

1 **Integrated APC-GAN and AttuNet framework for automated pavement crack
2 pixel-level segmentation:
3 A new solution to small training datasets**

5 **Tianjie Zhang**

6 Department of Computer Science
7 Boise State University

8 Boise, ID

9 tjzhang@u.boisestate.edu

10 **ABSTRACT**

11 Pavement crack segmentation using deep learning methods can improve crack segmentation accuracy, but
12 in many cases, the training dataset is lacking or uneven, making it insufficient to train an accurate
13 segmentation model. In this work, an integrated APC-GAN and AttuNet framework is proposed as an
14 automated pavement surface crack pixel-level segmentation strategy for small training datasets. First, an
15 automated pavement cracks generative adversarial network (APC-GAN) is designed for the pavement
16 cracks data as an image augmentation method, which is modified and improved from a traditional deep
17 convolutional generative adversarial network (DCGAN). Then, a novel pixel-level semantic segmentation
18 structure, Attunet, is proposed by introducing the attention module into the convolutional network work
19 structure. Another AttuNet-min is modified by replacing the max pooling layer and activation function in
20 AttuNet. In order to assess the performance of APC-GAN and AttuNet framework, an open-source dataset
21 DeepCrack is used, which only contains 300 training images. The results show that APC-GAN can produce
22 clearer and more distinct pavement images than DCGAN and more diversity than the traditional
23 augmentation method. The AttuNet model with APC-GAN can reach higher accuracy in evaluation metrics
24 than other augmentation methods. As for the segmentation model comparison, APC-GAN and AttuNet
25 framework gain the highest value in recall, F1 score, mean Intersection over Union (mIoU) and mean pixel
26 accuracy (mPA) among all models, including U-Net, DeepLabv3, FCN and LRASPP, while the AttuNet-
27 min gain the highest mean precision.

28
29
30 **Keywords:** Crack segmentation, GAN, attention module, deep learning

32 **1. INTRODUCTION**

33
34 Crack has become one of the primary defects in pavement, which seriously affects the service life of
35 pavement [1]. The traditional crack detection methods like counting cracks manually are labor-intensive
36 and time-consuming [2]. The automation of crack detection and segmentation has become the focus of
37 research in civil engineering and highway agencies[3]. Researchers have proposed a series of automated
38 detection methods for pavement crack detection. Image processing and machine learning are typical and
39 traditional methods in crack detection and extraction [4, 5]. The main idea of popular machine learning-
40 based crack inspection methods is to learn and summarize the characteristics of cracks and then build a
41 model to make predictions. Support Vector Machine (SVM) [6] and K-means algorithm [7] have already
42 been used in the classification and segmentation of crack images. A crack detection approach based on
43 Local Binary Patterns (LBP) with SVM proposed by Cheng [6], can extract the LBP feature of cracks and
44 make a segmentation from each frame of the video taken from the road. Ai, D. et al. [8] used multi-
45 neighborhood information to segment cracks and resulting in an F1 score of 0.8. The accuracy of the
46 algorithm proposed by Kaddah, W. et al. [9] is about 75%, which is based on the improved minimum path
47 method, an image processing method, to segment pavement cracks. However, image processing and
48 traditional machine learning heavily depend on an engineer's knowledge, which may limit the overall
49 performance. Deep learning is becoming one of the most advanced pixel-level target detection methods in
50 road condition inspection as it learns from large-scale data and requires little human involvement during
51 training. Convolution Neural Network (CNN), a deep learning method, has been gradually utilized in road
52 crack detection and segmentation [10, 11]. For example, Bang, S et al. [12] used ResNet-152 to classify
53 cracks and got the results with a precision of 77.68% and a recall of 71.98%. Cao Vu dung et al. [13] used
54 a Full Convolution Neural Network (FCN) to detect and segment cracks and got an average accuracy of
55 90%. Yahui Liu[14] proposed a DeepCrack dataset for crack segmentation and six DeepCrack Structures
56 which mainly consisted of the extended FCN and the Deeply-Supervised Nets (DSN), and it showed a
57 comparable result when compared with typical segmenting methods like AutoCrack and SegNet. Liu J W
58 et al. [15] proposed a two-step pavement crack detection and segmentation method by firstly using a YOLO
59 v3 to locate the crack area and then applying a modified U-Net model to segment the crack from this area.
60 Chengjia Han[16] proposed a U-Net-based CNN model, CrackW-Net, by adding a skip-level round-trip
61 sampling block to segment the pavement images from the Crack500 dataset and a self-built dataset, and
62 shows a good result.

63

64 Although deep learning is the most advanced pixel-level segmentation method, it requires a large amount
65 and a wide diversity of annotated data to train the network[17]. A small training dataset may cause the
66 neural network overfitting and lousy performance in robustness. However, the cost of obtaining a large

67 number of training samples is very high [18]. To address this issue, some researchers developed deep
68 learning architectures which can work with very few training images but still yields precise segmentations.
69 The main idea of U-Net[19] is to replace pooling operators with upsampling operators to increase the output
70 resolution. Also, it combines the high-resolution features with the upsampled output to learn more precise
71 information based on a small dataset. The attention module is popular in nature language process (NLP)
72 and now it has started to be applied in the computer vision area[20]. It can improve the sensitivity and
73 efficiency of the network to get rid of large amount of data[21]. For example, Wenjun Wang[3] proposed
74 a pyramid attention network that used pre-trained DenseNet121 and a feature pyramid attention module. It
75 was tested on the Crack500 and MCD datasets and achieved an IoU of 0.6235. Xuezhi Xiang[22] proposed
76 a pavement crack segmentation network based on the BAM attention module and it got an mPA of 0.831,
77 much higher compared to other networks like SegNet and CrackForest.

78

79 Another alternative method is data augmentation. The most common strategy for data augmentation is the
80 traditional augmentation methods like image random crop, image flip and adding noise[23]. In 2014,
81 Goodfellow[24] proposed the concept of generative adversarial networks (GANs), which can produce real-
82 like images through a battle between a generator and a discriminator. Alec Radford[25] proposed deep
83 convolutional generative adversarial networks (DCGANs) based on the conception of GAN, and it showed
84 good representations of images. Because of using the convolution structure, DCGAN is popular in computer
85 vision and has been applied in many areas including pavement crack data augmentation. For example, Lili
86 Pei[26] used a variational autoencoder (VAE) to encode crack images and the results from VAE were input
87 to the DCGAN model to generate the fake images. Boqiang Xu[17] reached quite a high identification
88 accuracy in pavement crack classification tasks based on DCGAN and VGG16. However, there are still
89 some problems with DCGAN. For example, the original DCGAN structure is more suitable for small-size
90 images like an image with a resolution of 32 * 32 pixels. In pavement crack segmentation tasks, the
91 resolution of images usually is larger than 256 * 256 pixels. Another problem is that the discriminator in
92 DCGAN studies too fast which would lead to the loss of the discriminator to zero very rapidly, while the
93 generator does not study very well.

94

95 This paper proposes a framework for pavement crack segmentation which contains an automated pavement
96 crack generative adversarial network (APC-GAN) and a new pixel-level pavement crack segmentation
97 network, AttuNet. The APC-GAN is modified from DCGAN and designed to improve the generated road
98 image quality. The kernel size of APC-GAN is enlarged in the generator to capture more information.
99 Moreover, the convolutional layers are increased in both generator and discriminator to produce sharper
100 images. Gaussian noise is added at the top of the discriminator to slow down its convergence speed. The

101 AttuNet is modified from U-Net. An attention module is introduced to this structure which can extract
102 cracks' features by fusing different channel information from different layers. Batch normalization is used
103 both in APC-GAN and AttuNet to accelerate training. The proposed AttuNet combines the advantages of
104 U-Net and the attention module. Another AttuNet-min is presented by replacing the max pooling layer with
105 the min pooling layer to make the network focus more on the crack in the road image. In this paper, an
106 open-source dataset is used to verify the segmentation accuracy of the proposed method. The experimental
107 results show that the APC-GAN provides more distinct and diverse images than DCGAN and the traditional
108 augmentation method. The proposed framework is compared with four classic CNN models including U-
109 Net, DeepLab-resnet50[27], FCN-resnet50[28], and LRASPP_mobilenet_v3-large[29], and it shows higher
110 accuracy in crack segmentation. The organization of this study is outlined as follows. In METHODOLOGY,
111 the basic theory and structure of APC-GAN and AttuNet are presented. In RESULTS,
112 the proposed framework is applied to the crack dataset and compares the performance with various CNN
113 models and it shows the pixel-level segmentation results of the cracks.

114

115 **2. METHODOLOGY**

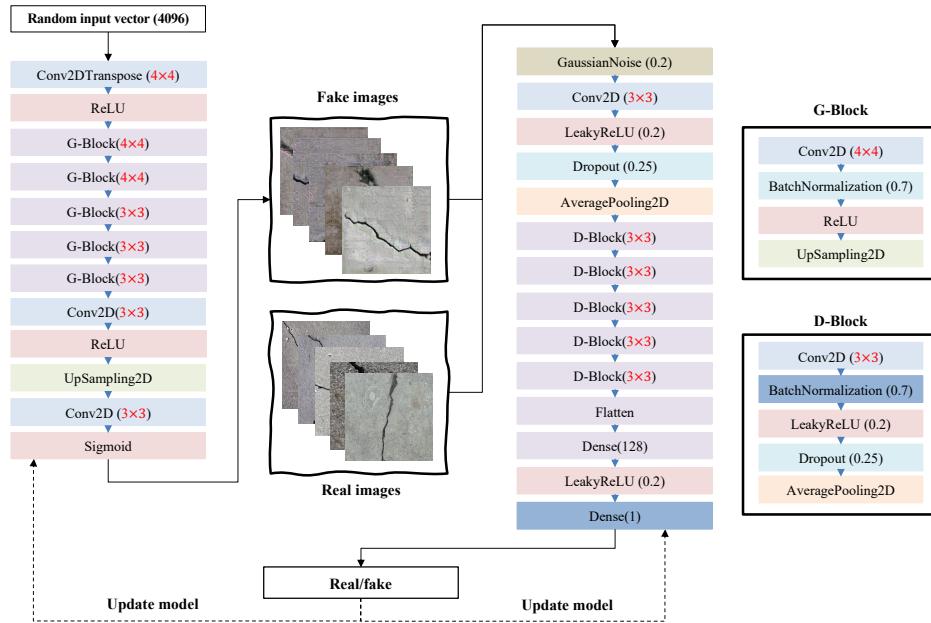
116

117 **2.1 APC-GAN**

118

119 APC-GAN is designed for pavement crack segmentation tasks according to the shortages of DCGAN. The
120 structure of APC-GAN is shown in Figure 1. The input is a random vector with a length of 4096. In the
121 Generator part, a G-Block (4×4) is comprised of a convolution layer with a kernel size of 4×4 , followed by
122 batch normalization with momentum 0.7, an activation function Rectified Linear Unit (ReLU) and an up-
123 sampling layer. In the Discriminator part, a D-Block (3×3) is comprised of a convolutional layer whose
124 kernel size is 3×3 , followed by batch normalization with momentum 0.7, an activation function leaky
125 Rectified Linear Unit (leakyReLU), a Dropout layer with parameter 0.25 and an average pooling layer. The
126 convolution layer can process the image to produce a set of feature maps while the activation functions
127 used in this APC-GAN model include ReLU, leakyReLU and Sigmoid, which can make the network learn
128 a non-linear task. The average pooling in the discriminator is utilized to translate invariance and reduce the
129 parameter size of the networks.

130



131

132

133

Figure 1. A structure diagram of the proposed APC-GAN model.

134

135 Typically, the DCGAN only works well at images with a low resolution, like 32×32 pixels or 64×64 pixels.

136 However, in a pavement crack segmentation task, a generated image with a resolution of 256×256 pixels is

137 required. Another problem in DCGAN is that the discriminator studies too fast leading to the loss of the

138 discriminator to 0 very rapidly. It would lead to the situation that the loss cannot be used to update the

139 generator although it did not learn well. In order to make the APC-GAN architecture better results, some

140 modifications are made. The contributions of this work can be summarized as follows:

141

142 1) Large kernel size is used. The kernel size is increased to 4×4 in the generator and to 3×3 in the
143 discriminator. For the generator, a large kernel at the top convolutional layers could cover more area and
144 thus, could capture more information, which could maintain the smoothness of the image. For the
145 discriminator, a small kernel may cause the discriminator loss rapidly approaches zero while a larger kernel
146 size can ease this situation.

147

148 2) The number of convolutional layers is increased in APC-GAN compared to the original DCGAN. A
149 small number of convolution operators, especially in the generator, would make the produced images very
150 blurry while more layers can help capture additional information which can eventually add sharpness to the
151 final produced images.

152

153 3) A batch normalization layer is followed by the convolutional layer. Batch normalization acts as a
154 regularizer which can reduce the accelerating training and improve the generated image quality. The function
155 can be described in Eq. (1).

156

$$157 \quad y = \frac{x - E[x]}{\sqrt{Var[x] + \epsilon}} * \gamma + \beta \quad (1)$$

158

159 Where γ and β are learnable parameter vectors ($\gamma = 1, \beta = 0$), ϵ is a value added to the denominator for
160 numerical stability ($\epsilon = 1e - 5$). $E[x]$ and $Var[x]$ are the mean and variance of input x .

161

162 4) A Gaussian noise layer is added as the first layer of the discriminator. It can prevent the discriminator
163 from studying too quickly.

164

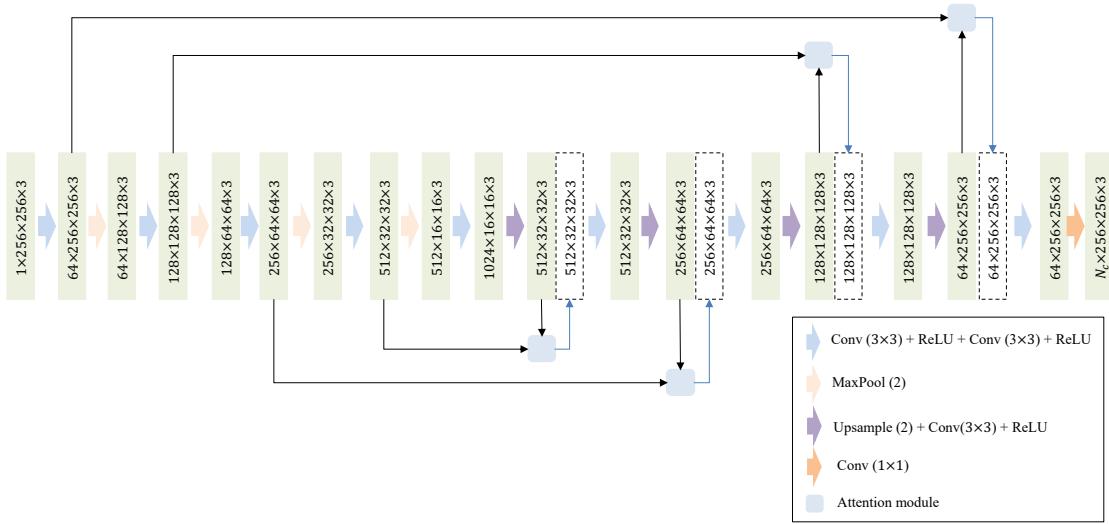
165 When applied the APC-GAN model in the data augmentation, the initial learning rate is set to 0.0001, and
166 adjusted during training, where the learning rate would decrease by 15% for each 10000 steps. Binary Cross
167 Entropy Loss (BCELoss) is used as the loss function and Adam was utilized as the optimizer to update the
168 network. The batch size of the dataset is set to 16.

169

170 **2.2 AttuNet**

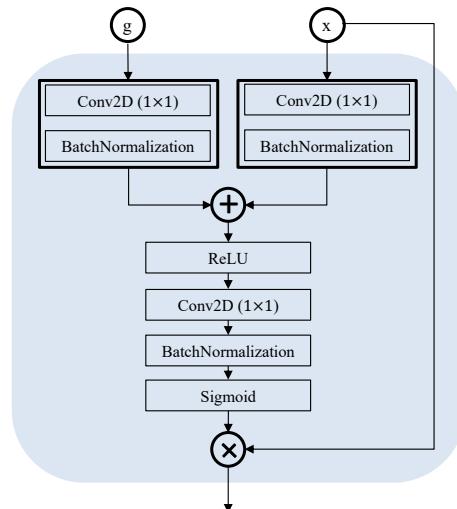
171

172 An AttuNet is proposed as a novel crack segmentation approach in this work. Figure 2 shows the main
 173 architecture of the proposed structure, in which the details of each operation are presented. This structure
 174 is modified from U-Net. The output is an N_c channel map of the probabilities where N_c is the number of
 175 classes ($N_c=1$ in this work). Four attention modules are introduced to connect different layers. Each
 176 attention module has two inputs, one from the current layer and the other from a previous layer. Then, the
 177 output from the attention module would be concatenated with the current layer. The attention module filters
 178 the features from different layers rather than connecting the features directly. Each convolutional operation
 179 follows a batch normalization to standardize the inputs and stabilize the learning procedure.
 180



190 3. The attention module gets two inputs g and x from the former layers. Each input is processed through a
191 convolution layer with a kernel size of 1×1 and batch normalization. They are added together where the
192 module can fuse the features under the two different scales. Then, the fused features are processed with
193 ReLU activation, a convolution operation with a kernel size 1×1 , a batch normalization and a sigmoid
194 function. At last, the result from the attention module would concatenate with the input x . By doing this,
195 the attention module can fuse the different features from different scale layers to improve the consistency
196 of the feature map and thus improve the model performance, as well as decrease the data usage. An attention
197 module can refine the pavement crack to make it effectively guide the AttuNet training. Moreover, the
198 convolution layer in the attention module can extract cracks' features by fusing different channel
199 information from different layers.

200



201

202

203 Figure 3. The structure of the attention module.

204

205 2) Each convolution layer is followed by a batch normalization layer, which can standardize the parameters
206 and speed up the training procedure.

207

208 3) Root Mean Squared Propagation (RMSProp) with a momentum of 0.9 was utilized as the optimizer to
209 update the network. BCEWithLogitsLoss is used as the loss function, which combines a sigmoid layer and
210 the Binary Cross Entropy in one single class. This loss is more numerically stable than a plain Sigmoid
211 followed by a BCELoss.

212

213 The proposed AttuNet structure combines the attention modules within the CNN network which can make
214 the network works well with a small dataset. This is because the attention module can fuse the features from
215 different layers which can make the model study faster.

216

217 For the crack segmentation task, another version of AttuNet, called AttuNet-min, is designed in this work.
218 In this version, the max pooling layer is replaced by the min pooling layer. This is because the crack pixels
219 always have a relatively small value in an image, using a min pooling layer can keep the crack information
220 accurate when downsizing the images. At the same time, the ReLU is replaced by the LogSigmoid function
221 as shown in Eq. (2).

222

$$223 \quad \text{LogSigmoid}(x) = \log\left(\frac{1}{1+\exp(-x)}\right) \quad (2)$$

224

225 When the input x gets smaller, the absolute value of output from LogSigmoid would be more significant.
226 Thus, the LogSigmoid function would give more weight to the small pixels while the ReLU function pays
227 more attention to the brighter part.

228

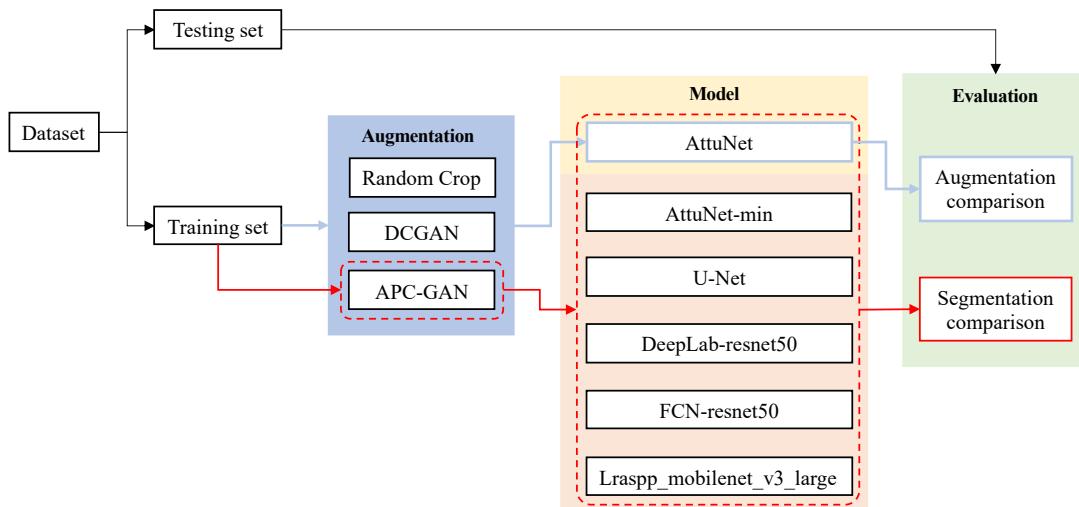
229 **2.3 Overall procedure**

230

231 An overall procedure to evaluate and validate the usage of the proposed framework is shown in Figure 4.
232 A DeepCrack dataset is utilized in this work to test the performance of the CNN networks. The DeepCrack

233 dataset is an open-source dataset published on GitHub (<https://github.com/yhlleo/DeepCrack>). This dataset
234 consists of 537 RGB crack images with manually annotated segmentations. The image has a resolution of
235 544 * 384 pixels. The images are divided into two subsets: 300 images for training and 237 images for
236 testing. As we can see, the number of images is relatively low for deep learning training.

237



238

239

240

Figure 4. The overall procedure of experiments.

241

242 The blue line in Figure 4 shows the procedure for comparing and evaluating the augmentation methods
243 including Random Crop, DCGAN and APC-GAN. Each augmentation method in this work generated 300
244 images in total. The generated images would be annotated manually and then combined with the DeepCrack
245 training dataset. Then, the proposed AttuNet would be trained on the augmented dataset. The data
246 augmentation methods are compared and assessed by the results from the segmentation.

247

248 The red line in Figure 4 shows the process for comparing and evaluating different segmentation models
249 including AttuNet, AttuNet-min, U-Net, DeepLab-resnet50, FCN-resnet50, and LRASPP_mobilenet_v3-

250 large. In order to evaluate the performance of the proposed AttuNet and AttuNet-min in the pavement cracks
251 segmentation work, other four segmentation methods including U-Net, DeepLab_resnet50, FCN-resnet50
252 and lraspp_mobilenet_v3_large are introduced and compared. We fine-tune these four methods. The
253 number of classes is set to 1. In the training procedure, the initial learning rate is set to 0.00001 and the
254 batch size of the dataset is set to 16. The training epochs are set to 300. BCEWithLogitsLoss is used as the
255 loss function and the RMSProp is utilized as the optimizer to update the network parameters. Before
256 importing to the deep learning structure, all the images including the crack image and its label, would be
257 reshaped to 256 * 256 pixels.

258

259 The data augmentation methods and CNN models are implemented in Python and computed under the
260 following machine speculations: Windows 10, Intel(R) Core (TM) i9-10900X CPU, NVIDIA RTX A4000
261 with 16 GB memory, and 64GB RAM.

262

263 **2.4 Evaluation metrics**

264

265 Precision (P), Recall (R), F1 score (F1), Intersection over Union (IoU) and pixel accuracy (PA) are utilized
266 to evaluate the semantic segmentation results.

267

268 1) P can measure how accurate your predictions are. The precision can be calculated by Eq. (3) where TP
269 is true positive and FP is false positive.

270

271
$$P = \frac{TP}{TP+FP} \quad (3)$$

272

273 2) R suggests the level of sensitivity for prediction results. Recall can be calculated by Eq. (4) where FN is
274 false negative.

275

276

$$R = \frac{TP}{TP+FN} \quad (4)$$

277

278 3) F1 is defined based on the harmonic average of Precision and Recall. It can be calculated using Eq. (5).

279

280

$$F_1 = \frac{2PR}{P+R} \quad (5)$$

281

282 4) IoU measures the overlap between 2 areas. It is used to measure how much the predicted areas overlap
283 with the ground truth. IoU is calculated according to Eq. (6).

284

285

$$IoU = \sum_{i,j}^k \frac{p_{ii}}{p_{ij} + p_{ji} - p_{ii}} \quad (6)$$

286

287 Where p_{ij} represents the number of pixels belonging to class i but predicted as class j.

288 The mean Intersection over Union (mIoU) is calculated according to Eq. (7).

289

290

$$mIoU = \frac{1}{2} \sum_{i,j=2}^k \frac{p_{ii}}{p_{ij} + p_{ji} - p_{ii}} \quad (7)$$

291

292 5) PA is a semantic segmentation metric that denotes the percentage of pixels that are accurately classified
293 in the image. It can be calculated by Eq. (8).

294

295

$$PA = \sum_i^k \frac{p_{ii}}{t_i} \quad (8)$$

296

297 Where t_i is the total number of pixels that are labeled as class i .

298 Since there are two classes present in this work: crack and background, the mean pixel accuracy (mPA) is
299 calculated to represent the class average accuracy, shown in Eq. (9).

300

301 $mPA = \frac{1}{2} \sum_i^{k=2} \frac{p_{ii}}{t_i}$ (9)

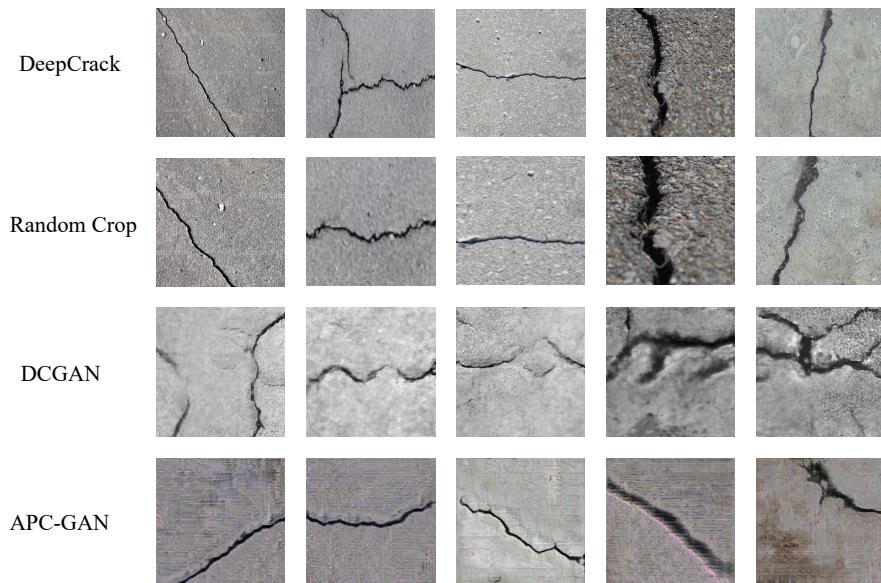
302

303 **3. RESULTS**

304

305 A traditional image augment method, random crop, and a DCGAN are used in this work to compare with
306 the proposed APC-GAN method. Some samples of the original image and the generated images from a
307 random crop, DCGAN, and APC-GAN are shown in Figure 5.

308



309

310

311 Figure 5. The raw images from DeepCrack and the generated images from the random crop, DCGAN and APC-
312 GAN.

313

314 As we can see from Figure 5., compared to the DCGAN, the images generated from APC-GAN are more
315 distinct and sharper. The images produced from the random crop are clear and more distinct than images
316 generated from GANs as the image is cropped from the original image directly. However, these images are
317 not as much diversity as the images produced by DCGAN and APC-GAN. The generated images are added
318 to the original training data respectively. And the augmented training dataset is used to train the AttuNet
319 model. The precision, recall, F1 score, mIoU, mPA are calculated and shown in Table 1.

320

321 TABLE 1. Comparison of different image augmentation methods using the same segmentation model.

Model	Data	Augmentation	P	R	F1	mIoU	mPA
AttuNet	DeepCrack	None	0.950	0.839	0.892	0.812	0.839
		APC-GAN	0.947	0.868	0.906	0.836	0.868
		DCGAN	0.949	0.851	0.897	0.822	0.851
		Random Crop	0.950	0.856	0.900	0.827	0.856

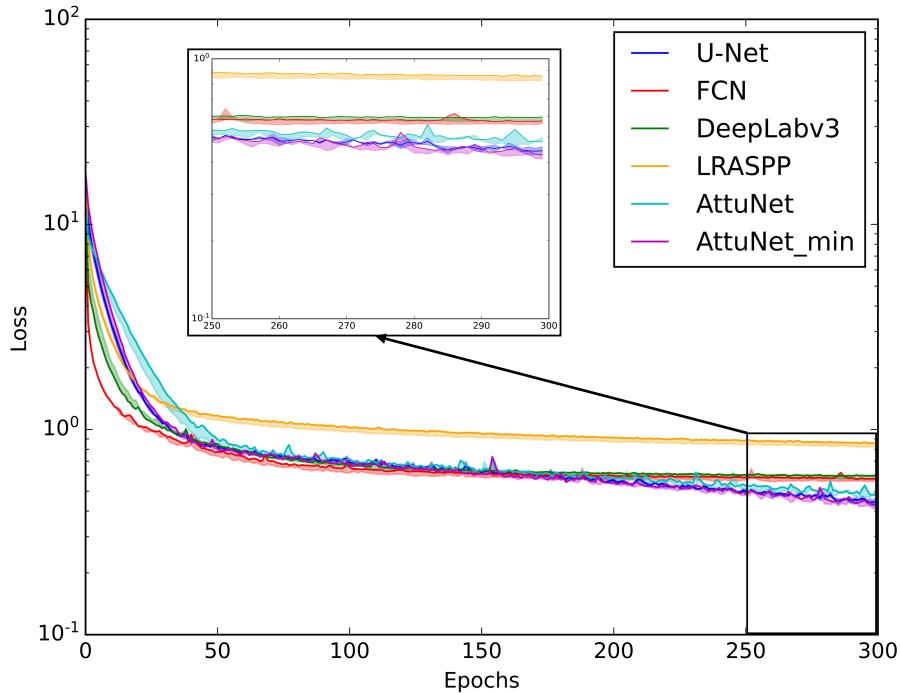
322

323 As we can see from TABLE 1, the dataset DeepCrack with an APC-GAN augmentation can make the
324 segmentation neural network gain the highest recall, F1 score, mIoU and mPA value among all the
325 augmentation methods used in this work. It means that by using the proposed APC-GAN, the accuracy of
326 the semantic segmentation method can be improved in a small dataset pavement crack detection task.

327

328 In order to evaluate that the proposed AttuNet and AttuNet-min have good performance in the small data
329 pavement cracks segmentation work. Other four segmentation methods including U-Net,
330 DeepLab_resnet50, FCN-resnet50 and lraspp_mobilenet_v3_large are introduced and compared. These
331 methods are trained on the DeepCrack train dataset with APC-GAN augmentation. The training procedure
332 is shown in Figure 6.

333



334

335

336

Figure 6. The Loss changes during the training procedure.

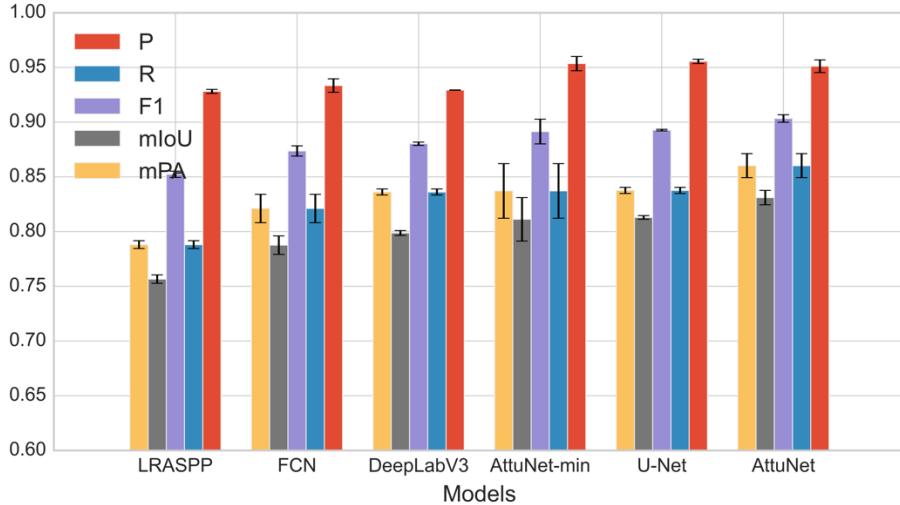
337

338 Figure 6 shows the training procedure where the loss of the model decreases with the epochs. And as we
 339 can see, the LRASPP does not decrease rapidly around 50 epochs as other models all decrease to around
 340 0.8 at 50 epochs. Compared to the FCN and DeepLabv3, the loss of U-Net, AttuNet and AttuNet-min
 341 continuously decreases after 200 epochs and finally reaches at about 0.4. Because these models are all using
 342 the same loss function BCEWithLogitsLoss, the loss can actually reflect the performance of the models. It
 343 indicates that the U-Net, AttuNet and AttuNet-min learn more during the training than the other three methods.

344

345 In order to evaluate and compare the performance between different deep learning segmentation models
 346 statistically, each model is trained and tested three times. The mean value and standard deviation are
 347 calculated. The evaluation metrics of each model are shown in Figure 7.

348



349

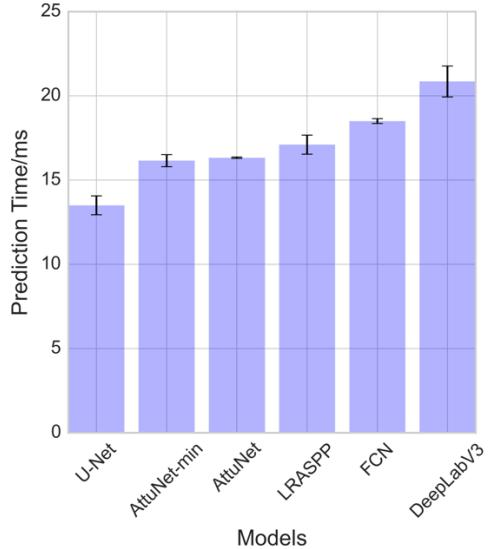
350

Figure 7. Comparison of different models on the test data.

351

352 The models in Figure 7 are ranked by the mIoU from lowest to highest. It shows that the AttuNet gets the
 353 highest mIoU (0.831) among all the models. It also gets the highest value in mPA, followed by AttuNet-
 354 min and U-Net, which means that the AttuNet has the highest percentage of predicting the pixels accurately
 355 as labeled in the image. The mIoU and mPA are the two most important indexes showing the segment
 356 ability of the method as both of them counted and compared each pixel. A higher value in mIoU and mPA
 357 means more pixels were classified accurately. The AttuNet-min gets the highest precision (0.96) among all
 358 CNN models. It means that 96% of the predicted cracks or background are labeled initially as cracks or
 359 background.

360



361

362

363

Figure 8. The comparison of the prediction time of each model.

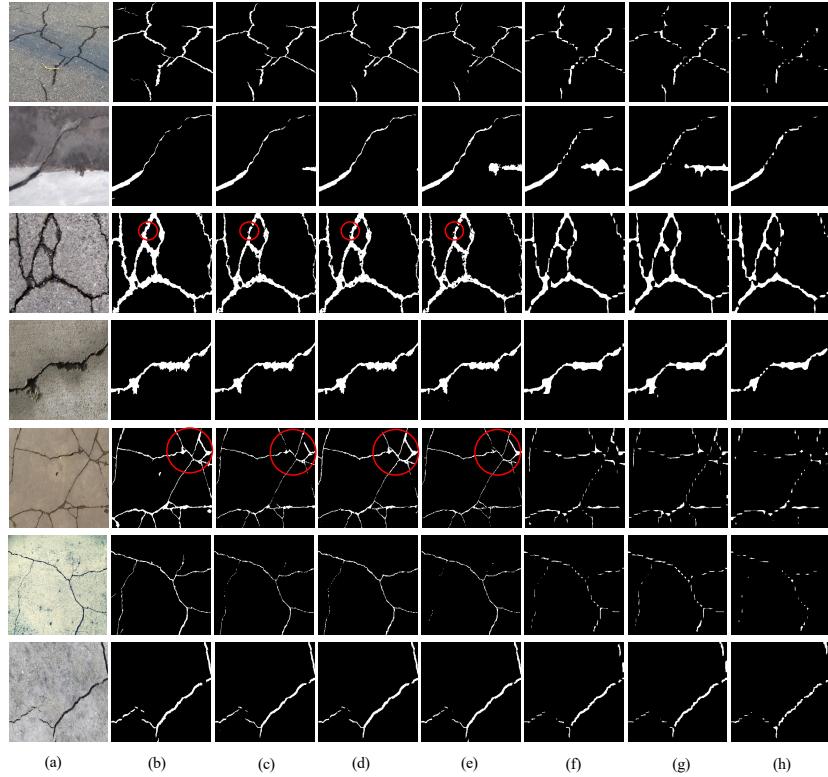
364

365 If we want to use the segmentation model in real-time crack detection work, the prediction time per image
 366 is an important factor. A faster prediction time of one model means this model is more suitable for real-
 367 time jobs. As shown in Figure 8, the models are ranked by the prediction time. As we can see, the U-Net
 368 model consumes the least time while the DeepLabv3-resnet50 consumes the largest time. The mean
 369 prediction time of AttuNet and AttuNet-min is 16.32 ms and 16.15 ms, respectively, which perform better
 370 than LRASPP, FCN and DeepLabV3.

371

372 Figure 9 shows some samples with cracks in various scenes and their segmentation results using different
 373 methods.

374



375

376

377 Figure 9. Several samples with cracks in various scenes and their segmentation results using different methods: (a)
 378 Original image (b) Ground Truth (c) AttuNet (d) AttuNet_min (e) U-Net (f) FCN_resnet50 (g) Deeplabv3_resnet50
 379 (h) LRASPP_mobilenet_v3_large

380

381 As shown in Figure 9, it is evident that the results from AttuNet, AttuNet_min and U-Net are more
 382 continuous and complete than the segmented results from FCN-resnet50, DeepLab-resnet50 and
 383 LRASPP_mobilenet_v3_large. The segmented part in the red circle shows that the cracks are segmented
 384 more entirely by AttuNet_min than by AttuNet and U-Net. It shows that the AttuNet-min has a good
 385 performance in the continuous of the cracks as the segmentation image is much closer to the ground truth.
 386 This is mainly because in AttuNet_min, the max pooling layer is replaced by a min pooling layer and the
 387 activation function logsigmoid focuses more on the small pixel value. Thus, the dark part in the real image
 388 would be paid more attention in AttuNet-min which makes the segmentation results continuously.

389

390 **4. CONCLUSIONS**

391

392 This paper proposed a novel pixel-level crack segmentation strategy for pavement crack inspection with a
393 small dataset. This situation is very common in road maintenance as the cost of obtaining a large number
394 of pavement top-view images and labeling these manually is very high. However, a small training dataset
395 may cause the neural network overfitting and bad model performance in robustness. The proposed
396 framework consists of APC-GAN and AttuNet, which can accurately segment images with a small dataset.

397

398 In this paper, only a total of 300 color images are used for training. The performance of APC-GAN is
399 evaluated and it shows a better ability to produce sharper contrast and more diverse images compared to
400 DCGAN and Random Crop. It improves the recall, F1 score, mIoU and mPA of the segmentation model.
401 The proposed AttuNet model combines the attention module and batch normalization layer with the CNN
402 network. It gets the testing mean IoU (0.831), which is higher than the classic CNN models including U-
403 Net, DeepLabv3, FCN and LRASPP. Compared with AttuNet, the AttuNet-min achieved a more continuous
404 segmentation result by applying the min pooling layer and LogSigmoid activation function.

405

406

407 REFERENCES

408

- 409 [1] W. Wang, M. Wang, H. Li, H. Zhao, K. Wang, C. He, J. Wang, S. Zheng, J. Chen, Pavement crack
410 image acquisition methods and crack extraction algorithms: A review, *Journal of Traffic and*
411 *Transportation Engineering (English Edition)*, 6 (2019) 535-556.
- 412 [2] Y. Yu, M. Rashidi, B. Samali, A.M. Yousefi, W. Wang, Multi-image-feature-based hierarchical
413 concrete crack identification framework using optimized SVM multi-classifiers and D-S fusion algorithm
414 for bridge structures, *Remote Sensing*, 13 (2021) 240.
- 415 [3] W. Wang, C. Su, Convolutional neural network-based pavement crack segmentation using pyramid
416 attention network, *IEEE Access*, 8 (2020) 206548-206558.
- 417 [4] W. Lin, Y. Sun, Q. Yang, Y. Lin, Real-time comprehensive image processing system for detecting
418 concrete bridges crack, *Computers and Concrete, An International Journal*, 23 (2019) 445-457.
- 419 [5] W. Cao, Q. Liu, Z. He, Review of pavement defect detection methods, *Ieee Access*, 8 (2020) 14531-
420 14544.
- 421 [6] C. Chen, H. Seo, C.H. Jun, Y. Zhao, Pavement crack detection and classification based on fusion
422 feature of LBP and PCA with SVM, *International Journal of Pavement Engineering*, (2021) 1-10.
- 423 [7] J. Huyan, W. Li, S. Tighe, R. Deng, S. Yan, Illumination compensation model with k-means
424 algorithm for detection of pavement surface cracks with shadow, *Journal of Computing in Civil*
425 *Engineering*, 34 (2020) 04019049.
- 426 [8] D. Ai, G. Jiang, L.S. Kei, C. Li, Automatic pixel-level pavement crack detection using information of
427 multi-scale neighborhoods, *IEEE Access*, 6 (2018) 24452-24463.
- 428 [9] W. Kaddah, M. Elbouz, Y. Ouerhani, V. Baltazart, M. Desthieux, A. Alfallou, Optimized minimal path
429 selection (OMPS) method for automatic and unsupervised crack segmentation within two-dimensional
430 pavement images, *The Visual Computer*, 35 (2019) 1293-1309.
- 431 [10] C. Chen, H. Seo, Y. Zhao, A novel pavement transverse cracks detection model using WT-CNN and
432 STFT-CNN for smartphone data analysis, *International Journal of Pavement Engineering*, (2021) 1-13.
- 433 [11] S.L. Lau, E.K. Chong, X. Yang, X. Wang, Automated pavement crack segmentation using u-net-
434 based convolutional neural network, *IEEE Access*, 8 (2020) 114892-114899.
- 435 [12] S. Bang, S. Park, H. Kim, H. Kim, Encoder-decoder network for pixel-level road crack detection in
436 black-box images, *Computer-Aided Civil and Infrastructure Engineering*, 34 (2019) 713-727.
- 437 [13] C.V. Dung, Autonomous concrete crack detection using deep fully convolutional neural network,
438 *Automation in Construction*, 99 (2019) 52-58.
- 439 [14] Y. Liu, J. Yao, X. Lu, R. Xie, L. Li, DeepCrack: A deep hierarchical feature learning architecture for
440 crack segmentation, *Neurocomputing*, 338 (2019) 139-153.
- 441 [15] J. Liu, X. Yang, S. Lau, X. Wang, S. Luo, V.C.S. Lee, L. Ding, Automated pavement crack detection
442 and segmentation based on two-step convolutional neural network, *Computer-Aided Civil and*
443 *Infrastructure Engineering*, 35 (2020) 1291-1305.
- 444 [16] C. Han, T. Ma, J. Huyan, X. Huang, Y. Zhang, CrackW-Net: A novel pavement crack image
445 segmentation convolutional neural network, *IEEE Transactions on Intelligent Transportation Systems*,
446 (2021).
- 447 [17] B. Xu, C. Liu, Pavement crack detection algorithm based on generative adversarial network and
448 convolutional neural network under small samples, *Measurement*, 196 (2022) 111219.
- 449 [18] Y. Zhang, K.V. Yuen, Crack detection using fusion features-based broad learning system and image
450 processing, *Computer-Aided Civil and Infrastructure Engineering*, 36 (2021) 1568-1584.
- 451 [19] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image
452 segmentation, *International Conference on Medical image computing and computer-assisted intervention*,
453 Springer, 2015, pp. 234-241.
- 454 [20] H. Wan, L. Gao, M. Su, Q. Sun, L. Huang, Attention-based convolutional neural network for
455 pavement crack detection, *Advances in Materials Science and Engineering*, 2021 (2021).

- 456 [21] O. Oktay, J. Schlemper, L.L. Folgoc, M. Lee, M. Heinrich, K. Misawa, K. Mori, S. McDonagh, N.Y.
457 Hammerla, B. Kainz, Attention u-net: Learning where to look for the pancreas, arXiv preprint
458 arXiv:1804.03999, (2018).
- 459 [22] X. Xiang, Y. Zhang, A. El Saddik, Pavement crack detection network based on pyramid structure
460 and attention mechanism, IET Image Processing, 14 (2020) 1580-1586.
- 461 [23] D. Mazzini, P. Napoletano, F. Piccoli, R. Schettini, A novel approach to data augmentation for
462 pavement distress segmentation, Computers in Industry, 121 (2020) 103225.
- 463 [24] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y.
464 Bengio, Generative adversarial nets, Advances in neural information processing systems, 27 (2014).
- 465 [25] A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional
466 generative adversarial networks, arXiv preprint arXiv:1511.06434, (2015).
- 467 [26] L. Pei, Z. Sun, L. Xiao, W. Li, J. Sun, H. Zhang, Virtual generation of pavement crack images based
468 on improved deep convolutional generative adversarial network, Engineering Applications of Artificial
469 Intelligence, 104 (2021) 104376.
- 470 [27] L.-C. Chen, G. Papandreou, F. Schroff, H. Adam, Rethinking atrous convolution for semantic image
471 segmentation, arXiv preprint arXiv:1706.05587, (2017).
- 472 [28] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation,
473 Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431-3440.
- 474 [29] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V.
475 Vasudevan, Searching for mobilenetv3, Proceedings of the IEEE/CVF international conference on
476 computer vision, 2019, pp. 1314-1324.

477