
VitCrackSeg:Crack Segmentation Using Vision Transformer

Oscar Poudel

New Jersey Institute of Technology
University Hts, Newark, NJ 07102
op72@njit.com

Abstract

1 This study explores the efficacy of a Vision Transformer (ViT) model for the task
2 of real-time crack detection in various surfaces which is an essential component
3 of structural health monitoring in civil engineering and developing a architecture
4 VitCrackSeg. Utilizing a dataset designed for crack segmentation, the model was
5 trained using Binary Cross-Entropy with Logits (BCELogit) loss further enhanced
6 with jitter transformations to improve its generalization capabilities on unseen data.
7 The model's performance was evaluated based on the Intersection over Union (IoU)
8 and Dice coefficient metrics, achieving moderate scores of 0.4229 and 0.4991,
9 respectively. Additionally, the model was trained using Mean Squared Error (MSE)
10 loss for comparison, which yielded a lower accuracy of 45%, indicating the superior
11 suitability of BCELogit. The findings demonstrate the potential of using advanced
12 deep learning architectures like ViT for real-time, accurate crack detection. This ap-
13 proach not only aids in the rapid assessment of structural integrity but also enhances
14 maintenance strategies through timely interventions. A prediction pipeline was
15 developed for further application of the model in real-world scenario. This study
16 establishes a foundation for leveraging cutting-edge Deep learning techniques in
17 critical infrastructure monitoring, contributing to safer and more resilient civil struc-
18 tures. The codes are available at https://github.com/oscarpoudel/DL_VitCrackSeg

19

1 Introduction

20 Deep learning (DL) is transforming various fields of engineering, including civil engineering, by
21 providing powerful techniques to enhance data analysis, automation, and accuracy in inspection and
22 maintenance tasks. In particular, Vision Transformers (ViTs) have emerged as a groundbreaking
23 architecture with significant potential in visual inspection tasks due to their scalability, self-attention
24 mechanisms, and ability to process complex structural images.

25 The classification of building superstructures, such as aprons and ceilings, is crucial for assessing
26 structural integrity. Traditional manual inspections are subjective and inefficient, while convolutional
27 neural networks (CNNs) like VGG16 and advanced models like ViT and Token-to-Token Vision
28 Transformer (T2T-ViT) deliver a more objective, scalable approach. ViT uses self-attention heads to
29 efficiently analyze smaller image patches while maintaining accuracy and reducing computational
30 costs, making it suitable for classifying structural components and detecting critical damage[1]. A
31 novel ViT-based approach has also been proposed for achieving high-resolution segmentation in
32 civil engineering inspections. Eltouny et al. (2022) highlight a hybrid framework combining ViT
33 with learnable downsample-upsample modules for managing global and local semantics trade-offs
34 in high-resolution inspection images, while retaining both fine local details and global contextual
35 information. This combination ensures high fidelity in structural analysis, enabling accurate detection
36 of cracks, spalling, and other structural defects [2]. Another key challenge is defect classification
37 in various environments. Chow et al. (2020) proposed an artificial intelligence pipeline using

38 deep learning for anomaly detection and defect classification in concrete structures. This pipeline
 39 combines deep classifiers with anomaly extraction methods to identify and classify defects under
 40 varying conditions like lighting, camera angles, and distances, achieving an average testing accuracy
 of 95.6 Furthermore, Hoskere et al. (2022) emphasized the importance of ViT in 3D synthetic

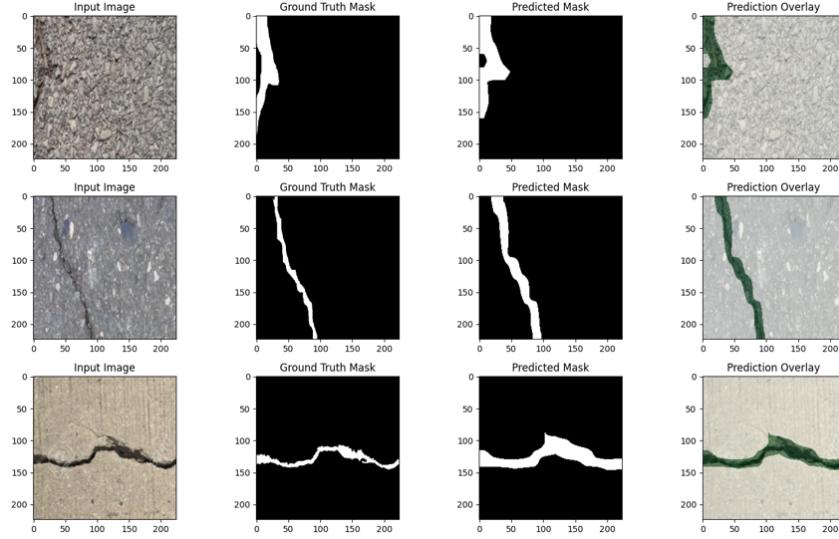


Figure 1: Cracksegmentation results on testing set after feeding them through VitCrackSeg models (first column being input, second column groundtruth, third column predicted mask and fourth column being predicted mask over input image)

41 environments. They developed an autonomous inspection testbed using Vision Transformers to
 42 classify damage and structural conditions in virtual earthquake-damaged buildings, demonstrating
 43 the improvement of deep learning models in real-world damage inference [4]. Another innovative
 44 application involves active learning strategies to tackle imbalanced datasets. Frick et al. (2022)
 45 proposed a method using a binary discriminator to prioritize data annotation tasks for training ViT
 46 models in civil engineering inspections. Their method improved 38% over traditional active learning
 47 methods[5].

48 This research project particularly emphasizes the identification of concrete cracks, corrosion, spalling,
 49 and other damage types to ensure the structural health of buildings and infrastructure. By leveraging
 50 the powerful features of deep learning and Vision Transformer (ViT) models, engineers can improve
 51 the precision of defect detection and classification. Additionally, the incorporation of real-time
 52 monitoring through camera systems enables immediate detection and response, ultimately leading to
 53 safer and more resilient infrastructure. The segmentation result from the model are shown in Figure 1

55 2 Method

56 2.1 Dataset

57 The Crack Segmentation Dataset available on Kaggle represents a significant resource for researchers
 58 and engineers dedicated to the advancement of structural health monitoring and the development of
 59 cutting-edge image processing technologies. This dataset combines together approximately 11,200
 60 high-resolution images sourced from 12 distinct crack segmentation datasets.

61 **Comprehensive Data Collection:** Each image within this dataset has been standardized to a uniform
 62 resolution of 448x448 pixels which ensured consistency across data inputs. This uniformity allows
 63 for streamlined processing and analysis, eliminating the need for additional preprocessing steps to
 64 normalize image sizes from various sources.

65 **Detailed Annotations and Masking** Accompanying each image are corresponding masks that outline
 66 the crack patterns. These masks provide labels that guide the algorithm in distinguishing between

67 cracked and uncracked regions, a critical factor in developing reliable and accurate crack detection
68 systems.

69 **Data Structure and Accessibility** The dataset is organized into four main folders: Images: Contains
70 all the high-resolution images. Masks: Includes the corresponding masks for each image, highlighting
71 the crack patterns. Train: A subset of images and masks designated for training purposes. The
72 selection of images in this folder ensures a representative mix of various conditions and types of
73 cracks. Test: Contains images and masks for testing and validating the performance of developed
74 models. This split is stratified, maintaining a proportional representation of each source dataset to
75 ensure comprehensive testing conditions.

76 **Non-Crack Data Inclusion** Moreover, the dataset includes images labeled with the prefix "non-
77 crack*", indicating the absence of crack features within these images. This inclusion is important
78 for training models capable of correctly identifying negative instances where no cracks are present,
79 thereby enhancing the model's accuracy and generalability in real-world applications.

80 **Use Cases and Applications** The Crack Segmentation Dataset is designed not only for academic re-
81 search but also for practical applications in civil engineering, architecture, infrastructure maintenance,
82 and disaster response. By facilitating the development of algorithms that can automatically detect
83 and analyze structural cracks, this dataset plays a crucial role in preventative maintenance strategies
84 and safety inspections of buildings, bridges, and other critical structures.

85 **2.2 Preprocessing**

86 To ensure reproducibility in experiments that involve stochastic processes such as data shuffling and
87 weight initialization, seeds for the random number generators are initialized at the beginning of the
88 script. Specific seed values are set for both NumPy and PyTorch, with a chosen seed of 1010. This
89 step is crucial as it makes the model's behavior deterministic, ensuring that results can be replicated
90 exactly across different runs and on different machines. The configuration of the model's training
91 parameters is established early in the script. The batch size is set to 32 which is the number of samples
92 processed before the model updates its internal parameters. The model is defined to recognize one
93 class as a segmentation is a binary classification task. Training is set to run for 60 epochs, allowing
94 the model sufficient time to learn from the data. A DataFrame is constructed to associate each image
95 with its corresponding mask, simplifying the process of accessing the data during model training.
96 The dataset is split into training, validation, and test sets, with 30% of the data allocated to the test
97 set and 15% of the remaining training data allocated to validation. This stratified split ensures the
98 model is evaluated on a representative sample of data, enhancing its ability to generalize well to new,
99 unseen data.

100 The images and masks undergo several preprocessing steps to prepare them for use in training the
101 neural network. All images and masks are resized to 224x224 pixels to ensure consistency in input
102 data size. A series of transformations is applied to both the training and validation datasets, which
103 includes converting the images to the PIL format, resizing, and normalizing. The normalization uses
104 mean values of [0.485, 0.456, 0.406] and standard deviations of [0.229, 0.224, 0.225], which are
105 used for models pre-trained on the ImageNet dataset. Custom dataset classes and DataLoaders are
106 set up to manage data loading and batching efficiently. These components are essential for handling
107 different aspects of data management during the training process, such as shuffling the training data
108 and batching it for each iteration. This setup not only streamlines the training process but also ensures
109 that the model is fed with correctly formatted and preprocessed data vital for the successful training
110 of deep learning models, especially in complex tasks such as crack segmentation in images.

111 **2.3 Model Architecture**

112 The core of our deep learning approach involves the Vision Transformer (ViT), specifically leveraging
113 the pre-trained ImageNet-21k ViT model available through Hugging Face's Transformers library,
114 configured for image sizes of 224x224. Vision Transformers are particularly suitable for this task
115 due to their ability to process images as sequences of patches, applying self-attention mechanisms
116 that consider the global context of the image, which is vital for accurately detecting and segmenting
117 cracks. The architecture of the model is shown in Figure 2.

118 **Removing the Classification Head**

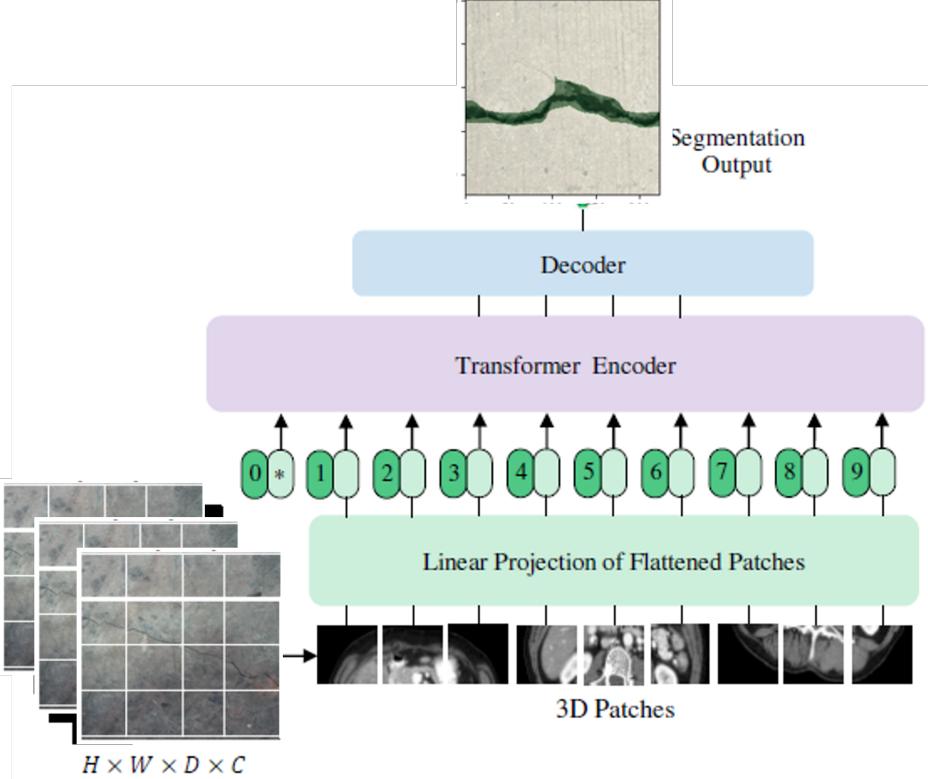


Figure 2: Representative architecture of the model

119 The standard ViT model is typically used for classification tasks. However, for the purpose of
 120 segmentation, the classification head of the pre-trained model is removed to adapt the model for
 121 generating image features rather than class outputs. This modification allows the transformer to
 122 output a sequence of embeddings corresponding to different patches of the image, which are then
 123 used for the segmentation task.

124 **Segmentation Head**

125 After processing the image through the transformer encoder and removing the classification head, the
 126 sequence of patch embeddings generated by the ViT is reshaped back into a two-dimensional spatial
 127 grid that mirrors the arrangement of patches in the original input image. This restructuring is crucial
 128 as it maps the learned deep features back to the spatial dimensions of the image, facilitating a more
 129 intuitive application of convolutional operations.

130 A convolutional layer with a 1x1 kernel size is subsequently applied to these reshaped features. This
 131 convolutional 'segmentation head' serves to map the deep feature representations of the patches to
 132 the desired number of output channels. For crack detection and segmentation, the output typically
 133 includes channels for different types of damage or different damage states (e.g., cracked, not cracked).

134 **Upsampling**

135 Due to the division of the original image into patches and the downscaling involved in matching
 136 them to the transformer's input resolution, the output from the segmentation head usually has a
 137 lower resolution than the original image. To address this, an upsampling step is incorporated, where
 138 the output from the segmentation head is up-scaled back to the dimensions of the original image.
 139 This upsampling is typically performed using bilinear interpolation, which helps in smoothing the
 140 output segmentation map and aligning it closely with the ground truth masks. This configuration not
 141 only preserves the high resolution of the input images but also ensures that the segmentation maps
 142 generated by the model are precise and detailed, enabling accurate pixel-wise classification of the
 143 structural integrity of the materials being analyzed.

```

class ViTForSegmentation(nn.Module):
    def __init__(self, num_classes, output_dim=(224, 224)):
        super(ViTForSegmentation, self).__init__()
        self.vit = ViTForImageClassification.from_pretrained('google/vit-base-patch16-224-in21k')
        self.segmentation_head = nn.Conv2d(768, num_classes, kernel_size=1)
        self.upsample = nn.Upsample(size=output_dim, mode='bilinear', align_corners=False)

    def forward(self, x):
        features = self.vit.vit(x).last_hidden_state

        batch_size, sequence_length, num_channels = features.size()
        # first token is the classification token and should be discarded for segmentation
        features = features[:, 1:, :]

        side_length = int((sequence_length - 1)**0.5) # Adjust for the removed classification token

        # Now the number of patches should form a perfect square
        assert side_length * side_length == (sequence_length - 1), "Number of patches does not form a perfect square"

        features = features.permute(0, 2, 1).reshape(batch_size, num_channels, side_length, side_length)

        # Apply the segmentation head and then upsample
        output = self.segmentation_head(features)
        output = self.upsample(output)
        return output

```

Figure 3: Code for the Model

144 By integrating the Vision Transformer with a tailored segmentation head and upsampling strategy,
145 this methodology leverages the strengths of transformers for feature extraction while adapting them
146 to the specific needs of high-resolution image segmentation tasks in civil engineering. This approach
147 ensures that the spatial relationships and detailed textural information crucial for crack detection are
148 effectively captured and utilized, resulting in highly accurate segmentation outputs that are critical
149 for structural health monitoring and assessment.

150 **Code**

151 As shown in the Figure 3, the ViTForSegmentation class, implemented in Python using the PyTorch
152 framework, is central to our model's architecture. It adapts a pre-trained Vision Transformer for
153 image segmentation by attaching a convolutional segmentation head that classifies each image patch
154 based on its features.

155 In the forward pass, after processing the input through the Vision Transformer to obtain the last hidden
156 states, the classification token is discarded. The patch embeddings are then reshaped into a 2D spatial
157 grid. This structured arrangement of features is crucial for maintaining spatial coherence, which
158 allows the subsequent convolutional layer to effectively perform segmentation based on localized
159 image features. The convolutional segmentation head then maps these features to the number of
160 desired output channels for each class in the segmentation task. This process effectively transforms
161 the deep feature representation into class-specific segmentation predictions for each patch.

162 The output from the segmentation head usually has a lower resolution due to the division into
163 patches. To address this, the output is upsampled back to the original image dimensions using bilinear
164 interpolation. This upsampling step is vital as it ensures that the final segmentation map is both precise
165 and aligned with the high-resolution input images, facilitating accurate pixel-wise segmentation.
166 This part of the methodology leverages the detailed textural information captured by the Vision
167 Transformer, applying it effectively for the precise localization required in the segmentation tasks.
168 This approach not only retains the integrity of the structural analysis but also enhances the model's
169 ability to detect minute details crucial for assessing the condition of infrastructure materials.

170 **2.4 Training**

171 In the training process, the study utilized the Binary Cross-Entropy with Logits (BCEWithLogitsLoss)
172 as the loss function. The optimizer of choice was the Adam optimizer with a set learning rate of
173 0.0001, facilitating fine-grained adjustments in the model weights. Data was processed in manageable
174 batches of 32, which balanced the computational load and allowed for efficient gradient updates.

175 To ensure robustness and prevent overfitting, several regularization techniques were employed
176 including gradient clipping and early stopping. Gradient clipping was used to prevent the exploding
177 gradient problem by capping the gradients during backpropagation to not exceed a defined threshold.

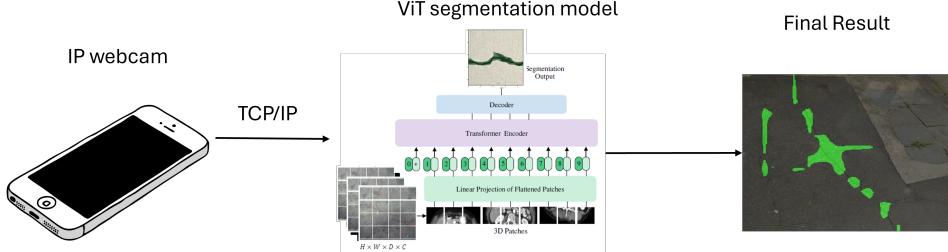


Figure 4: Prediction pipeline for the crack segmentation application

178 Early stopping was triggered after 10 epochs indicating that continuing training would likely lead to
 179 overfitting.
 180 The training was logged and monitored using TensorBoard, which provided a visual interface to
 181 track and analyze metrics like loss and accuracy over time, helping to fine-tune the training phases.
 182 The model's training was conducted on a high-performance hardware setup consisting of 32 GB
 183 of RAM and an NVIDIA RTX 3090 16 GB GPU. The final model comprised a significant number
 184 of trainable parameters, totaling 85,800,963, which underscores the complexity and depth of the
 185 network architecture used. Overall, the training was completed in approximately 2 hours.

186 2.5 Prediction Pipeline

187 The primary objective of the study is to implement a real-time crack detection system leveraging an IP
 188 webcam setup on a smartphone and a deep learning model running on a laptop. The prediction pipeline
 189 is shown in Figure 4. Data acquisition involves the smartphone configured as an IP webcam capturing
 190 live video footage, streamed over a TCP/IP network to a laptop for processing. Each video frame
 191 is analyzed using a Vision Transformer (ViT) segmentation model resulting precise segmentation
 192 masks overlayed on the original video frames. This segmented output is immediately displayed on the
 193 laptop, offering real-time visualization crucial for applications like in-field structural assessments or
 194 continuous infrastructure monitoring. This methodology combines advanced deep learning techniques
 195 with accessible technology to offer a robust tool for real-time safety and maintenance assessments
 196 across diverse industries.

197 The evaluation was conducted using two different metrics: Intersection over Union (IoU) and Dice
 198 coefficient, alongside traditional accuracy measures from model training with Mean Squared Error
 199 (MSE) loss.

200 Training and Validation Loss and Accuracy The training process was documented through a series of
 201 plots showing both training and validation loss and accuracy under two conditions: with and without
 202 jitter transformation. These conditions are crucial for understanding the model's robustness to input
 203 variations.

204 3 Results

205 Figure 5plot (a) and (b) illustrate the training and validation loss and accuracy when trained with
 206 Binary Cross-Entropy with Logits (BCELogit) loss without jitter. It is observed that the training
 207 loss decreases steadily, indicating good model convergence. Correspondingly, the training accuracy
 208 increases which reaches a plateau, which suggests that the model effectively learns from the training
 209 data without overfitting, as shown by the parallel trend in validation accuracy. Figure 5plot (c) and (d)
 210 display the results for training with BCELogit loss with jitter transformations. The inclusion of jitter,
 211 which introduces variations in the training images (such as changes in brightness or contrast), shows a
 212 slightly more erratic progression in both loss and accuracy curves. However, the model still improves
 213 over time, demonstrating adaptability to augmented data. Notably, the validation accuracy here does
 214 not plateau as distinctly as in the non-jitter case, possibly indicating a model that generalizes better
 215 on unseen data due to the diversity in training examples.

216 Performance Metrics

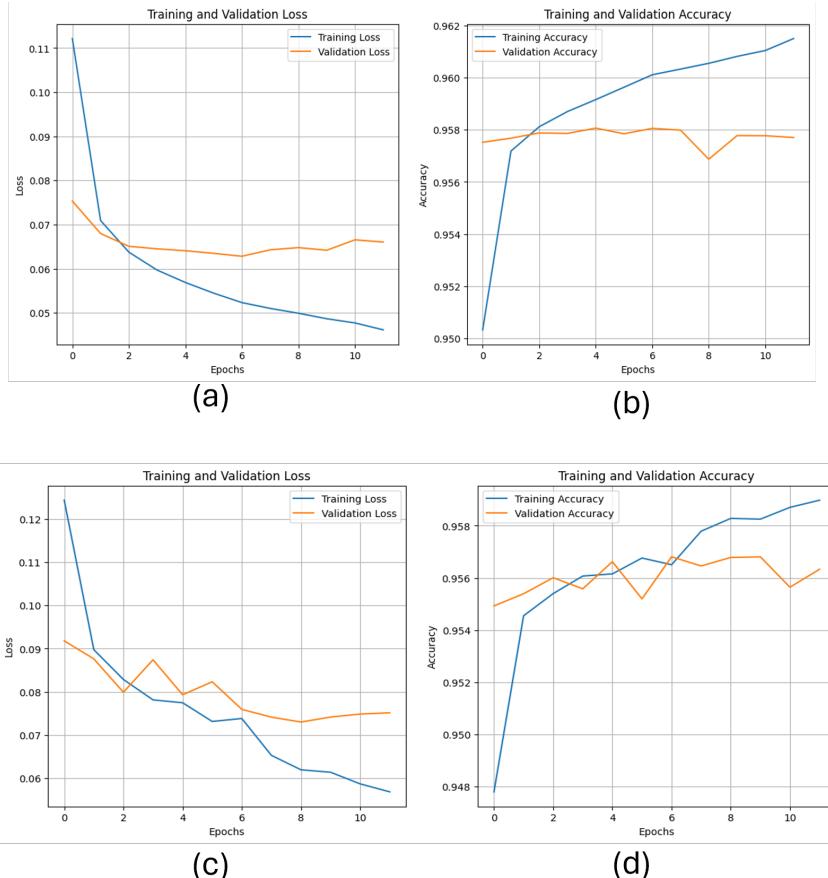


Figure 5: Training and validation (a)loss (b)accuracy for training with BCELOGIT loss without jitter. Training and validation (a)loss (b)accuracy for training with BCELOGIT loss with jitter transforms

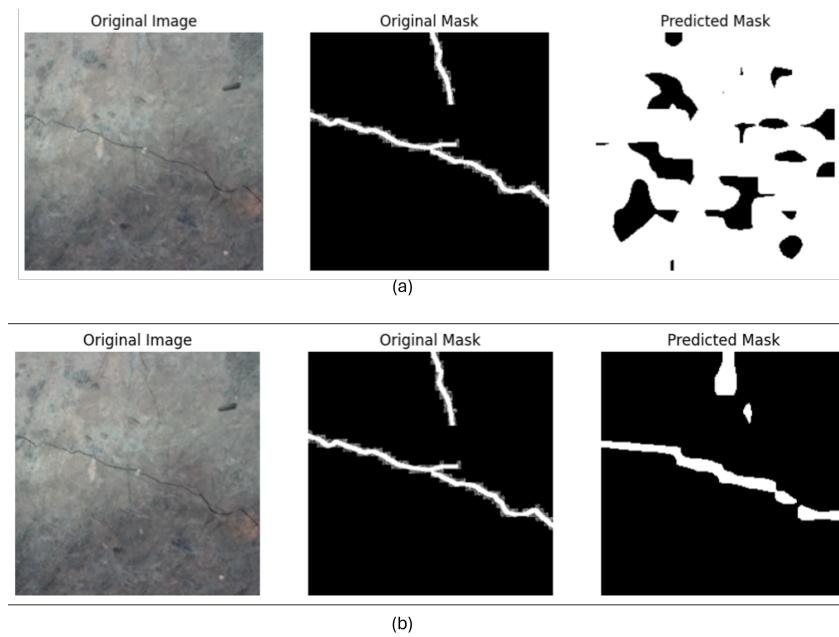


Figure 6: (a) Prediction from final model with MSE loss (b) Prediction from final model with BCElogitloss

217 Intersection over Union (IoU): The average IoU achieved by the model is 0.4229. This metrics
218 indicates a moderate level of accuracy in the model's segmentation capability.
219 Dice Coefficient: Similarly, the Dice coefficient, which is another method of measuring the efficacy
220 of the segmentation by measuring the overlap twice the intersection over the sum of the sizes of two
221 labels, is recorded at 0.4991. This value further supports the findings of the IoU, showcasing a nearly
222 50% effectiveness in identifying the correct regions of interest in the images. Model
223 Training with MSE Loss: When the model was trained using MSE loss, a commonly used loss
224 function for regression tasks, it achieved an accuracy of 45%. This result is notably lower compared
225 to the training with BCELogit loss, suggesting that MSE might not be as effective for the segmentation
226 tasks handled by this model, potentially due to its emphasis on averaging the errors in pixel values
227 rather than focusing on the categorical outcomes. The prediction results for the model with MSE loss
228 and BCElogit loss are shown in Figure 6

229 4 Discussion and Conclusion

230 The outcomes of the project hence shows a nuanced understanding of the performance dynamics with
231 a Vision Transformer (ViT) model under a changed operational framework for image segmentation
232 task focusing in the crack detection. Applying BCELogit with augmentation of jitter transformations,
233 the model learned to generalize well enough to present an average IoU of 0.4229 and Dice coefficient
234 of 0.4991 over unseen data. These metrics are of medium value; however, it can be inferred that the
235 model might help in successful crack pattern recognition performance, which is critical for real-time
236 structural health monitoring applications to successfully identify and segment out crack patterns.
237 Training the model with Mean Squared Error (MSE) loss gives an accuracy level that is drastically
238 less effective; it gives a reduced 45% accuracy level. This further shows that there is a need for the
239 right loss function, not only meant for the minimization of simple error but also one that is suited for
240 the specific given segmentation task.

241 In conclusion, the study has successfully demonstrated the application of a deep learning model,
242 leveraging a Vision Transformer architecture, for the segmentation of cracks in images—a task
243 with significant implications in the fields of civil engineering and infrastructure maintenance and
244 developed a model VitCrackSeg. The training process was optimized through strategic choices in loss
245 functions and regularization techniques, ensuring that the model was not only accurate but also robust
246 against overfitting and efficient in terms of computational resources. The use of advanced hardware,
247 along with systematic training approaches, resulted in a model that was trained in a relatively short
248 time frame, showcasing the feasibility of deploying such models in real-world scenarios where quick,
249 reliable crack detection is essential. The early stopping of training at 10 epochs due to non-improving
250 validation loss metrics is indicative of a well-tuned training process that prioritizes quality of learning
251 over mere completion of arbitrary training cycles. The prediction pipeline gives a tangible method
252 for the application of the deep learning model in the real world. Future work should focus on further
253 enhancing the model's accuracy and generalization capabilities, possibly by exploring more diverse
254 datasets and implementing more complex transformations and augmentation strategies. Additionally,
255 experimenting with other loss functions like dice loss and optimization techniques could provide
256 deeper insights into the model's performance and lead to higher precision in crack segmentation
257 tasks. The insights gained from this project underscore the potential of using advanced machine
258 learning techniques for critical applications in infrastructure health monitoring, with a view towards
259 improving safety and longevity of civil structures.

260 References

- 261 [1] J. T. S. Phang, K. H. Lim, and C. H. Lim, "Deep Learning Neural Networks in Building Superstructure
262 Classification," in 2023 International Conference on Digital Applications, Transformation Economy (ICDATE),
263 Miri, Sarawak, Malaysia: IEEE, Jul. 2023, pp. 1–5. doi: 10.1109/ICDATE58146.2023.10248658.
- 264 [2] K. Eltouny, S. Sajedi, and X. Liang, "High-Fidelity Visual Structural Inspections through Transformers and
265 Learnable Resizers." arXiv, Oct. 21, 2022. doi: 10.48550/arXiv.2210.12175.
- 266 [3] J. K. Chow et al., "Artificial intelligence-empowered pipeline for image-based inspection of concrete
267 structures," Autom. Constr., vol. 120, p. 103372, Dec. 2020, doi: 10.1016/j.autcon.2020.103372.

268 [4] V. Hoskere, Y. Narazaki, and B. F. Spencer, "Physics-Based Graphics Models in 3D Synthetic Environ-
269 ments as Autonomous Vision-Based Inspection Testbeds," Sensors, vol. 22, no. 2, p. 532, Jan. 2022, doi:
270 10.3390/s22020532.

271 [5] T. Frick, D. Antognini, M. Rigotti, I. Giurgiu, B. Grewe, and C. Malossi, "Active Learning for Imbalanced
272 Civil Infrastructure Data." arXiv, Oct. 19, 2022. doi: 10.48550/arXiv.2210.10586.