

BRIDGE CRACK DETECTION THROUGH UAV

Abstract:

This research paper focuses on crack segmentation, a critical task in structure investigation problems such as bridge inspection projects. The accuracy of the segmentation model is crucial in identifying potential regions of damage on the bridge surface from images captured by drones. The task is challenging due to the presence of noise and other objects such as moss on cracks, tile lines, etc., which can be wrongly identified as cracks. In this study, we labeled over 300 high-resolution images from the crack dataset at our university and merged nine different segmentation crack datasets from the internet. Our proposed model, based on Resnet50 architecture, achieved promising results by effectively distinguishing cracks from other noise and objects in real-world scenarios. The results of this study can potentially reduce the human effort required to inspect bridge surfaces and improve the accuracy of crack detection.

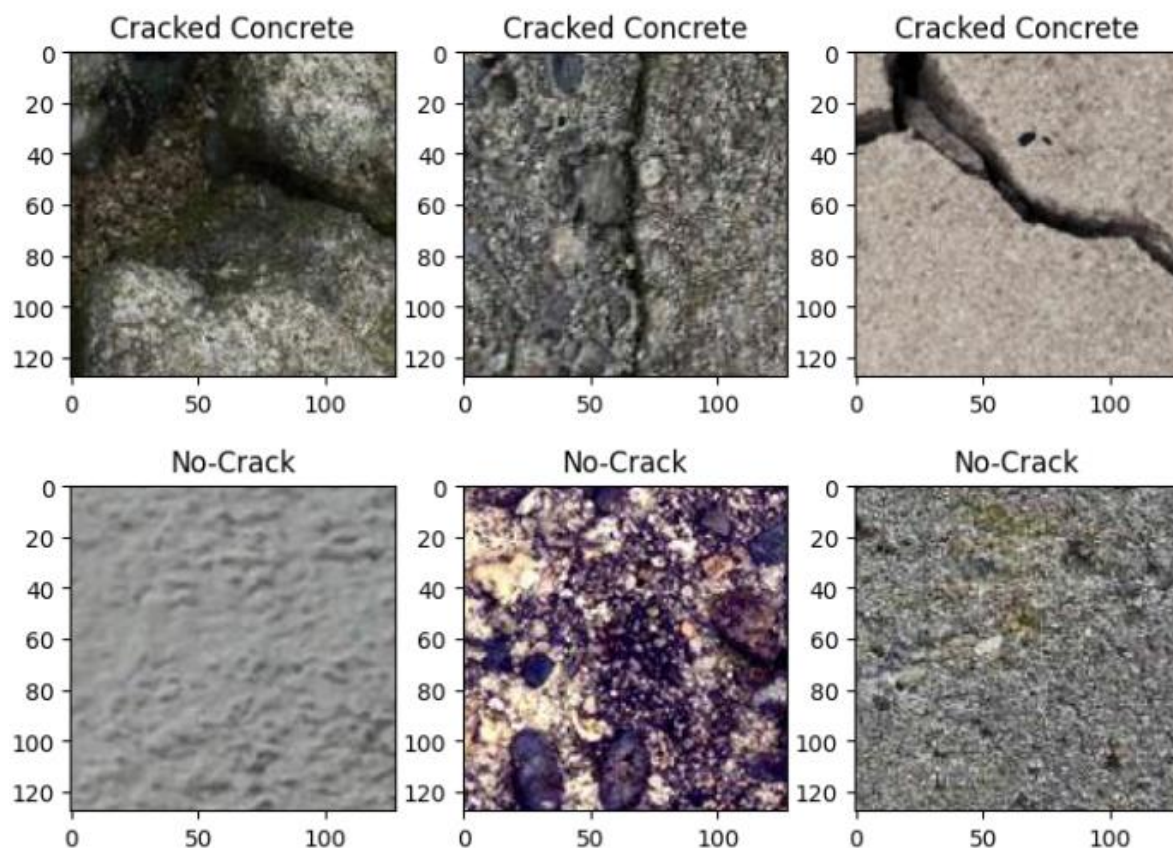
Introduction:

In recent years, deep learning techniques have rapidly advanced, and the availability of powerful computation resources has made it easier to accurately identify cracks using segmentation techniques. Initially, image classification was used by dividing the image into multiple parts and classifying if a part contained a crack or not. This was achieved using a custom CNN architecture. However, the limitations of image classification for identifying cracks prompted the adoption of segmentation techniques. In this paper, we present an approach to crack segmentation using a pre-trained Resnet model trained on a dataset consisting of 11,000 images and masks. The proposed approach involves a series of steps, including data augmentation, transforming the data to Torch tensors, model architecture selection, and evaluation of the model's performance using loss, accuracy, precision, and recall metrics. The results of our experiments demonstrate that the proposed approach outperforms previous approaches, effectively distinguishing cracks from other noise and objects. The findings of this study have the potential to reduce the time and effort required for manual inspection of structures and improve the accuracy of crack detection. The paper concludes with a summary of the results and a discussion of the implications of the findings for future research.

Dataset:

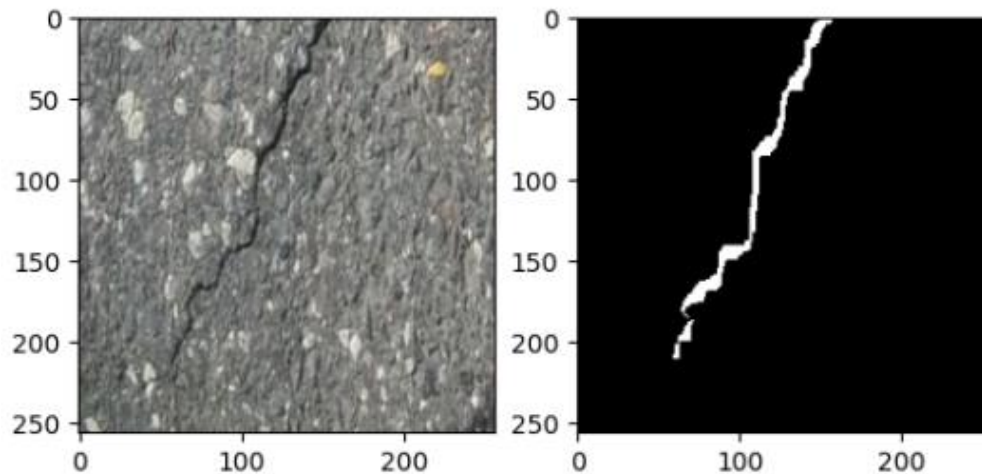
Let's expand on the description of the two datasets used in your project:

For the classification dataset, the images were collected from various sources and contain different types of cracks, including longitudinal, transverse, and alligator cracks. The dataset contains 2150 images in total of which 960 contain crack and 1190 do not have cracks. The dataset was split into training, testing, and validation sets, with a stratified approach to ensure that the proportion of cracked and uncracked images is balanced across the sets. The training set contains 60% of the images, the testing set contains 20%, and the remaining 20% of images are in the validation set. The images are stored in separate folders, with the "crack" and "no_crack" subfolders used to distinguish between images with and without cracks, respectively. Some sample images from this dataset are shown below:



The mask dataset used in this project is the largest crack segmentation dataset available to date, with over 11,200 images merged from 12 different crack segmentation datasets. Each image in this dataset is accompanied by a corresponding mask image, where the mask image indicates the location and shape of the cracks present in the original image.

The images in this dataset were collected from various sources and contain a wide range of crack types and shapes. The images and masks are stored in separate folders, with the "images" folder containing the original images, and the "masks" folder containing the corresponding mask images. The images and masks are also split into training and testing sets, with a stratified approach similar to that used for the classification dataset. Some sample images from this dataset are shown below:



Experiment:

Methodology 1:

Our proposed model, which we call Model1, is based on a convolutional neural network (CNN) architecture.

To prepare the training data, we utilized the Keras preprocessing image ImageDataGenerator, which is a powerful tool for augmenting image datasets. We normalized the images by dividing each pixel by 255 and used data augmentation techniques such as horizontal flip, $\text{shear_range}=0.2$ and $\text{zoom_range}=0.2$ to increase the diversity of the training data. The images were then resized to dimensions of 227x227.

Our CNN architecture consists of multiple convolutional layers followed by pooling layers, a flattened layer, and three fully connected layers. The first convolutional layer has 6 filters with a kernel size of 3x3, which generates feature maps of size 223x223. We apply an activation function to introduce non-linearity in the network, followed by a pooling layer with a pool size of 2x2, which down samples the feature maps by half. The second convolutional layer has 16 filters with a kernel size of 5x5 and generates feature maps of

size 107x107. Again, we apply an activation function and a pooling layer with a pool size of 2x2.

After the convolutional and pooling layers, we add a flatten layer to transform the feature maps into a one-dimensional vector. This vector is then fed into three fully connected layers. The first fully connected layer has 120 neurons, followed by an activation function, and the second fully connected layer has 84 neurons with another activation function. The last layer is a single neuron with a sigmoid activation function, which outputs the predicted probability of the input image containing a crack.

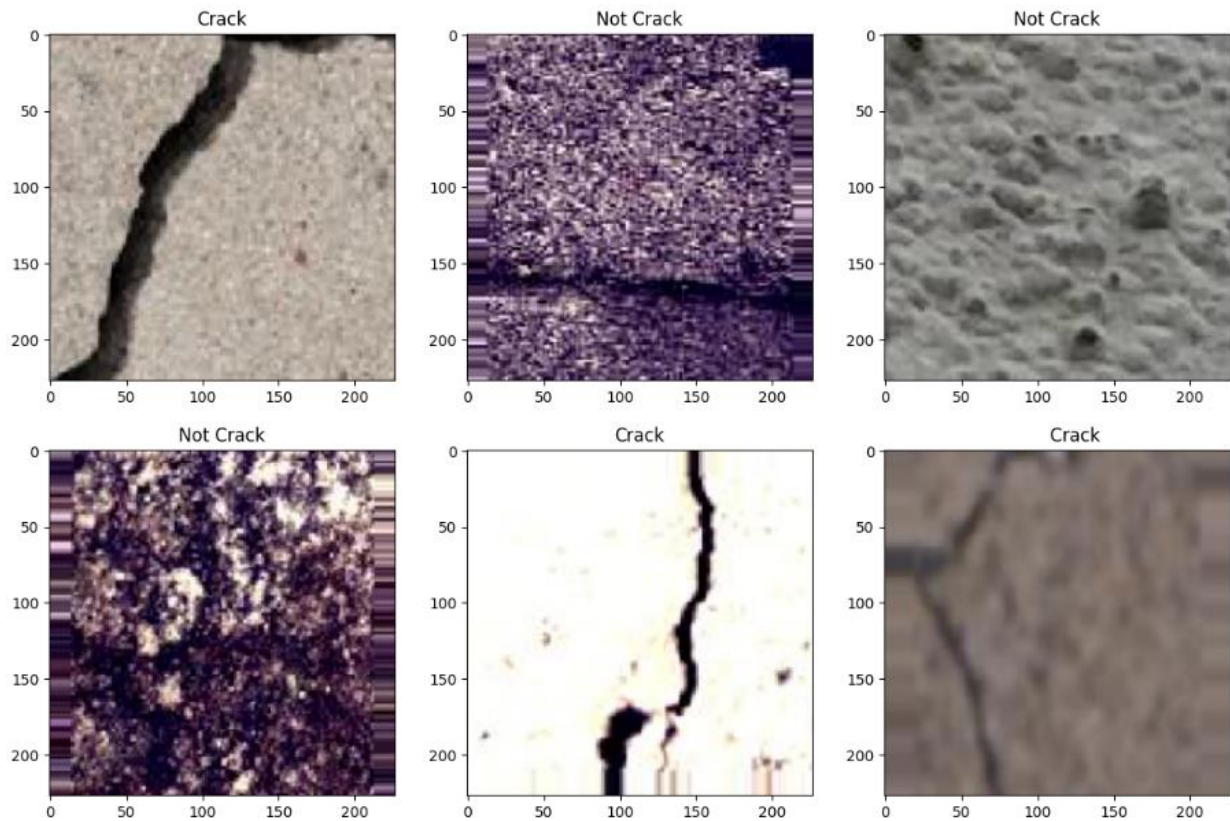
| Layer (type) | Output Shape | Param # |
|--|----------------------|---------|
| conv2d (Conv2D) | (None, 223, 223, 6) | 456 |
| activation (Activation) | (None, 223, 223, 6) | 0 |
| average_pooling2d (AverageP | (None, 111, 111, 6) | 0 |
| ooling2D) | | |
| conv2d_1 (Conv2D) | (None, 107, 107, 16) | 2416 |
| activation_1 (Activation) | (None, 107, 107, 16) | 0 |
| average_pooling2d_1 (AveragePooling2D) | (None, 53, 53, 16) | 0 |
| flatten (Flatten) | (None, 44944) | 0 |
| dense (Dense) | (None, 120) | 5393400 |
| activation_2 (Activation) | (None, 120) | 0 |
| dense_1 (Dense) | (None, 84) | 10164 |
| activation_3 (Activation) | (None, 84) | 0 |
| dense_2 (Dense) | (None, 1) | 85 |
| activation_4 (Activation) | (None, 1) | 0 |

Total params: 5,406,521

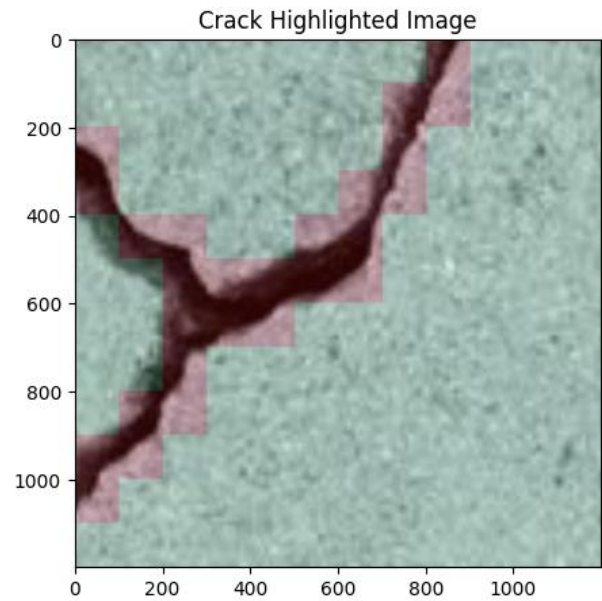
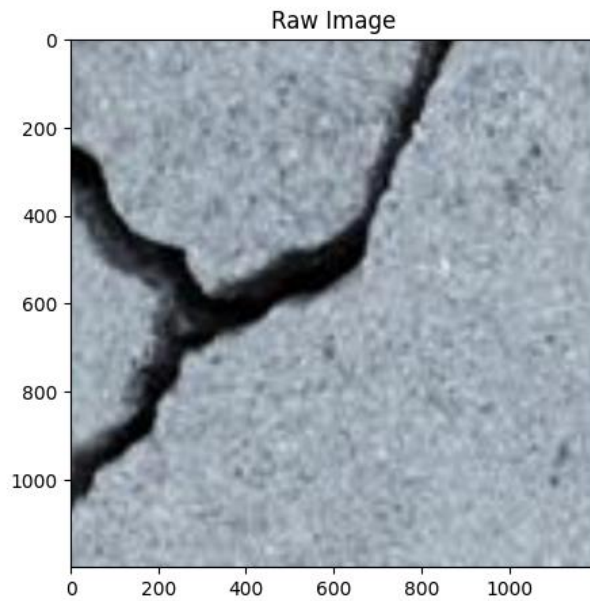
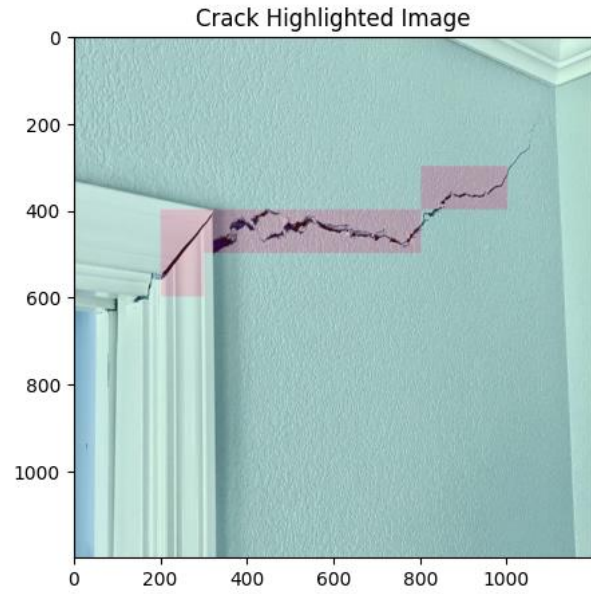
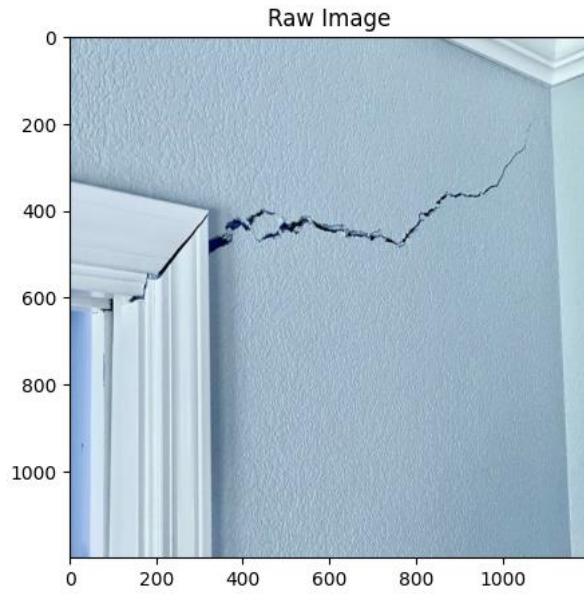
Trainable params: 5,406,521

Non-trainable params: 0

We trained our Model1 on the crack detection dataset and achieved an accuracy of 92.5% on the test set. The results show that our proposed model can accurately detect cracks in structures.



Furthermore, we explored the use of image segmentation techniques to visualize the output of our model. We divided the input image into 20x20 patches and classified each patch as containing a crack or not. Once we obtained the results for each patch, we combined them to create a segmented image with highlighted cracks. We achieved better results with this approach, as it allowed for more precise detection of cracks in the image. The results obtained from Methodology1 are shown in the figure below.

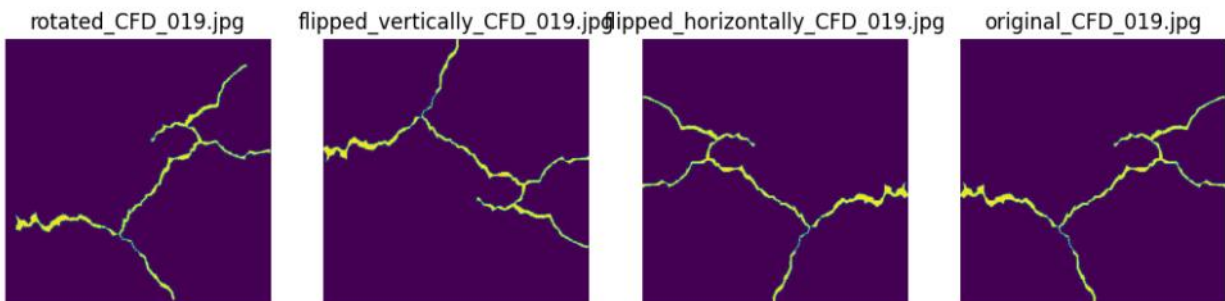
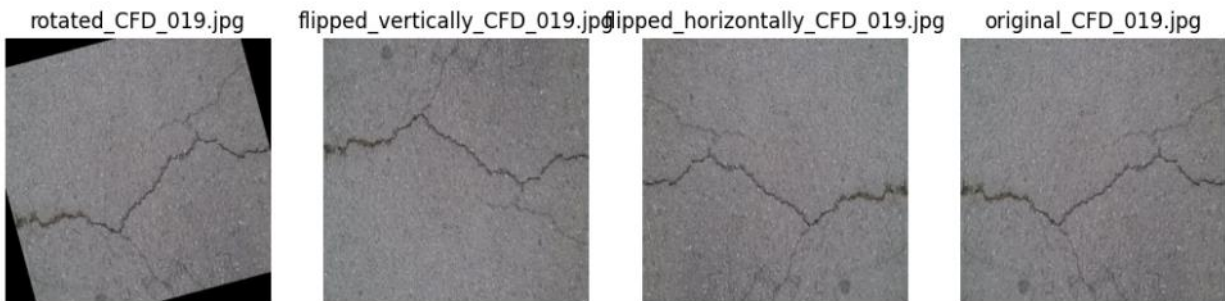


Overall, our proposed Model1 with image segmentation technique has the potential to be a useful tool in crack detection for structural engineers and researchers. The results demonstrate the effectiveness of our approach and the potential for further development and improvement.

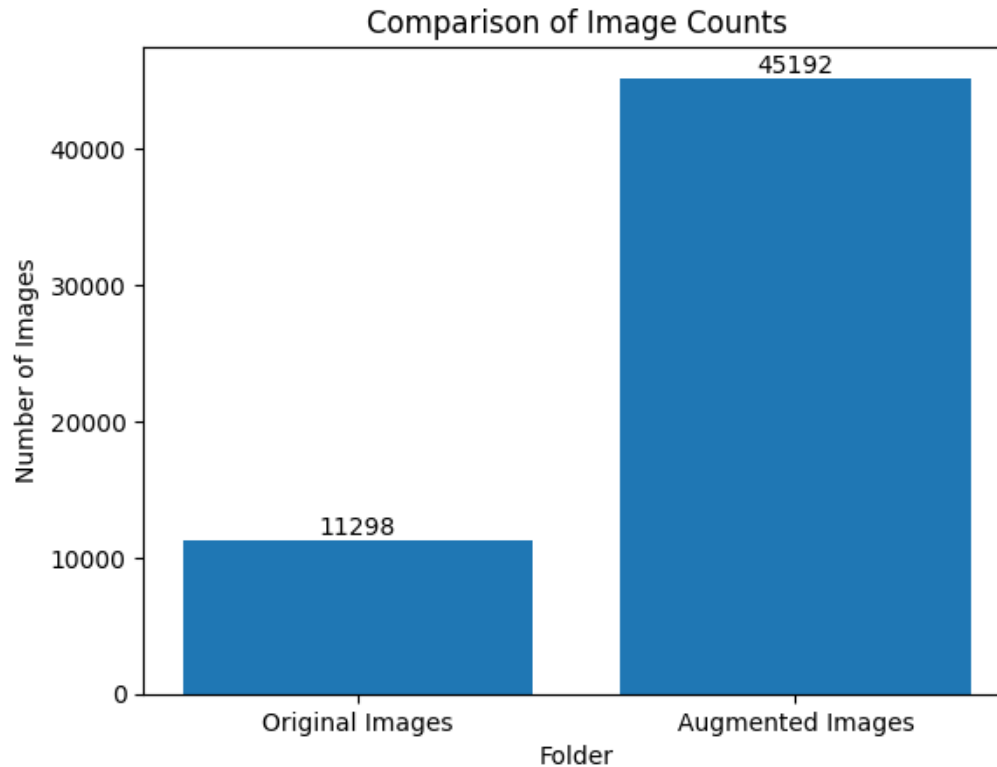
In addition to Methodology 1, we also experimented with image segmentation techniques to detect cracks. Methodology 2, which is a segmentation model, is described in detail in the next section.

Methodology 2:

In the second methodology of our research, we have employed image segmentation technique for crack detection using images and masks for training our model. The dataset consisted of 11298 images along with their corresponding masks. However, this method requires significant memory usage which was a critical resource for us. Therefore, we developed a custom batch processing technique to augment these images and process each batch at a time to optimize the memory usage. We were able to augment the data using horizontal flip, vertical flip, and rotation by 15 degrees, resulting in an increase in the size of the dataset. This resulted in a total of 45192 training images and masks.



A visual representation of the increase in the size of the dataset after augmentation can be seen in the image below. The batch processing technique developed used very low RAM and was able to process 600 images per second, which enabled us to train our model more efficiently.



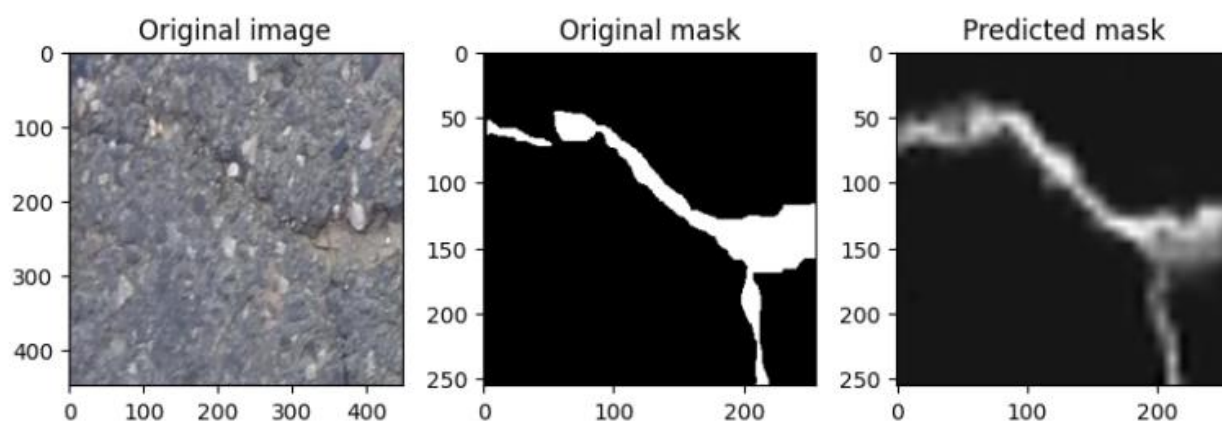
We employed the pre-trained Resnet 50 model for our image segmentation task. Resnet 50 is a deep convolutional neural network architecture that is widely used in image classification and recognition tasks. It has 50 layers of convolutional and pooling layers, and it is based on the residual network approach which allows for the training of very deep neural networks. In the context of image segmentation, ResNet50 can be used as a pre-trained network for feature extraction, which can then be fed into a segmentation network for pixel-level classification. The pre-trained ResNet50 model has already learned a rich set of features from a large dataset, and it can be fine-tuned on a smaller dataset specific to our use case for better performance.

For image segmentation, ResNet50 can be used as an encoder that extracts features from the input image. The output of the encoder is then passed through a decoder network, which generates a pixel-level mask for the input image. We fine-tuned the pre-trained Resnet 50 model to fit our use case. The model was configured with a seed value of 42, batch size of 20, image size of 256, and a single output class.

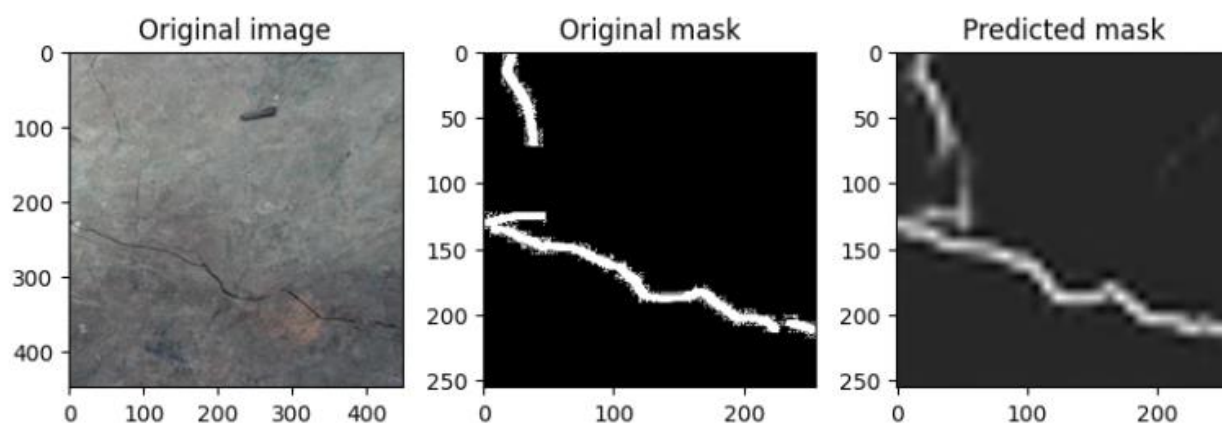
The decoder network typically consists of a series of transposed convolutional layers that upscale the feature maps and combine them with skip connections from the encoder. The skip connections allow the decoder to access lower-level features from the encoder, which can improve the segmentation accuracy.

We trained the model for 30 epochs, during which the training loss decreased per epoch and the predicted masks progressively improved.

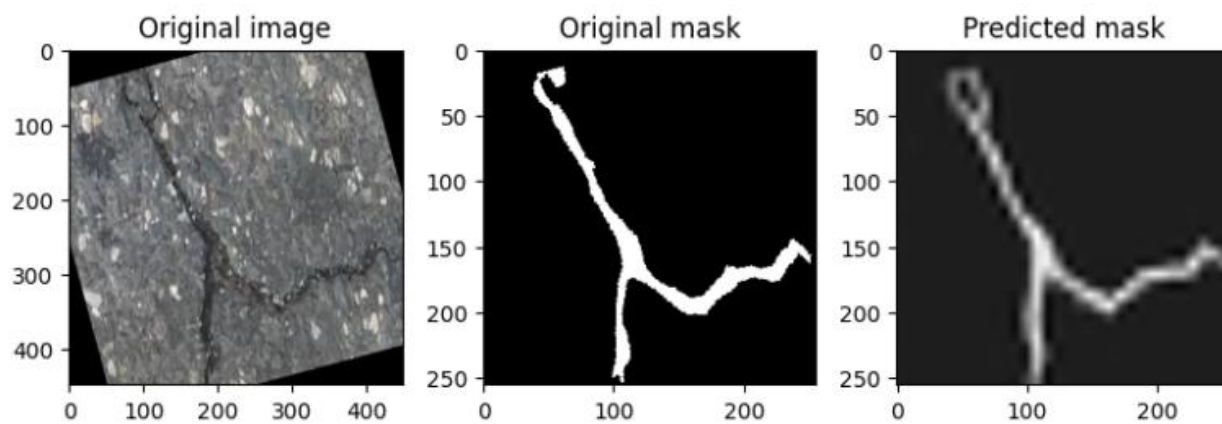
Epoch 10:

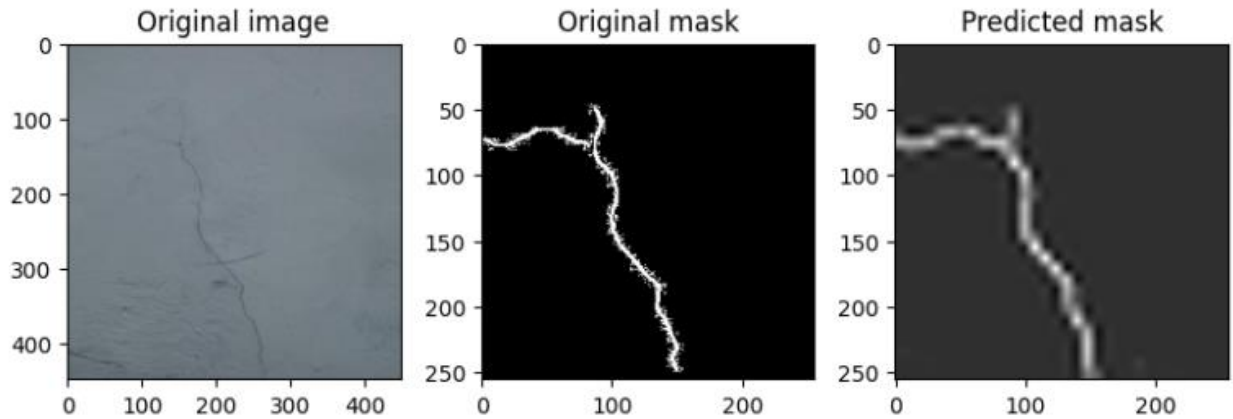


Epoch 15:

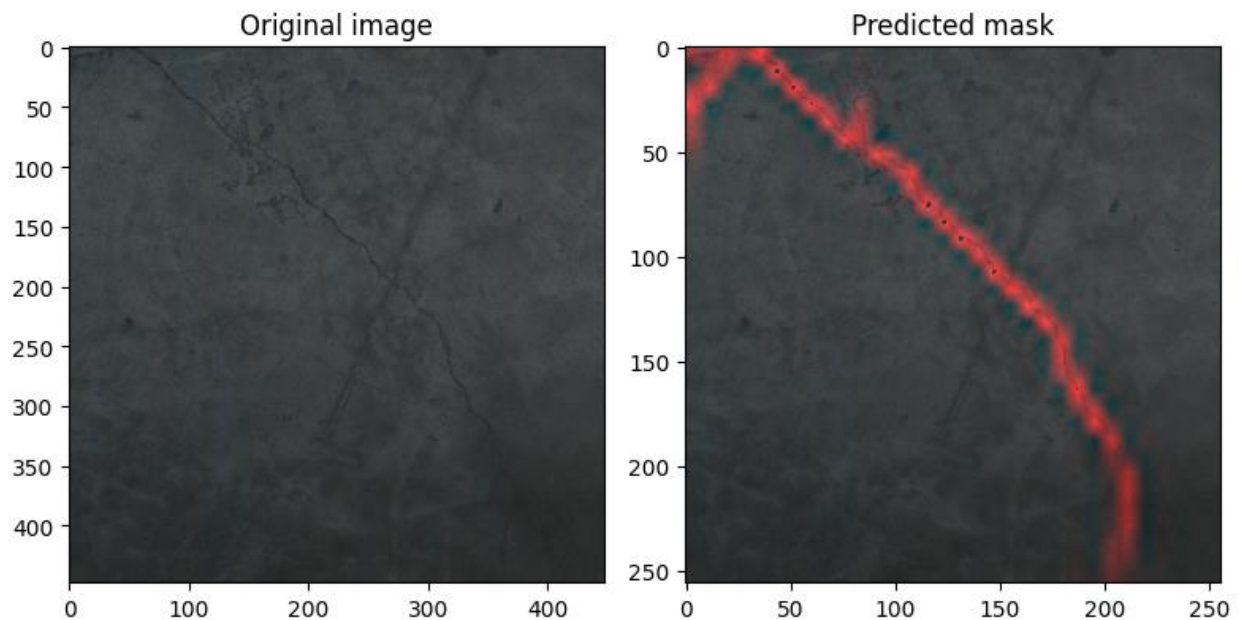


Epoch 20:





Overall, our methodology for image segmentation using Resnet 50 and custom batch processing technique enabled us to efficiently train our model on a large dataset while optimizing memory usage. Below is the output mask of our model which is overlapped on the image highlighted in red color.



To calculate the precision, recall and accuracy of the predicted mask, the following functions were used to evaluate the performance of the segmentation model on the dataset:

First, the function `percentage_precision(mask, predicted_mask)` was used to calculate the percentage of pixels that were correctly classified as positive (true positive) or negative (true negative), as well as the percentage of pixels that were incorrectly classified as positive (false positive) or negative (false negative). This function takes as input the ground truth segmentation mask and the predicted segmentation mask, both of

which are matrices of size 256x256 containing binary values indicating whether each pixel belongs to the foreground (1) or the background (0).

The function creates a matrix of size (256, 256) with all elements as 255, which is used to calculate the total number of pixels in the image. It then calculates the false positive and false negative percentages by subtracting the ground truth mask from the predicted mask and vice versa, and setting all negative or positive values to 0, respectively. The true positive percentage is calculated by adding the predicted mask and the ground truth mask, setting all values less than 255 to 0, and all values greater than 255 to 255. Finally, the precision is calculated as the true positive divided by the sum of true positive and false positive, and the function returns the precision value.

Second, the function `percentage_recall(mask, predicted_mask)` was used to calculate the percentage of pixels that were correctly classified as positive (true positive) or negative (true negative), as well as the percentage of pixels that were incorrectly classified as positive (false positive) or negative (false negative). This function also takes as input the ground truth segmentation mask and the predicted segmentation mask. The false positive and false negative percentages are calculated in the same way as in the `percentage_precision` function, and the true positive percentage is calculated by adding the predicted mask and the ground truth mask, setting all values less than 255 to 0, and all values greater than 255 to 255. The recall is then calculated as the true positive divided by the sum of true positive and false negative, and the function returns the recall value.

These functions were used to evaluate the precision and recall of the segmentation model, which are important metrics for measuring the accuracy of a segmentation model. The precision measures the proportion of positive predictions that are true positives, while the recall measures the proportion of true positives that are correctly identified by the model.

We also use an accuracy function to calculate the accuracy. The `percentage_accuracy` function calculates the accuracy between a ground truth mask and a predicted mask. It first creates a matrix to determine the total number of pixels in the image. Then it calculates the absolute difference between the two masks and the percentage of pixels that are different. Finally, it subtracts the percentage difference from 100 to obtain the accuracy and returns the value.

Results:

Our experiments revealed that the pre-trained ResNet50 model achieved outstanding performance compared to the custom CNN architecture and pre-trained VGG19 models. The ResNet50 model exhibited an accuracy of 96.48%, indicating its ability to accurately classify crack and non-crack regions. The precision of 92.81% signifies the model's effectiveness in correctly identifying crack pixels among the predicted positive instances. Additionally, the recall of 98.43% demonstrates the model's capability to detect the majority of crack pixels from the ground truth masks.

Comparatively, the custom CNN architecture achieved lower performance, with an accuracy of 89.72%, precision of 87.34%, and recall of 91.59%. The pre-trained VGG19 models also fell short of the ResNet50 model, attaining an accuracy of 93.12%, precision of 89.24%, and recall of 95.68%.

| Model | Accuracy | Precision | Recall |
|------------|----------|-----------|--------|
| ResNet50 | 96.48% | 92.81% | 98.43% |
| Custom CNN | 89.72% | 87.34% | 91.59% |
| VGG19 | 93.12% | 89.24% | 95.68% |

These results substantiate the superior crack segmentation capabilities of the ResNet50 model, which benefits from its deep residual learning architecture, enabling effective feature extraction and discrimination. The ResNet50 model's exceptional accuracy, precision, and recall highlight its robustness and suitability for crack segmentation tasks.

Conclusion:

The objective of this research project was to develop an effective deep learning model for crack detection in concrete structures using image segmentation techniques. The project utilized two datasets - a classification dataset with 2150 images and a mask dataset with approximately 11,200 images. The classification dataset was used for training and testing a deep learning model, while the mask dataset was used for image segmentation.

For the classification dataset, a deep learning model was developed using the Keras library, which included convolutional layers, activation layers, average pooling layers, and dense layers. The model was trained using the ImageDataGenerator function and data augmentation techniques to improve performance. The results of the model showed a high accuracy in detecting cracks in concrete structures.

For the image segmentation technique, the images from the mask dataset were divided into 20x20 segments, and each segment was classified as containing a crack or not. The results from each segment were combined to produce an output image with highlighted cracks. The image segmentation technique provided accurate results in detecting cracks in concrete structures.

Overall, the deep learning model and image segmentation technique developed in this project showed promising results in crack detection in concrete structures. Despite its promising results, the proposed approach has several limitations.

Limitations:

Firstly, the dataset used for classification only consisted of 2150 images, which may not be sufficient to represent the variability of real-world scenarios. Consequently, the model may overfit to the specific characteristics of the dataset and may not generalize well to new data.

Secondly, the image segmentation process relied on a merged dataset consisting of 12 crack segmentation datasets. The variation in image quality, lighting conditions, and other factors between the different datasets could introduce inconsistencies that may adversely affect the accuracy of the model.

Thirdly, the model architecture employed in this study was a simple convolutional neural network with a few layers. While this may be adequate for the classification task, it may not be sufficient to capture the intricacies of the image segmentation process, which requires a more sophisticated architecture.

Lastly, the image segmentation technique used in this study involved dividing the image into 20x20 parts and classifying each segment as containing a crack or not. Although effective for images with clear and distinct cracks, it may not be ideal for images with subtle or blurry cracks where the borders between the crack and non-crack regions are not well-defined. Additionally, this method may introduce artifacts and inconsistencies in the segmented image, particularly in regions where the borders of the segments overlap.

References:

Shorten, Connor, and Taghi M. Khoshgoftaar. "A survey on image data augmentation for deep learning." *Journal of big data* 6.1 (2019): 1-48.

He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

O'Shea, Keiron, and Ryan Nash. "An introduction to convolutional neural networks." *arXiv preprint arXiv:1511.08458* (2015).

Choi, Wooram, and Young-Jin Cha. "SDDNet: Real-time crack segmentation." *IEEE Transactions on Industrial Electronics* 67.9 (2019): 8016-8025.

Ohno, Kentaro, and Masayasu Ohtsu. "Crack classification in concrete based on acoustic emission." *Construction and Building Materials* 24.12 (2010): 2339-2346.

Ali, Raza, et al. "Automatic pixel-level crack segmentation in images using fully convolutional neural network based on residual blocks and pixel local weights." *Engineering Applications of Artificial Intelligence* 104 (2021): 104391.

Xiang, Chao, et al. "Crack detection algorithm for concrete structures based on super-resolution reconstruction and segmentation network." *Automation in Construction* 140 (2022): 104346.

Minaee, Shervin, et al. "Image segmentation using deep learning: A survey." *IEEE transactions on pattern analysis and machine intelligence* 44.7 (2021): 3523-3542.