, and
829 the ability to mitigate execution bottlenecks. For example, although Mamba-based
830 models demonstrate favorable parameter efficiency, they exhibit suboptimal inference
831 speed on edge devices. This discrepancy can be attributed to the current lack of mature
832 GPU-level parallelization and hardware-specific optimization support for Mamba
833 operators.

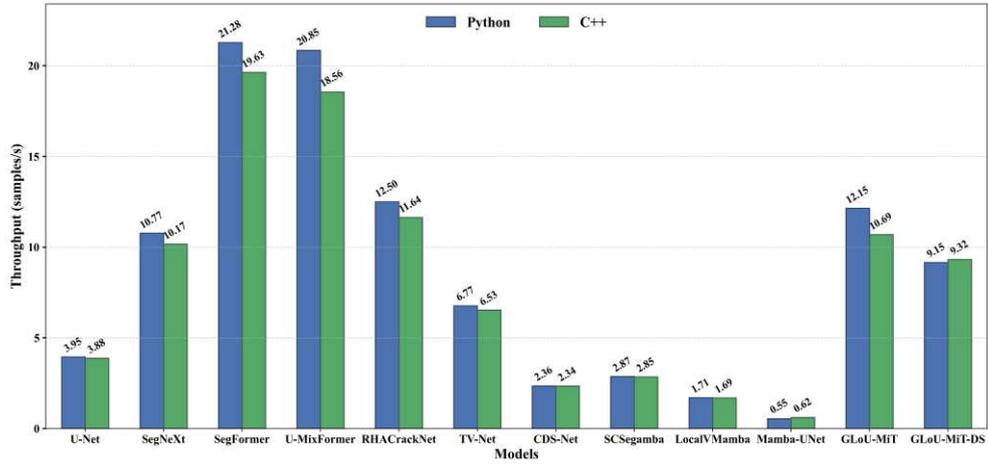
834 Therefore, the extensive deployment of Mamba modules is not advisable in
 835 latency-critical edge scenarios. Instead, our proposed architecture adopts a more
 836 pragmatic and effective approach by integrating Mamba modules within skip
 837 connections. This selective incorporation strategy enables the architecture to retain the
 838 modeling strengths of Mamba while mitigating its negative impact on inference latency.
 839 Furthermore, as Mamba currently lacks comprehensive GPU-oriented optimization and
 840 support strategies, future efforts to enhance inference speed may benefit from the
 841 integration of conventional model compression techniques, including quantization,
 842 pruning, and knowledge distillation.



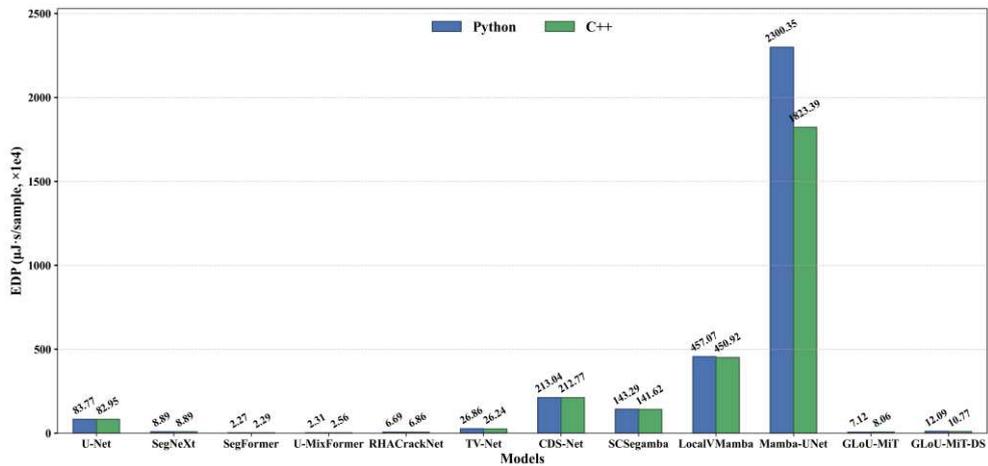
(a) Energy Per Sample (EPS)



(b) Latency Per Sample (LPS)



(c) Throughput



(d) Energy Delay Product (EDP)

843 **Fig. 12 Energy-Delay efficiency of Python and C++ on Jetson Orin Nano (8G)**

844 **5.4 Ablation Studies**

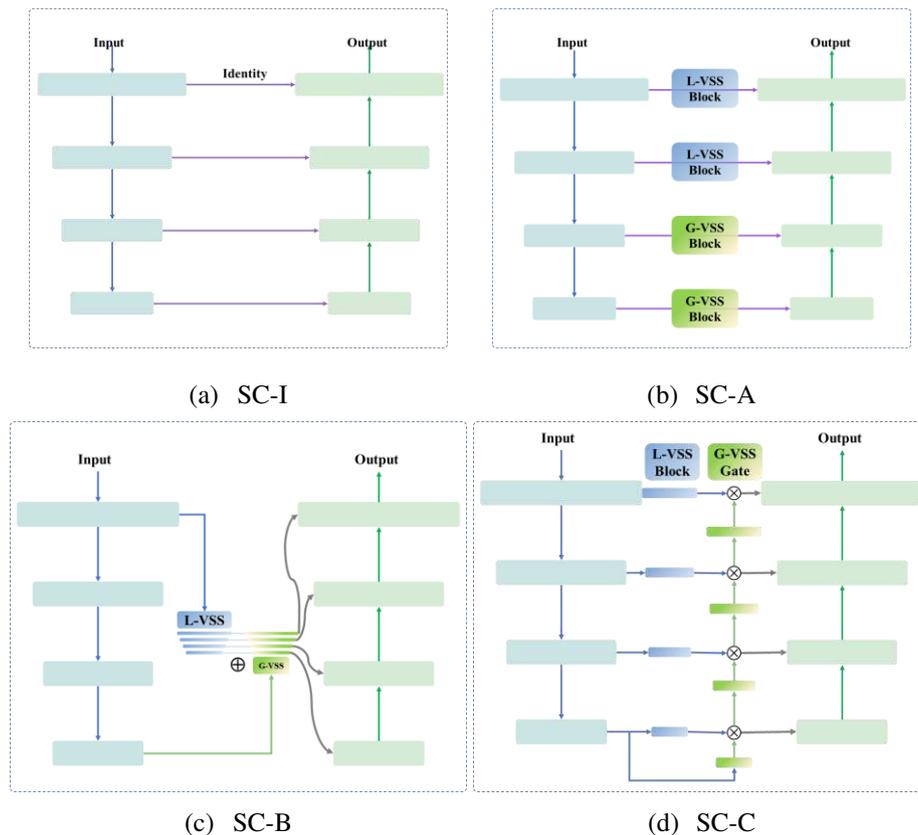
845 The carefully designed skip connections can effectively integrate low-level feature
 846 maps rich in detail with feature maps containing high-level semantic information. To
 847 validate the superiority of our model, we explored various designs for skip connections
 848 and proposed four different types (Fig.13):

- 849 (1) **SC-I**: The feature maps generated by the encoder are directly added to the
 850 decoder's feature maps without any additional operations.
- 851 (2) **GLo-SC-A**: In the shallow layers (Stage 1 and Stage 2), we use a local VSS
 852 block to align the detailed feature maps with the semantic feature maps. In the
 853 deeper layers (Stage 3 and Stage 4), a global VSS block is used to align the

854 semantic-rich feature maps from the encoder with the semantic feature maps in
855 the decoder.

856 (3) **GLo-SC-B**: In this design, the local VSS block is applied to the high-resolution
857 layer (Stage 1) of the encoder to extract local features, while the global VSS
858 block is used for the low-resolution layer (Stage 4) to extract global features.
859 The detailed and semantic feature maps are added directly at corresponding
860 levels.

861 (4) **GLo-SC-C**: As discussed in Section 3.4, we extract local features from each
862 layer using local VSS blocks, with the low-resolution and semantically rich
863 Layer 4 serving as the guiding layer. Through a progressive gating mechanism,
864 we gradually control the amount of local feature information passed from each
865 layer to its corresponding layer in the decoder.



866 **Fig. 13 Visualization of our proposed skip connection (SC) architectures**

867 As shown in Table 5, the performance improves when the Global-Local Mamba
868 skip connection (GLo-SC) modules A, B, and C are added, compared to the identity

869 add method. This indicates that by extracting both global and local features from the
870 skip connections, the semantic representation capability of the feature maps is enhanced,
871 allowing them to align more effectively with the feature maps in the decoder. The GLo-
872 SC-C module achieves the greatest improvement by innovatively integrating
873 hierarchical global feature extraction with an adaptive gating mechanism. This design
874 dynamically balances the incorporation of critical global context with the suppression
875 of redundant or noisy information, thereby enhancing semantic richness and ensuring
876 optimal fusion between encoder and decoder features. As a result, the model attains
877 superior predictive accuracy and improves robustness in capturing fine-grained details.

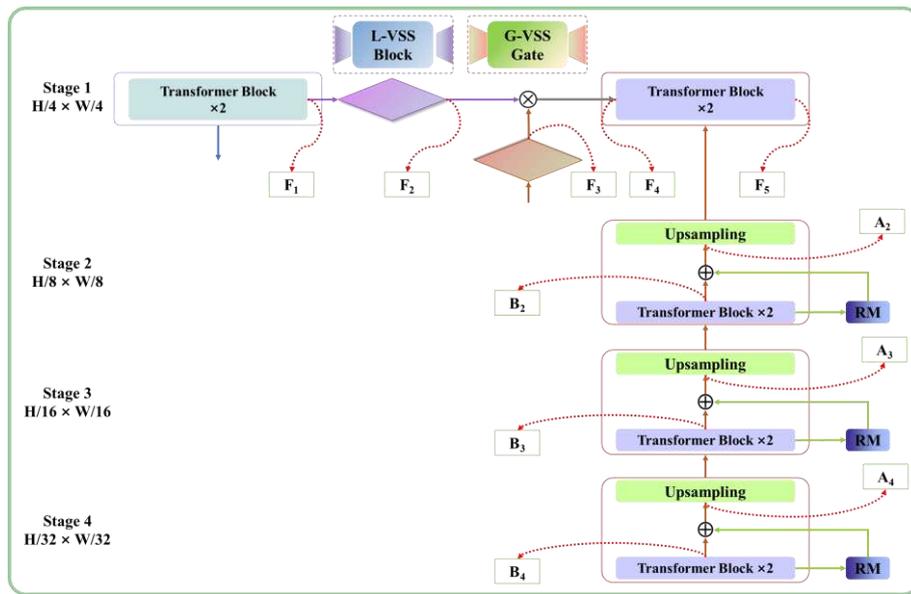
878 **Table 5** Ablation results on UAV-Crack500, CrackSC and Crack500

Model	Skip Connection				DS	Params (M)	Flops (G)	UAV-Crack500		CrackSC		Crack500	
	I	A	B	C				F ₁	IoU	F ₁	IoU	F ₁	IoU
U-MiT	✓					6.2	6.2	80.98	68.04	71.31	55.41	80.01	66.68
GLoU-MiT-A		✓				6.4	6.3	81.65	68.99	71.98	56.22	<u>80.11</u>	<u>66.82</u>
GLoU-MiT-B			✓			6.6	6.3	81.39	68.62	72.22	56.52	79.78	66.36
GLoU-MiT				✓		6.6	6.5	<u>82.41</u>	<u>70.08</u>	<u>74.21</u>	<u>59.00</u>	80.59	67.49
U-MiT-DS	✓				✓	6.6	7.0	81.51	68.79	71.91	56.14	79.56	66.05
GLoU-MiT-DS-A		✓			✓	6.7	7.1	81.69	69.05	73.19	57.71	80.08	66.77
GLoU-MiT-DS-B			✓		✓	6.9	7.1	82.32	69.95	72.62	57.01	79.90	66.53
GLoU-MiT-DS				✓	✓	7.0	7.2	83.00	70.94	74.52	59.39	79.99	66.66

879 Furthermore, we applied deep supervision to the different skip connection models
880 mentioned above by introducing a Boundary/Semantic Deep Supervision Refinement
881 Module. The results show significant improvement, particularly for low-resolution,
882 low-contrast datasets such as UAV-Crack500, and datasets with complex backgrounds
883 such as CrackSC. This improvement is attributed to the addition of boundary and
884 semantic supervision at each layer, which helps accelerate the learning of semantic
885 information and enhances the representation of deeper features. As a result, the model
886 can better understand and distinguish complex scenes or objects. However, for datasets
887 where cracks are relatively obvious (e.g., Crack500), the addition of deep supervision
888 leads to a decline in performance. This may be because the cracks are already prominent,
889 and the base model structure is sufficient to capture the necessary features. In this case,
890 deep supervision may over-constrain the learning of intermediate layers, thereby
891 negatively affecting the overall performance.

892 **5.5 Visualization of Activation Maps**

893 To understand the impact of our designed skip connection modules on model
 894 performance, we employed LayerCAM [76] to visualize the class activation maps of
 895 two images from the publicly available CrackSC and Crack500 datasets at different
 896 stages of our proposed GLoU-MiT models and the SegFormer model. Specifically, we
 897 visualized the feature maps in stage 1 (high-resolution layer) of our designed model,
 898 with the visualization locations shown in Fig. 14. In this figure, F_1 represents the output
 899 of stage 1 in the encoder, F_2 is the output after the local VSS block, F_3 is the output
 900 from the previous layer followed by the global VSS block, F_4 is the output after
 901 applying sigmoid to the element-wise multiplication of F_2 and F_3 , and F_5 is the final
 902 output of the decoder (corresponding to the MLP stage in SegFormer).



903
 904 **Fig.14 Diagram of LayerCAM visualization positions**

905 As can be seen from the Fig.15, the presence of environmental interference results
 906 in suboptimal activation of crack regions during the shallow stage (F_1), where only low-
 907 level features are extracted, leading to significant noise around the edges. After
 908 introducing the local VSS block, the activation maps F_2 become more focused around
 909 the crack regions. Compared to GLoU-MiT-B (SC-B) model that directly compute
 910 global features in stage 4, GLoU-MiT (SC-C), which gradually guides the feature
 911 extraction process through the global VSS gate, shows improved edge activation in its

912 feature maps F_3 . The feature maps F_4 processed by the global VSS gate reveals more
 913 comprehensive crack region activation. As a result, the final segmentation map
 914 demonstrates better continuity and improved performance in segmenting fine cracks in
 915 complex backgrounds.

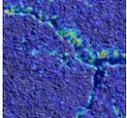
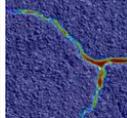
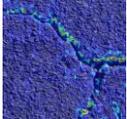
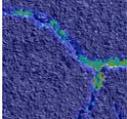
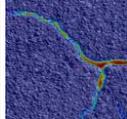
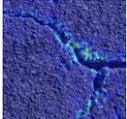
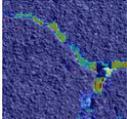
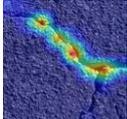
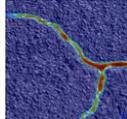
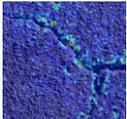
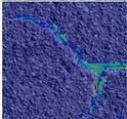
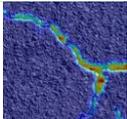
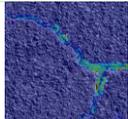
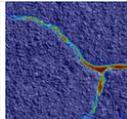
Model	F_1	F_2	F_3	F_4	F_5/MLP	Output
SegFormer-MiT-B0		N/A	N/A	N/A		
U-MiT		N/A	N/A	N/A		
GLoU-MiT-A			N/A	N/A		
GLoU-MiT-B				N/A		
GLoU-MiT						

(a) Image from the UAV-Crack500 dataset

Model	F_1	F_2	F_3	F_4	F_5/MLP	Output
SegFormer-MiT-B0		N/A	N/A	N/A		
U-MiT		N/A	N/A	N/A		
GLoU-MiT-A			N/A	N/A		
GLoU-MiT-B				N/A		
GLoU-MiT						

(b) Image from the CrackSC dataset

SegFormer-MiT-B0		N/A	N/A	N/A		
------------------	--	-----	-----	-----	--	--

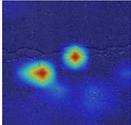
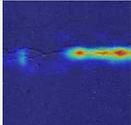
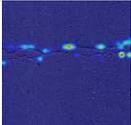
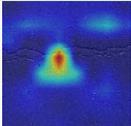
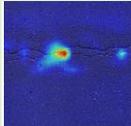
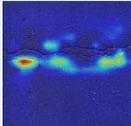
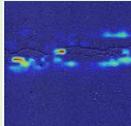
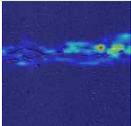
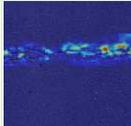
U-MiT		N/A	N/A	N/A		
GLoU-MiT-A			N/A	N/A		
GLoU-MiT-B				N/A		
GLoU-MiT						

(c) Image from the Crack500 dataset

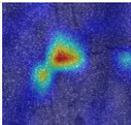
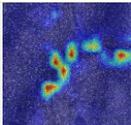
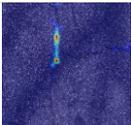
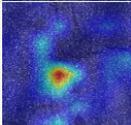
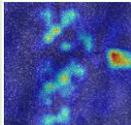
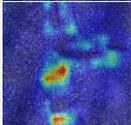
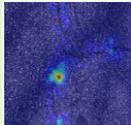
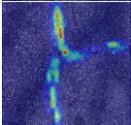
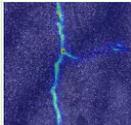
916 **Fig.15 LayerCAM visualizations: Comparing of GLoU-MiT and SegFormer at**
917 **corresponding layers**

918 We further compare feature maps before and after the insertion of the DS module,
919 as well as with and without the DS module. Specifically, we compare the LayerCAM
920 visualizations across decoder Stage 2 to Stage 4 (Fig. 16). For the model with the Deep
921 Supervision Refinement Module (GLoU-MiT-DS), the feature maps before DS
922 insertion are denoted as B_k and those after insertion as A_k (where k indicates the stage)
923 (Fig. 14). For the model without the DS module (GLoU-MiT), the corresponding
924 feature maps are similarly denoted as B_k . The LayerCAM results indicate that, after DS
925 insertion, the activation maps become more focused towards the crack centers and
926 exhibit finer boundary delineation, which enhances the segmentation accuracy.
927 Moreover, for datasets with finer cracks (e.g., CrackSC and UAVCrack500), the
928 activation in the shallower B_2 layer increases with the DS module, indicating that the
929 additional boundary and semantic supervision improves the overall crack detection. For
930 datasets with more prominent cracks (e.g., Crack500), the activation regions post-DS
931 insertion are more concentrated within the crack regions. These observations clearly
932 demonstrate that the deep supervision (DS) module plays a crucial role in refining the
933 boundary and semantic features, particularly for fine and narrow cracks that are
934 challenging to detect.

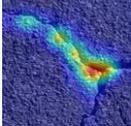
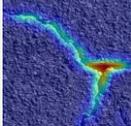
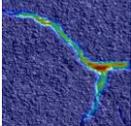
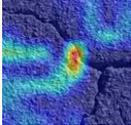
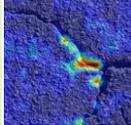
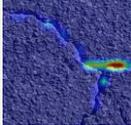
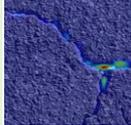
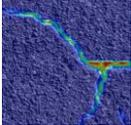
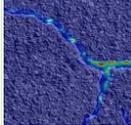
935

Model	B ₄	A ₄	B ₃	A ₃	B ₂	A ₂	
GLoU-MiT		N/A		N/A		N/A	
GLoU-MiT-DS							

(a) Image from the UAV-Crack500 dataset

Model	B ₄	A ₄	B ₃	A ₃	B ₂	A ₂	
GLoU-MiT		N/A		N/A		N/A	
GLoU-MiT-DS							

(b) Image from the CrackSC dataset

Model	B ₄	A ₄	B ₃	A ₃	B ₂	A ₂	
GLoU-MiT		N/A		N/A		N/A	
GLoU-MiT-DS							

(c) Image from the Crack500 dataset

936 **Fig.16 LayerCAM Visualizations: Comparison of GLoU-MiT (without DS) and GLoU-**
937 **MiT-DS (before and after DS insertion)**

938 6. Conclusion and Future Research

939 In this paper, we present a novel lightweight pavement crack segmentation model,
940 GLoU-MiT, designed to address the unique challenges posed by UAV-captured
941 pavement images, such as low resolution, fine crack structures, and low contrast. The
942 proposed model integrates three main components: a U-shaped segmentation
943 framework based on Mix Transformer, a Global-Local Mamba-Guided Skip
944 Connection mechanism, and a Deep Supervision Refinement Module. These
945 innovations contribute to improved feature extraction, efficient computation, and
946 precise segmentation of crack boundaries.

947 The U-Mix Transformer framework efficiently combines hierarchical feature
948 extraction and attention mechanisms to enhance segmentation accuracy while reducing
949 computational complexity. By replacing concatenation with direct addition in skip
950 connections, the model achieves more effective multi-level feature fusion. The
951 introduction of the Global-Local Mamba-Guided Skip Connection further improves
952 semantic representation by dynamically filtering and fusing global and local features.
953 Additionally, the deep supervision refinement module ensures accurate boundary and
954 semantic supervision, particularly for fine and narrow cracks that are often difficult to
955 detect. Comparative experiments on UAV-Crack500, CrackSC, and Crack500 datasets
956 demonstrate that GLoU-MiT outperforms state-of-the-art CNN, Transformer, and
957 Mamba-based models in terms of F_1 -score and Crack IoU, particularly in complex
958 scenarios with challenging crack structures.

959 Furthermore, while the absolute performance improvement of GLoU-MiT over
960 existing models appears to be around 1% or less, this gain remains both theoretically
961 and practically significant, particularly in UAV-based crack segmentation, where
962 challenges such as low resolution, fine crack structures, and environmental noise make
963 accurate detection inherently difficult. Even minor improvements in F_1 -score and IoU
964 can lead to more reliable crack identification, reduced false positives, and better
965 decision-making in automated road maintenance, ultimately enhancing infrastructure
966 monitoring efficiency. Although Mamba-based models excel in long-range feature
967 modeling, their high computational cost limits their feasibility for real-time edge
968 deployment. Our approach effectively balances segmentation accuracy, computational
969 efficiency, and inference speed, making it better suited for UAV-based applications.
970 Additionally, the incorporation of the DS module further enhances fine crack and
971 boundary segmentation, reinforcing the practical advantages of our method in complex
972 and challenging environments.

973 Looking ahead, future research will focus on further optimizing the inference
974 speed of Mamba-based models, improving their computational efficiency to enhance

975 the feasibility of real-time edge deployment on UAVs, ensuring faster and more
976 efficient crack detection in real-world infrastructure monitoring.

977 **CRedit author statement**

978 **Jinhuan Shan:** Methodology, Software, Data Curation, Writing - Original Draft.

979 **Yue Huang:** Conceptualization, Validation, Writing - Review & Editing. **Wei Jiang:**

980 Conceptualization, Resources, Supervision, Funding acquisition. **Dongdong Yuan:**

981 Writing - Review & Editing, Validation. **Feiyang Guo:** Software, Validation.

982

983 **Declaration of competing interest**

984 The authors declare that they have no known competing financial interests or

985 personal relationships that could have influenced the work reported in this study.

986

987 **Acknowledgments**

988 This work was supported by the National Natural Science Foundation of China

989 (Grant No. 52478432), Qin Chuangyuan “Scientist + Engineer” Team Construction

990 Project of Shaanxi Province (Grant No. 2024QCY-KXJ-020), and the Scientific

991 Research Project of Department of Transportation of Shaanxi Province (Grant No. 24-

992 35K).

993

994

995 **References**

- 996 [1] M. Zheng, W. Chen, X. Ding, W. Zhang, S. Yu, Comprehensive Life Cycle Environmental
997 Assessment of Preventive Maintenance Techniques for Asphalt Pavement, *Sustainability* 13
998 (2021) 4887. <https://doi.org/10.3390/su13094887>.
- 999 [2] T. Zhang, D. Wang, Y. Lu, ECSNet: An Accelerated Real-Time Image Segmentation CNN
1000 Architecture for Pavement Crack Detection, *IEEE Transactions on Intelligent Transportation*
1001 *Systems* (2023) 1–8. <https://doi.org/10.1109/TITS.2023.3300312>.
- 1002 [3] Pixel-wise crack defect segmentation with dual-encoder fusion network, *Constr. Build. Mater.*
1003 426 (2024) 136179. <https://doi.org/10.1016/j.conbuildmat.2024.136179>.
- 1004 [4] A comparison study of semantic segmentation networks for crack detection in construction
1005 materials, *Constr. Build. Mater.* 414 (2024) 134950.
1006 <https://doi.org/10.1016/j.conbuildmat.2024.134950>.
- 1007 [5] A. Ragnoli, M.R. De Blasiis, A. Di Benedetto, Pavement Distress Detection Methods: A
1008 Review, *Infrastructures* 3 (2018) 58. <https://doi.org/10.3390/infrastructures3040058>.
- 1009 [6] Y. Li, P. Che, C. Liu, D. Wu, Y. Du, Cross-scene pavement distress detection by a novel
1010 transfer learning framework, *Computer-Aided Civil and Infrastructure Engineering* 36 (2021)
1011 1398–1415. <https://doi.org/10.1111/mice.12674>.
- 1012 [7] Z. Tong, T. Ma, W. Zhang, J. Huyan, Evidential transformer for pavement distress
1013 segmentation, *Computer-Aided Civil and Infrastructure Engineering* 38 (2023) 2317–2338.
1014 <https://doi.org/10.1111/mice.13018>.
- 1015 [8] G. Zhu, J. Liu, Z. Fan, D. Yuan, P. Ma, M. Wang, W. Sheng, K.C.P. Wang, A lightweight
1016 encoder–decoder network for automatic pavement crack detection, *Computer-Aided Civil and*
1017 *Infrastructure Engineering* 39 (2023) 1743–1765. <https://doi.org/10.1111/mice.13103>.
- 1018 [9] Z. Han, J. Tang, L. Hu, W. Jiang, A. Sha, Automated measurement of asphalt pavement rut
1019 depth using smartphone imaging, *Autom. Constr.* 174 (2025) 106124.
1020 <https://doi.org/10.1016/j.autcon.2025.106124>.
- 1021 [10] J. Shan, W. Jiang, X. Feng, Bridging cross-domain and cross-resolution gaps for UAV-based
1022 pavement crack segmentation, *Autom. Constr.* 174 (2025) 106141.
1023 <https://doi.org/10.1016/j.autcon.2025.106141>.
- 1024 [11] Y. Zhang, Z. Zuo, X. Xu, J. Wu, J. Zhu, H. Zhang, J. Wang, Y. Tian, Road damage detection
1025 using UAV images based on multi-level attention mechanism, *Automation in Construction*
1026 144 (2022) 104613. <https://doi.org/10.1016/j.autcon.2022.104613>.
- 1027 [12] M. Chai, G. Li, W. Ma, D. Chen, Q. Du, Y. Zhou, S. Qi, L. Tang, H. Jia, Damage characteristics
1028 of the Qinghai-Tibet Highway in permafrost regions based on UAV imagery, *International*
1029 *Journal of Pavement Engineering* 0 (2022) 1–12.
1030 <https://doi.org/10.1080/10298436.2022.2038381>.
- 1031 [13] Xiao Jiang, Shanjun Mao, Mei Li, Hui Liu, Haoyuan Zhang, Shuwei Fang, Mingze Yuan, Chi
1032 Zhang, MFPA-Net: An efficient deep learning network for automatic ground fissures

- 1033 extraction in UAV images of the coal mining area, *Int. J. Appl. Earth Obs. Geoinf.* (2022).
 1034 <https://doi.org/10.1016/j.jag.2022.103039>.
- 1035 [14] X. Lei, C. Liu, L. Li, G. Wang, Automated Pavement Distress Detection and Deterioration
 1036 Analysis Using Street View Map, *IEEE Access* 8 (2020) 76163–76172.
 1037 <https://doi.org/10.1109/ACCESS.2020.2989028>.
- 1038 [15] L. Song, X. Wang, Faster region convolutional neural network for automated pavement
 1039 distress detection, *Road Materials and Pavement Design* 22 (2021) 23–41.
 1040 <https://doi.org/10.1080/14680629.2019.1614969>.
- 1041 [16] Y. Du, N. Pan, Z. Xu, F. Deng, Y. Shen, H. Kang, Pavement distress detection and
 1042 classification based on YOLO network, *International Journal of Pavement Engineering* 22
 1043 (2021) 1659–1672. <https://doi.org/10.1080/10298436.2020.1714047>.
- 1044 [17] Z. Qu, C.-Y. Wang, S.-Y. Wang, F.-R. Ju, A Method of Hierarchical Feature Fusion and
 1045 Connected Attention Architecture for Pavement Crack Detection, *IEEE Transactions on*
 1046 *Intelligent Transportation Systems* 23 (2022) 16038–16047.
 1047 <https://doi.org/10.1109/TITS.2022.3147669>.
- 1048 [18] X. Sun, Y. Xie, L. Jiang, Y. Cao, B. Liu, DMA-Net: DeepLab With Multi-Scale Attention for
 1049 Pavement Crack Segmentation, *IEEE Transactions on Intelligent Transportation Systems* 23
 1050 (2022) 18392–18403. <https://doi.org/10.1109/TITS.2022.3158670>.
- 1051 [19] Q. Zhou, Z. Qu, S.-Y. Wang, K.-H. Bao, A Method of Potentially Promising Network for
 1052 Crack Detection With Enhanced Convolution and Dynamic Feature Fusion, *IEEE*
 1053 *Transactions on Intelligent Transportation Systems* 23 (2022) 18736–18745.
 1054 <https://doi.org/10.1109/TITS.2022.3154746>.
- 1055 [20] P.R.T. Peddinti, H. Puppala, B. Kim, Pavement Monitoring Using Unmanned Aerial Vehicles:
 1056 An Overview, *Journal of Transportation Engineering, Part B: Pavements* 149 (2023)
 1057 03123002. <https://doi.org/10.1061/JPEODX.PVENG-1291>.
- 1058 [21] H. Liu, J. Yang, X. Miao, C. Mertz, H. Kong, CrackFormer network for pavement crack
 1059 segmentation, *IEEE Trans. Intell. Transp. Syst.* 24 (2023) 9240–9252.
 1060 <https://doi.org/10.1109/TITS.2023.3266776>.
- 1061 [22] J. Wang, Z. Zeng, P.K. Sharma, O. Alfarraj, A. Tolba, J. Zhang, L. Wang, Dual-path network
 1062 combining CNN and transformer for pavement crack segmentation, *Autom. Constr.* 158 (2024)
 1063 105217. <https://doi.org/10.1016/j.autcon.2023.105217>.
- 1064 [23] F. Guo, Y. Qian, J. Liu, H. Yu, Pavement crack detection based on transformer network, *Autom.*
 1065 *Constr.* 145 (2023) 104646. <https://doi.org/10.1016/j.autcon.2022.104646>.
- 1066 [24] J. Shan, Y. Huang, W. Jiang, DCUFormer: enhancing pavement crack segmentation in
 1067 complex environments with dual-cross/upsampling attention, *Expert Syst. Appl.* 264 (2025)
 1068 125891. <https://doi.org/10.1016/j.eswa.2024.125891>.
- 1069 [25] A. Gu, T. Dao, Mamba: Linear-Time Sequence Modeling with Selective State Spaces, (2024).
 1070 <http://arxiv.org/abs/2312.00752> (accessed August 25, 2024).

- 1071 [26] Z. Zhang, B. Peng, T. Zhao, An ultra-lightweight network combining mamba and frequency-
1072 domain feature extraction for pavement tiny-crack segmentation, *Expert Syst. Appl.* 264
1073 (2025) 125941. <https://doi.org/10.1016/j.eswa.2024.125941>.
- 1074 [27] X. Zuo, Y. Sheng, J. Shen, Y. Shan, Topology-aware mamba for crack segmentation in
1075 structures, *Autom. Constr.* 168 (2024) 105845. <https://doi.org/10.1016/j.autcon.2024.105845>.
- 1076 [28] C. Han, H. Yang, Y. Yang, Enhancing pixel-level crack segmentation with visual mamba and
1077 convolutional networks, *Automation in Construction* 168 (2024) 105770.
1078 <https://doi.org/10.1016/j.autcon.2024.105770>.
- 1079 [29] H. Liu, C. Jia, F. Shi, X. Cheng, S. Chen, SCSegamba: lightweight structure-aware vision
1080 mamba for crack segmentation in structures, (2025).
1081 <https://doi.org/10.48550/arXiv.2503.01113>.
- 1082 [30] Z. He, Y.-H. Wang, Mamba meets crack segmentation, (2024).
1083 <https://doi.org/10.48550/arXiv.2407.15714>.
- 1084 [31] J. Zhong, J. Zhu, J. Huyan, T. Ma, W. Zhang, Multi-scale feature fusion network for pixel-
1085 level pavement distress detection, *Automat Constr* 141 (2022) 104436.
1086 <https://doi.org/10.1016/j.autcon.2022.104436>.
- 1087 [32] Y. Wang, Z. He, X. Zeng, J. Zeng, Z. Cen, L. Qiu, X. Xu, Q. Zhuo, GGMNet: pavement-crack
1088 detection based on global context awareness and multi-scale fusion, *Remote Sens.* 16 (2024)
1089 1797. <https://doi.org/10.3390/rs16101797>.
- 1090 [33] S. Bai, L. Yang, Y. Liu, H. Yu, DMF-net: a dual-encoding multi-scale fusion network for
1091 pavement crack detection, *IEEE Trans. Intell. Transp. Syst.* 25 (2024) 5981–5996.
1092 <https://doi.org/10.1109/TITS.2023.3331769>.
- 1093 [34] J. Dong, N. Wang, H. Fang, W. Guo, B. Li, K. Zhai, MFAFNet: an innovative crack intelligent
1094 segmentation method based on multi-layer feature association fusion network, *Adv. Eng. Inf.*
1095 62 (2024) 102584. <https://doi.org/10.1016/j.aei.2024.102584>.
- 1096 [35] Q. Zhou, Z. Qu, Y. Li, F. Ju, Tunnel Crack Detection With Linear Seam Based on Mixed
1097 Attention and Multiscale Feature Fusion, *IEEE Trans. Instrum. Meas.* (2022).
1098 <https://doi.org/10.1109/tim.2022.3184351>.
- 1099 [36] J. Liang, X. Gu, D. Jiang, Q. Zhang, CNN-based network with multi-scale context feature and
1100 attention mechanism for automatic pavement crack segmentation, *Autom. Constr.* 164 (2024)
1101 105482. <https://doi.org/10.1016/j.autcon.2024.105482>.
- 1102 [37] M. Ma, L. Yang, Y. Liu, H. Yu, An attention-based progressive fusion network for pixelwise
1103 pavement crack detection, *Measurement* 226 (2024) 114159.
1104 <https://doi.org/10.1016/j.measurement.2024.114159>.
- 1105 [38] A.M. Roy, J. Bhaduri, DenseSPH-YOLOv5: an automated damage detection model based on
1106 DenseNet and swin-transformer prediction head-enabled YOLOv5 with attention mechanism,
1107 *Adv. Eng. Inf.* 56 (2023) 102007. <https://doi.org/10.1016/j.aei.2023.102007>.
- 1108 [39] Z. Al-Huda, B. Peng, R.N.A. Algburi, M.A. Al-antari, R. AL-Jarazi, D. Zhai, A hybrid deep

- 1109 learning pavement crack semantic segmentation, *Eng. Appl. Artif. Intell.* 122 (2023) 106142.
1110 <https://doi.org/10.1016/j.engappai.2023.106142>.
- 1111 [40] J.-M. Guo, H. Markoni, J.-D. Lee, BARNet: boundary aware refinement network for crack
1112 detection, *IEEE Trans. Intell. Transp. Syst.* 23 (2022) 7343–7358.
1113 <https://doi.org/10.1109/TITS.2021.3069135>.
- 1114 [41] J. Long, E. Shelhamer, T. Darrell, Fully Convolutional Networks for Semantic Segmentation,
1115 (2015). <https://doi.org/10.48550/arXiv.1411.4038>.
- 1116 [42] O. Ronneberger, P. Fischer, T. Brox, U-net: convolutional networks for biomedical image
1117 segmentation, in: N. Navab, J. Hornegger, W.M. Wells, A.F. Frangi (Eds.), *Medical Image*
1118 *Computing and Computer-Assisted Intervention – MICCAI 2015*, Springer International
1119 Publishing, Cham, 2015: pp. 234–241. https://doi.org/10.1007/978-3-319-24574-4_28.
- 1120 [43] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-Decoder with Atrous
1121 Separable Convolution for Semantic Image Segmentation, in: *In Proceedings of the European*
1122 *Conference on Computer Vision (ECCV)*, Springer International Publishing, Cham, 2018: pp.
1123 833–851. https://doi.org/10.1007/978-3-030-01234-2_49.
- 1124 [44] J. Hu, L. Shen, S. Albanie, G. Sun, E. Wu, Squeeze-and-Excitation Networks, (2019).
1125 <https://doi.org/10.48550/arXiv.1709.01507>.
- 1126 [45] J. Park, S. Woo, J.-Y. Lee, I.S. Kweon, BAM: Bottleneck Attention Module, (2018).
1127 <https://doi.org/10.48550/arXiv.1807.06514>.
- 1128 [46] S. Woo, J. Park, J.-Y. Lee, I.S. Kweon, CBAM: Convolutional Block Attention Module, in: V.
1129 Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), *Computer Vision – ECCV 2018*,
1130 Springer International Publishing, Cham, 2018: pp. 3–19. https://doi.org/10.1007/978-3-030-01234-2_1.
- 1132 [47] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, B. Guo, Swin Transformer:
1133 Hierarchical Vision Transformer using Shifted Windows, (2021).
1134 <https://doi.org/10.48550/arXiv.2103.14030>.
- 1135 [48] Y. Liu, Y. Tian, Y. Zhao, H. Yu, L. Xie, Y. Wang, Q. Ye, Y. Liu, VMamba: Visual State Space
1136 Model, (2024). <http://arxiv.org/abs/2401.10166> (accessed August 15, 2024).
- 1137 [49] T. Huang, X. Pei, S. You, F. Wang, C. Qian, C. Xu, LocalMamba: Visual State Space Model
1138 with Windowed Selective Scan, (2024). <http://arxiv.org/abs/2403.09338> (accessed August 14,
1139 2024).
- 1140 [50] Z. Wang, J.-Q. Zheng, Y. Zhang, G. Cui, L. Li, Mamba-UNet: UNet-Like Pure Visual Mamba
1141 for Medical Image Segmentation, (2024). <http://arxiv.org/abs/2402.05079> (accessed August
1142 14, 2024).
- 1143 [51] T. Takikawa, D. Acuna, V. Jampani, S. Fidler, Gated-SCNN: Gated Shape CNNs for Semantic
1144 Segmentation, (2019). <https://doi.org/10.48550/arXiv.1907.05740>.
- 1145 [52] X. Qin, D.-P. Fan, C. Huang, C. Diagne, Z. Zhang, A.C. Sant’Anna, A. Suárez, M. Jagersand,
1146 L. Shao, Boundary-Aware Segmentation Network for Mobile and Web Applications, (2021).

- 1147 <https://doi.org/10.48550/arXiv.2101.04704>.
- 1148 [53] A. Kirillov, Y. Wu, K. He, R. Girshick, PointRend: Image Segmentation as Rendering, (2020).
- 1149 <https://doi.org/10.48550/arXiv.1912.08193>.
- 1150 [54] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, Y. Wei, Deformable Convolutional Networks,
- 1151 (2017). <https://doi.org/10.48550/arXiv.1703.06211>.
- 1152 [55] H. Kervadec, J. Bouchtiba, C. Desrosiers, E. Granger, J. Dolz, I. Ben Ayed, Boundary loss for
- 1153 highly unbalanced segmentation, *Medical Image Analysis* 67 (2021) 101851.
- 1154 <https://doi.org/10.1016/j.media.2020.101851>.
- 1155 [56] A. Gu, A. Gupta, K. Goel, C. Ré, On the parameterization and initialization of diagonal state
- 1156 space models, (2022). <https://doi.org/10.48550/arXiv.2206.11893>.
- 1157 [57] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, X. Wang, Vision mamba: efficient visual
- 1158 representation learning with bidirectional state space model, (2024).
- 1159 <http://arxiv.org/abs/2401.09417> (accessed October 12, 2024).
- 1160 [58] J. Shan, W. Jiang, Y. Huang, D. Yuan, Y. Liu, Unmanned Aerial Vehicle (UAV)-Based
- 1161 Pavement Image Stitching Without Occlusion, Crack Semantic Segmentation, and
- 1162 Quantification, *IEEE Trans. Intell. Transport. Syst.* 25 (2024) 17038–17053.
- 1163 <https://doi.org/10.1109/TITS.2024.3424525>.
- 1164 [59] F. Yang, L. Zhang, S. Yu, D. Prokhorov, X. Mei, H. Ling, Feature pyramid and hierarchical
- 1165 boosting network for pavement crack detection, *IEEE Trans. Intell. Transp. Syst.* 21 (2020)
- 1166 1525–1535. <https://doi.org/10.1109/TITS.2019.2910595>.
- 1167 [60] MMSegmentation Contributors, OpenMMLab semantic segmentation toolbox and
- 1168 benchmark, (2020). <https://github.com/open-mmlab/msegmentation> (accessed March 8,
- 1169 2025).
- 1170 [61] H. Zhang, A.A. Zhang, Z. Dong, A. He, Y. Liu, Y. Zhan, K.C.P. Wang, Robust semantic
- 1171 segmentation for automatic crack detection within pavement images using multi-mixing of
- 1172 global context and local image features, *IEEE Trans. Intell. Transp. Syst.* 25 (2024) 11282–
- 1173 11303. <https://doi.org/10.1109/TITS.2024.3360263>.
- 1174 [62] H. Li, H. Zhang, H. Zhu, K. Gao, H. Liang, J. Yang, Automatic crack detection on concrete
- 1175 and asphalt surfaces using semantic segmentation network with hierarchical transformer, *Eng.*
- 1176 *Struct.* 307 (2024) 117903. <https://doi.org/10.1016/j.engstruct.2024.117903>.
- 1177 [63] B. Liu, J. Kang, H. Guan, X. Zhi, Y. Yu, L. Ma, D. Peng, L. Xu, D. Wang, RTCNet: a novel
- 1178 real-time triple branch network for pavement crack semantic segmentation, *Int. J. Appl. Earth*
- 1179 *Obs. Geoinf.* 136 (2025) 104347. <https://doi.org/10.1016/j.jag.2024.104347>.
- 1180 [64] M. Usman, H. Chen, M. Raza, A lightweight convolutional transformer architecture approach
- 1181 for crack segmentation in safety assessment, in: 2024 IEEE International Conference on
- 1182 Systems, Man, and Cybernetics (Smc), 2024: pp. 5195–5199.
- 1183 <https://doi.org/10.1109/SMC54092.2024.10831651>.
- 1184 [65] W. Zheng, X. Jiang, Z. Fang, Y. Gao, TV-net: a structure-level feature fusion network based

1185 on tensor voting for road crack segmentation, *IEEE Trans. Intell. Transp. Syst.* 25 (2024)
1186 5743–5754. <https://doi.org/10.1109/TITS.2023.3334266>.

1187 [66] T. Ruan, T. Liu, Z. Huang, Y. Wei, S. Wei, Y. Zhao, T. Huang, Devil in the details: towards
1188 accurate single and multiple human parsing, (2018). <http://arxiv.org/abs/1809.05996>
1189 (accessed September 15, 2024).

1190 [67] J. Xu, Z. Xiong, S.P. Bhattacharyya, PIDNet: a real-time semantic segmentation network
1191 inspired by PID controllers, in: *2023 IEEE/CVF Conference on Computer Vision and Pattern
1192 Recognition (CVPR)*, IEEE, Vancouver, BC, Canada, 2023: pp. 19529–19539.
1193 <https://doi.org/10.1109/CVPR52729.2023.01871>.

1194 [68] F. Panella, A. Lipani, J. Boehm, Semantic segmentation of cracks: data challenges and
1195 architecture, *Autom. Constr.* 135 (2022) 104110.
1196 <https://doi.org/10.1016/j.autcon.2021.104110>.

1197 [69] Y. Zhang, J. Wu, Q. Li, X. Zhao, M. Tan, Beyond crack: fine-grained pavement defect
1198 segmentation using three-stream neural networks, *IEEE Trans. Intell. Transp. Syst.* 23 (2022)
1199 14820–14832. <https://doi.org/10.1109/TITS.2021.3134374>.

1200 [70] X. Weng, Y. Huang, W. Wang, Segment-based pavement crack quantification, *Automation in
1201 Construction* 105 (2019) 102819. <https://doi.org/10.1016/j.autcon.2019.04.014>.

1202 [71] M.-H. Guo, C.-Z. Lu, Q. Hou, Z.-N. Liu, M.-M. Cheng, S.-M. Hu, SegNeXt: rethinking
1203 convolutional attention design for semantic segmentation, in: *Proceedings of the 36th
1204 International Conference on Neural Information Processing Systems*, Curran Associates Inc.,
1205 Red Hook, NY, USA, 2022: pp. 1140–1156.

1206 [72] G. Zhu, J. Liu, Z. Fan, D. Yuan, P. Ma, M. Wang, W. Sheng, K.C.P. Wang, A lightweight
1207 encoder–decoder network for automatic pavement crack detection, *Computer-Aided Civil and
1208 Infrastructure Engineering* 39 (2024) 1743–1765. <https://doi.org/10.1111/mice.13103>.

1209 [73] Q. Tan, A. Li, L. Dong, W. Dong, X. Li, G. Shi, CDS-net: contextual difference sensitivity
1210 network for pixel-wise road crack detection, *IEEE Trans. Circuits Syst. Video Technol.* (2025)
1211 1–1. <https://doi.org/10.1109/TCSVT.2025.3529039>.

1212 [74] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J.M. Alvarez, P. Luo, SegFormer: Simple and
1213 Efficient Design for Semantic Segmentation with Transformers, in: *Advances in Neural
1214 Information Processing Systems*, Curran Associates, Inc., 2021: pp. 12077–12090.
1215 [https://proceedings.neurips.cc/paper_files/paper/2021/hash/64f1f27bf1b4ec22924fd0acb550
1216 c235-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2021/hash/64f1f27bf1b4ec22924fd0acb550c235-Abstract.html) (accessed July 25, 2024).

1217 [75] S.-K. Yeom, J. von Klitzing, U-MixFormer: UNet-like Transformer with Mix-Attention for
1218 Efficient Semantic Segmentation, (2023). <http://arxiv.org/abs/2312.06272> (accessed
1219 December 27, 2023).

1220 [76] P.-T. Jiang, C.-B. Zhang, Q. Hou, M.-M. Cheng, Y. Wei, LayerCAM: exploring hierarchical
1221 class activation maps for localization, *IEEE Trans. Image Process.* 30 (2021) 5875–5888.
1222 <https://doi.org/10.1109/TIP.2021.3089943>.