# A FEW-SHOT ATTENTION RECURRENT RESIDUAL U-NET FOR CRACK SEGMENTATION

*Iason Katsamenis*⋆, *Eftychios Protopapadakis*†, *Nikolaos Bakalos*⋆,
*Anastasios Doulamis*⋆, *Nikolaos Doulamis*⋆ *and Athanasios Voulodimos*⋆

⋆ National Technical University of Athens, 9th Iroon Polytechniou str., 15773 Athens, Greece
† University of Macedonia, 156th Egnatia str., 54636 Thessaloniki, Greece

**ABSTRACT**

Recent studies indicate that deep learning plays a crucial role in the automated visual inspection of road infrastructures. However, current learning schemes are static, implying no dynamic adaptation to users' feedback. To address this drawback, we present a few-shot learning paradigm for the automated segmentation of road cracks, which is based on a U-Net architecture with recurrent residual and attention modules (R2AU-Net). The retraining strategy dynamically fine-tunes the weights of the U-Net as a few new rectified samples are being fed into the classifier. Extensive experiments show that the proposed few-shot R2AU-Net framework outperforms other state-of-the-art networks in terms of Dice and IoU metrics, on a new dataset, named CrackMap, which is made publicly available at https://github.com/ikatsamenis/CrackMap.

*Index Terms*— Semantic segmentation, U-Net, attention, recurrent residual convolutional unit, road cracks

## 1. INTRODUCTION

The development of cracks on the road surface is a frequently occurring defect and can constitute a safety hazard for road users. Cracking in its various types (longitudinal, oblique, alligator cracks, etc.) affects the traffic flow and safety, resulting in poor performance of the road infrastructure, accidents, as well as increased $CO_2$ emissions, fuel costs, and time delays. Indicatively, for 2006, the comprehensive cost of traffic crashes where road conditions contributed to crash occurrence or severity, in the United States alone, is estimated at $217.5 billion, which corresponds to 43.6% of the total crash costs [1]. More recent evidence highlights that approximately $400 billion is invested globally each year in pavement construction and maintenance [2]. Therefore, the adoption of effective monitoring strategies can lead to enormous economic, social, and environmental benefits to the community.

Recently, there is a great research interest in the automatic visual inspection of road distress, by analyzing visual

data. Generally, deteriorated pavement produces rough surfaces, which entails that various image processing methods such as thresholding [3], edge detection [4], and mathematical morphology [5] can be used to localize crack regions. The core idea behind these approaches is that cracked regions tend to demonstrate non-uniformity, while on the other hand, the color and textural characteristics of the non-deteriorated road surface are more consistent and smoother. However, even though such techniques are computationally efficient, they are susceptible to image noise and fail to generalize the differentiation between the defect and the surface background [6].
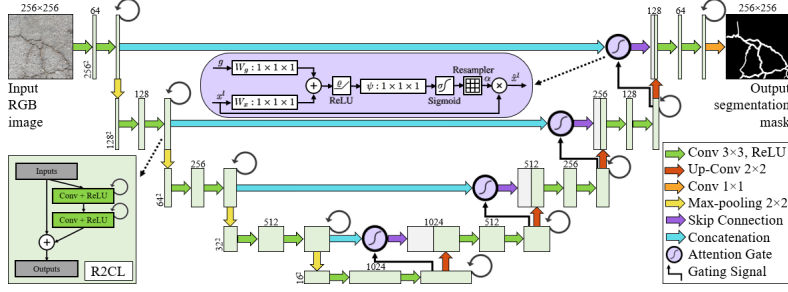
Current developments in deep learning and artificial intelligence technology have led Convolutional Neural Networks (CNNs) to be an effective tool for the automatic visual inspection of road infrastructures [7, 8, 9]. The main asset of such deep architectures, compared to the aforementioned conventional image processing methods, is the fact that they leverage throughout the learning procedure annotated ground truth data [10]. Thereby, these algorithms demonstrate high identification accuracy by effectively learning the essential features needed to classify a given pixel as defective or not.

Usually, Fully Convolutional Networks (FCNs), or their variants U-Nets, are considered for providing a precise pixel-based segmentation of damaged areas from RGB images of road infrastructures [6, 11, 12]. This is mainly due to the fact that FCNs have emerged as powerful segmentation tools, especially for performing accurate pixel-based classification tasks for challenging problems (e.g., biomedical imaging problems with data expressed either in 2D or 3D [13, 14], as well as crack segmentation [15]). To enhance the performance of the original U-Net, numerous elements have been introduced, such as residual convolutional units [16, 17] and attention gates instead of the typical skip connection [18, 19].
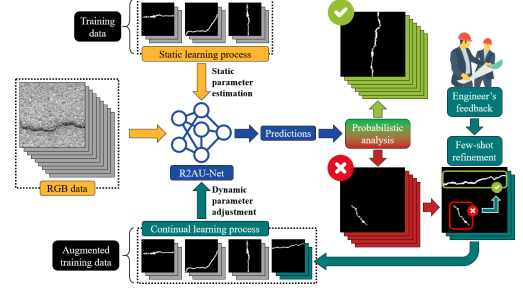
### 1.1. Current limitations and our contribution

Inspired by the above research work, we present R2AU-Net, which is a deep U-Net structure with recurrent residual and attention modules (see Fig. 1a). Compared to the standard U-Net, R2AU-Net incorporates recurrent residual convolutional layers (R2CL) that ensure better feature representation for the segmentation task and attention gates to highlight salient fea-

(a) The proposed U-Net architecture with recurrent residual and attention modules.

(b) The proposed few-shot retraining strategy.

**Fig. 1**. A schematic representation of the proposed (a) R2AU-Net and (b) few-shot learning scheme for road crack segmentation.

tures which are passed through the skip connections [17, 18].

Still, however, the most crucial hindrance of the aforementioned typical deep learning frameworks is that they treat the segmentation task as a static procedure. More specifically, they leverage knowledge derived from labeled data, but, nevertheless, it is not possible to further refine their outputs by exploiting user's interaction, especially in cases where the deep network underperforms [14]. To this end, we introduce a few-shot refinement scheme which is a semi-supervised learning paradigm, based on the R2AU-Net, that is able to adapt the model's behavior and weights according to user's feedback, to further increase the segmentation performance (see Fig. 1b).

## 2. PROPOSED ARCHITECTURE

### 2.1. R2AU-Net architecture for road crack segmentation

In this section, we present R2AU-Net, whose architecture is illustrated in Fig. 1a and is a combination of Recurrent Residual U-Net [17] and Attention U-Net [18]. It is emphasized that the model was designed for segmenting cracks in RGB images and, hence, provides meticulous information on a variety of metrics and properties that are critical for the automated and robotic-driven maintenance process, such as geometry, type, orientation, length, density, and shape of cracks.

In particular, the operations within the Recurrent Convolutional Layers (RCL) in R2CL (see Fig. 1a) are carried out based on the discrete time steps which are expressed according to the RCNN [20]. Suppose we have an input at layer $l$ within an R2CL block, and a pixel located at $(i, j)$ within an input on the feature map $k$ in the RCL. We denote $\mathcal{Y}_{ijk}^l(t)$ the output of the model at time step $t$, which can be expressed as:

$$\mathcal{Y}_{ijk}^l(t) = (w_k^\varrho)^T \cdot x_l^{\varrho(i,j)}(t) + (w_k^r)^T \cdot x_l^{r(i,j)}(t-1) + \beta_k \quad (1)$$

where $x_l^{\varrho(i,j)}(t)$ and $x_l^{r(i,j)}(t-1)$ represent respectively the inputs to the standard convolution layers and $l^{th}$ RCL. In parallel, $w_k^\varrho$ and $w_k^r$ denote respectively the weights of the standard convolutional layer and RCL that correspond to the $k^{th}$ feature map, whereas $\beta_k$ symbolizes the bias. The RCL's output is activated by the ReLU function $\varrho$ as follows:

$$\mathcal{R}(x_l, w_l) = \varrho(\mathcal{Y}_{ijk}^l(t)) = max(0, \mathcal{Y}_{ijk}^l(t)) \quad (2)$$

Let $x_l$ an input sample of the R2CL unit, then, the R2CL's output $x_{l+1}$, in both the downsampling layer in the encoding path and the upsampling layer in the decoding path of the R2AU-Net, can be calculated using the following equation:

$$x_{l+1} = x_l + \mathcal{R}(x_l, w_l) \quad (3)$$

In parallel, as one can observe in Fig. 1a we integrate into R2AU-Net an attention gate mechanism to focus on points and shapes of interest (i.e., road cracks) [18]. More specifically, for each pixel $i$, the attention coefficients $\alpha_i \in [0, 1]$ tend to yield higher values in target crack regions and lower values in background road areas. We obtain the output of the attention gate in layer $l$ by multiplying element-wise the input feature maps and attention coefficients: $\hat{x}_i^l = x_i^l \cdot \alpha_i^l$. Attention values are calculated for each pixel vector $x_i^l \in \mathbb{R}^{F_l}$, where $F_l$ denotes the number of feature maps in layer $l$. Also, a gating vector $g_i \in \mathbb{R}^{F_g}$ is utilized in order to determine the focus area per pixel. To achieve greater performance the attention coefficient is derived by leveraging the additive attention:

$$Q_\alpha^l = \psi^T(\varrho(W_x^T x_i^l + W_g^T g_i + \beta_g)) + \beta_\psi, \quad a_i^l = \sigma(Q_\alpha^l) \quad (4)$$

where $\sigma$ corresponds to the sigmoid activation function, $W_x \in \mathbb{R}^{F_l \times F_{int}}$ and $W_g \in \mathbb{R}^{F_g \times F_{int}}$ are linear transformations that are calculated by utilizing channel-wise 1×1×1 convolutions for the input tensors, and, lastly, $\beta_g \in \mathbb{R}^{F_{int}}$ and $\beta_\psi \in \mathbb{R}$ denote the bias.

### 2.2. Few-shot learning for segmentation refinement

As shown in Fig. 1b, we hereby propose a dynamic rectification scheme that leverages expert users' feedback on a small part of the data in order to improve the overall performance of the aforementioned R2AU-Net. The proposed retraining strategy dynamically updates the weights of the model, so that (a) the refined incoming samples are trusted as much as possible, while simultaneously (b) a minimal degradation of the already gained knowledge is achieved. To this end, let us denote $p_{ij}$ the soft label value of a pixel that is located in position $(i, j)$ of a given image. Then, for each input $n$ we calculate the average image confidence score $I_n$, defined as:

$$I_n = \frac{1}{\sum_{\forall i,j} \zeta_{ij}} \cdot \sum_{i=1}^{R} \sum_{j=1}^{C} \zeta_{ij} \cdot p_{ij} \quad (5)$$

where C and R correspond to the image's columns and rows respectively. In parallel, $\zeta_{ij} \in \{0, 1\}$ equals 1 when $p_{ij} > \vartheta$ and 0 otherwise, where $\vartheta$ is the detection acceptance threshold, which is set to 0.5. As such, the confidence score considers only the cracked regions over the image $n$, provided by the deep classifier. Subsequently, we rank the images according to $I_n$ scores. The 5% of the lower ranked images are provided to an engineer expert, who rectifies the model's segmentation outputs. Lastly, the refined few-shot annotated data are fed back to the network for updating the model's weights.

## 3. EXPERIMENTAL SETUP AND RESULTS

### 3.1. Dataset description

For the training procedure five datasets, consisting of 4,717 images in total, that depict crack defects, were utilized (see Table 1). During the data preprocessing step, the RGB data were resized to a resolution of 256×256 pixels. Lastly, 80% of the data was used for training the models (3,774 images), while the rest 20% was used for validation (943 images).

| Set | Name | Number of RGB samples |
|-----|------|----------------------|
| | CFD | 118 |
| | CRACK500 | 3,363 |
| **Train and** | Cracktree200 | 206 |
| **Validation** | DeepCrack | 521 |
| | GAPS384 | 509 |
| | **Total** | **4,717** |
| **Test** | **CrackMap** | **120 – 6** |

**Table 1**. Utilized datasets for training and evaluation tasks.

For the data collection process of the CrackMap dataset that constitutes the test dataset for the current study, a GoPro HERO9 Black was used. During the data acquisition process, the optical sensor was mounted on an inspection vehicle (see Fig. 2a). It is emphasized that the acquired data are RGB images with an aspect ratio of 4:3 and, in particular, with a pixel resolution of 5,184×3,888 (see Fig. 2b). Moreover, the RGB sensor was set to shoot at a high frame rate and, more specifically, at 50 frames per second in order to ensure sufficient data acquisition regarding both the positive (road surface with cracks) and negative (non-deteriorated road surface) events.

It is also highlighted that in order to deal with the severe class imbalance problem from the acquired RGB data, image patches with a resolution of 256×256 were manually extracted and then annotated. The patches were segmented and verified by engineer experts, within the framework of the H2020 HERON project [21]. As presented in Table 1, the CrackMap dataset contains 120 annotated images with a resolution of 256×256. We evaluate the comparative models that perform the crack segmentation task on the CrackMap data, minus the 6 extracted images that correspond to the 5% of the lower ranked images and were eventually utilized for the refinement
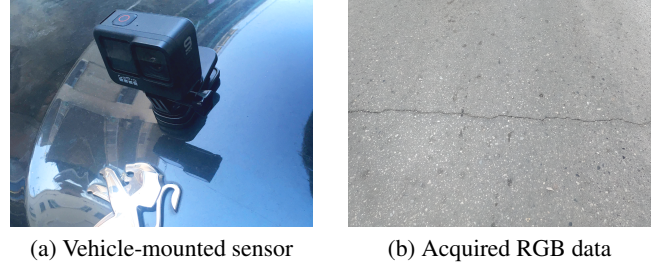


(a) Vehicle-mounted sensor        (b) Acquired RGB data

**Fig. 2**. Experimental setup for acquiring RGB road images.

process (see Section 2.2). CrackMap has been made available to the scientific community, for verifying the results and further research, at https://github.com/ikatsamenis/CrackMap.

### 3.2. Comparative algorithms and training configuration

The validation of the proposed methodology for the crack segmentation task is based on examining its performance against other state-of-the-art approaches that perform crack recognition and precise localization in a different way. In particular, we compare the proposed static and dynamically refined R2AU-Net models on the CrackMap dataset (minus the 6 extracted images) with the following segmentation algorithms: (i) U-Net [12], (ii) V-Net [22], (iii) ResU-Net [16], (iv) R2U-Net [17], (v) Attention U-Net [18], and (vi) ResUNet-a [23].

The aforementioned deep models were developed using Keras and TensorFlow libraries in Python. In parallel, they were trained and evaluated on NVIDIA Tesla T4 GPU provided by Google Colab. We trained the neural networks for 100 epochs with early stopping criteria set to 10 epochs in order to avoid overfitting, using mini-batches of size 8. The training processes started from scratch by randomly initializing the networks' weights. The models are trained end-to-end using the Adam algorithm to optimize the dice loss function, in order to deal with class imbalance problems. It is noted in parallel that the optimizer is set to its default parameters ($\beta_1 = 0.9$ and $\beta_2 = 0.999$), with an initial learning rate of $10^{-3}$ that is decayed by a factor of 10 each time there was no improvement in the validation loss, for 5 consecutive epochs. Lastly, for the rectification mechanism, we fine-tune the proposed R2AU-Net by retraining it for 5 epochs, with a learning rate reduced by a factor of 10, to avoid damaging its weights.

### 3.3. Experiments and comparisons

Fig. 3 depicts a visual comparison of the output masks generated by our method and the aforementioned comparative models. For a quantitative analysis of the experimental results, the performance of the implemented models is evaluated in terms of the Dice coefficient and IoU metric. In particular, we compute the aforementioned metrics for every RGB image of the CrackMap dataset (minus the 6 extracted images) and, thereby, we report the average values across all 114
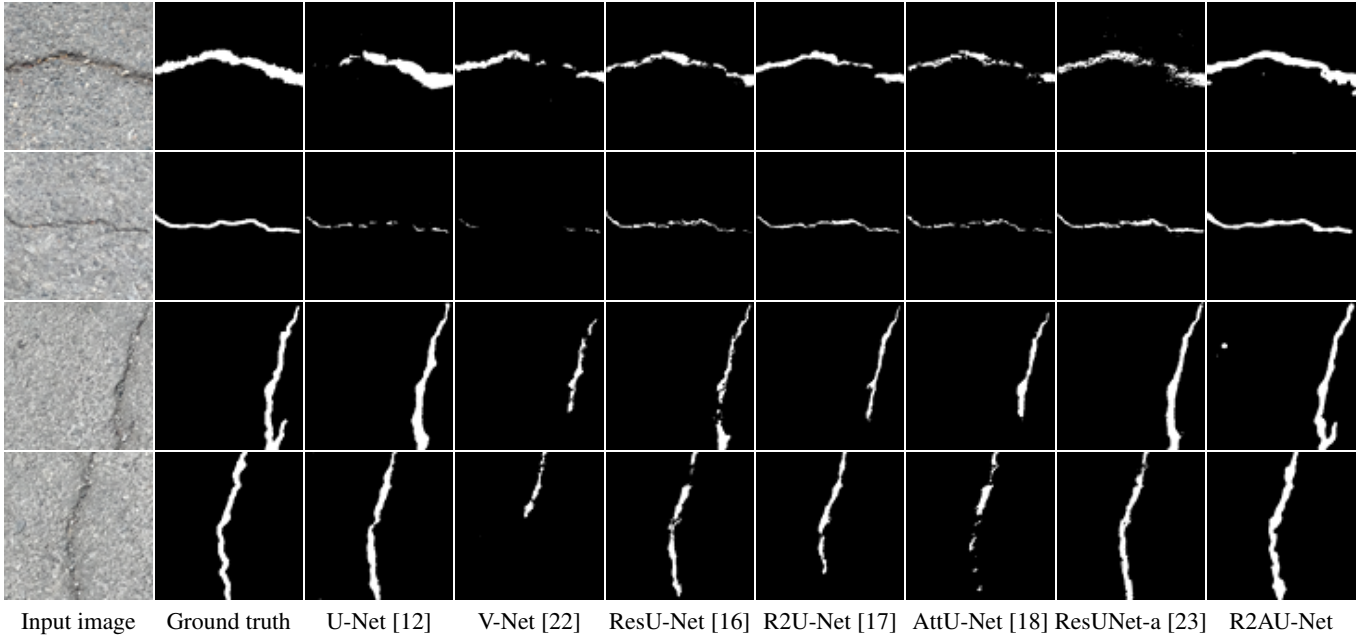
| Input image | Ground truth | U-Net [12] | V-Net [22] | ResU-Net [16] | R2U-Net [17] | AttU-Net [18] | ResUNet-a [23] | R2AU-Net |

**Fig. 3**. Visual comparison of the static deep models' segmentation outputs.

images of the test set with a confidence level of 95%. As can be observed in Table 2 the proposed R2AU-Net outperforms the various state-of-the-art algorithms by at least 1.69% and 1.89% in terms of Dice and IoU scores respectively.

| Model | Avg. Dice | Avg. IoU |
|---|---|---|
| U-Net [12] | 49.73% $\pm$ 3.75% | 35.54% $\pm$ 3.36% |
| V-Net [22] | 38.04% $\pm$ 4.65% | 26.57% $\pm$ 3.68% |
| ResU-Net [16] | 60.37% $\pm$ 2.50% | 44.53% $\pm$ 2.48% |
| R2U-Net [17] | 70.74% $\pm$ 1.59% | 55.39% $\pm$ 1.84% |
| AttU-Net [18] | 52.64% $\pm$ 3.70% | 38.06% $\pm$ 3.21% |
| ResUNet-a [23] | 63.28% $\pm$ 2.47% | 47.60% $\pm$ 2.49% |
| **R2AU-Net** | **72.43% $\pm$ 1.36%** | **57.28% $\pm$ 1.61%** |
| **FS R2AU-Net** | **77.06% $\pm$ 1.17%** | **63.11% $\pm$ 1.50%** |

**Table 2**. Performance evaluation and comparisons.

Finally, Fig. 4 illustrates indicative segmentation outputs of the R2AU-Net before and after applying the proposed few-shot rectification mechanism. As shown in Table 2, the rectified R2AU-Net model demonstrated increased performance of 4.63% and 5.83% in terms of Dice and IoU respectively, after the few-shot refinement procedure. To investigate whether this improvement is statistically significant, we exploit the Wilcoxon signed-rank test on the obtained scores of the two models, which is a nonparametric statistical test that compares two paired groups [24]. The obtained $p$-values for both metrics are lower than .001 and, thus, we can reject the null hypothesis, which entails that there is a statistically significant difference in the comparative results of the R2AU-Net, before and after the proposed few-shot refinement process.
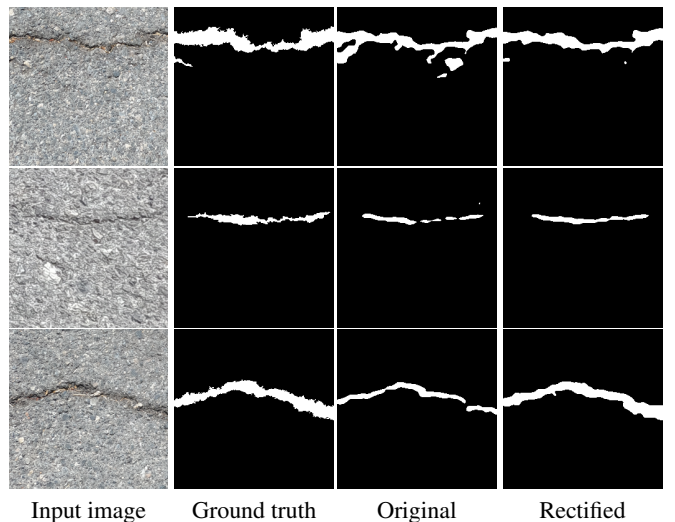


| Input image | Ground truth | Original | Rectified |

**Fig. 4**. Output masks produced by the original R2AU-Net against the proposed rectified few-shot learning paradigm.

## 4. CONCLUSION

This paper presents a few-shot learning strategy for road crack segmentation. The scheme is based on R2AU-Net which exploits recurrent residual and attention mechanisms to capture richer global context information and local semantic features. Also, the adopted few-shot refinement process, through which the network weights are dynamically updated as a few incoming rectified samples are fed into the algorithm, led to state-of-the-art performance on a new publicly available dataset.

# 5. REFERENCES

[1] E. Zaloshnja and T. R. Miller, "Cost of crashes related to road conditions, united states, 2006," in *Annals of Advances in Automotive Medicine/Annual Scientific Conference*. Association for the Advancement of Automotive Medicine, 2009, vol. 53, p. 141.

[2] C. Torres-Machí et al., "Sustainable pavement management: Integrating economic, technical, and environmental aspects in decision making," *Transportation Research Record*, vol. 2523, no. 1, pp. 56–63, 2015.

[3] H. Oliveira and P. L. Correia, "Automatic road crack segmentation using entropy and image dynamic thresholding," in *2009 17th European Signal Processing Conference*, 2009, pp. 622–626.

[4] H. Zhao, G. Qin, and X. Wang, "Improvement of canny algorithm based on pavement edge detection," in *2010 3rd International Congress on Image and Signal Processing*, 2010, vol. 2, pp. 964–967.

[5] N. Tanaka and K. Uematsu, "A crack detection method in road surface images using morphology.," *MVA*, vol. 98, pp. 17–19, 1998.

[6] S. L. Lau et al., "Automated pavement crack segmentation using u-net-based convolutional neural network," *IEEE Access*, vol. 8, pp. 114892–114899, 2020.

[7] T. U. Ahmed et al., "An integrated cnn-rnn framework to assess road crack," in *2019 22nd International Conference on Computer and Information Technology (ICCIT)*, 2019, pp. 1–6.

[8] N. Ogawa et al., "Distress level classification of road infrastructures via cnn generating attention map," in *2020 IEEE 2nd Global Conference on Life Sciences and Technologies (LifeTech)*, 2020, pp. 97–98.

[9] A. K. Pandey et al., "Convolution neural networks for pothole detection of critical road infrastructure," *Comp. and Electrical Engineering*, vol. 99, pp. 107725, 2022.

[10] I. Katsamenis et al., "Simultaneous precise localization and classification of metal rust defects for robotic-driven maintenance and prefabrication using residual attention u-net," *Automation in Construction*, vol. 137, pp. 104182, 2022.

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

[12] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Cham, 2015, pp. 234–241, Springer International Publishing.

[13] H. Huang et al., "Unet 3+: A full-scale connected unet for medical image segmentation," in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 1055–1059.

[14] A. Voulodimos et al., "A few-shot u-net deep learning model for covid-19 infected area segmentation in ct images," *Sensors*, vol. 21, no. 6, 2021.

[15] M. D. Jenkins et al., "A deep convolutional neural network for semantic pixel-wise segmentation of road and pavement surface cracks," in *2018 26th European Signal Processing Conference (EUSIPCO)*. IEEE, 2018, pp. 2120–2124.

[16] Z. Zhang, Q. Liu, and Y. Wang, "Road extraction by deep residual u-net," *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 5, pp. 749–753, 2018.

[17] M. Z. Alom et al., "Recurrent residual convolutional neural network based on u-net (r2u-net) for medical image segmentation," *arXiv preprint:1802.06955*, 2018.

[18] O. Oktay et al., "Attention u-net: Learning where to look for the pancreas," *arXiv preprint:1804.03999*, 2018.

[19] J. König et al., "A convolutional neural network for pavement surface crack segmentation using residual connections and attention gating," in *2019 IEEE International Conference on Image Processing (ICIP)*, 2019, pp. 1460–1464.

[20] M. Liang and X. Hu, "Recurrent convolutional neural network for object recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3367–3375.

[21] I. Katsamenis et al., "Robotic maintenance of road infrastructures: The heron project," in *Proceedings of the 15th International Conference on PErvasive Technologies Related to Assistive Environments*, New York, NY, USA, 2022, PETRA '22, p. 628–635, Association for Computing Machinery.

[22] F. Milletari, N. Navab, and S. A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 fourth international conference on 3D vision (3DV)*. Ieee, 2016, pp. 565–571.

[23] F. I. Diakogiannis et al., "Resunet-a: A deep learning framework for semantic segmentation of remotely sensed data," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 162, pp. 94–114, 2020.

[24] F. Wilcoxon, *Individual comparisons by ranking methods*, Springer, 1992.