

# Automated Monitoring of Cultural Heritage Artifacts Using Semantic Segmentation

A. Ranieri, G. Palmieri, S. Biasotti

CNR-IMATI, Via De Marini, 6 - 16149 Genova (GE), ITALY  
`{andrea.ranieri,giorgio.palmieri,silviamaria.biasotti}@cnr.it`

**Abstract.** This paper addresses the critical need for automated crack detection in the preservation of cultural heritage through semantic segmentation. We present a comparative study of U-Net architectures, using various convolutional neural network (CNN) encoders, for pixel-level crack identification on statues and monuments. A comparative quantitative evaluation is performed on the test set of the OmniCrack30k dataset [1] using popular segmentation metrics including Mean Intersection over Union (mIoU), Dice coefficient, and Jaccard index. This is complemented by an out-of-distribution qualitative evaluation on an unlabeled test set of real-world cracked statues and monuments. Our findings provide valuable insights into the capabilities of different CNN-based encoders for fine-grained crack segmentation. We show that the models exhibit promising generalization capabilities to unseen cultural heritage contexts, despite never having been explicitly trained on images of statues or monuments.

**Keywords:** Cultural Heritage · Monitoring · Deep Learning · U-Nets · Semantic Segmentation.

## 1 Introduction

The preservation of cultural heritage, encompassing historical statues and monuments, is paramount for understanding human history and artistic achievement. These invaluable artifacts are constantly exposed to environmental degradation, leading to structural damage such as cracks. Early and accurate detection of these cracks is crucial for timely intervention, preventing further deterioration, and ensuring their longevity. Traditional manual inspection methods are often labor-intensive, time-consuming, subjective, and can be limited by accessibility, particularly for large or complex structures. This necessitates the development of automated, efficient, and precise diagnostic tools.

Semantic segmentation, a fundamental task in computer vision, offers a powerful solution by enabling pixel-level classification of images. Unlike object detection, which provides bounding box localization, semantic segmentation precisely delineates the boundaries of objects or regions of interest, making it ideal for identifying fine-grained, irregular structures like cracks. Recent advancements in deep learning, particularly with Convolutional Neural Networks (CNNs), have

revolutionized image segmentation, demonstrating superior performance in capturing intricate spatial patterns and semantic representations [2].

A significant challenge in applying deep learning to cultural heritage preservation is the inherent domain shift and scarcity of annotated data. While large-scale datasets exist for crack detection in civil infrastructure, such as roads and buildings [1,3,4], these often present vastly different visual characteristics and crack morphologies compared to the intricate surfaces of statues and monuments. Road cracks, for instance, typically appear on asphalt or concrete, following linear or alligator patterns, whereas cracks in cultural heritage objects can be more subtle, follow material grain (e.g., stone, marble), or occur on metallic surfaces, frequently against complex, textured backgrounds. This domain gap, coupled with the practical impossibility of acquiring large, pixel-level annotated ground truth datasets for unique historical artifacts, mandates a specialized approach to data preparation and evaluation.

This paper explores the application of state-of-the-art deep learning architectures for semantic segmentation of cracks [3,5] at the domain of cultural heritage. Specifically, we conduct a comparative analysis of U-Nets [6], a widely adopted encoder-decoder architecture known for its efficacy in image segmentation and fine-grained detail preservation, when paired with various CNN backbones as well as in image generation (such as in diffusion models). We utilize the OmniCrack30k dataset [1], a comprehensive benchmark dataset for crack segmentation, for training and quantitative evaluation. Our experimental setup involves training U-Nets with ResNet-50, ResNet-101, ConvNeXt V2 Base, and ConvNeXt V2 Huge as encoders [7,8]. Given the scarcity of large, annotated datasets specifically for cultural heritage crack detection, we employ a two-fold evaluation approach: a quantitative assessment on the OmniCrack30k test set and a qualitative evaluation on an unlabeled test set of real-world cracked statues and monuments.

The main contributions of this work are:

- An analysis of the applicability and performance of U-Net architectures with diverse CNN encoders for semantic segmentation of cracks on cultural heritage artifacts, including different CNN backbone complexities.
- A detailed methodology for leveraging the OmniCrack30k dataset for training, addressing the challenges of data diversity in crack segmentation.
- A comparative quantitative evaluation of model performance on the OmniCrack30k test set using mIoU, Dice, and Jaccard metrics.
- A qualitative evaluation framework for assessing model performance on an unlabeled test set of real-world cracked statues and monuments, providing practical insights for conservators and highlighting generalization capabilities to unseen domains.
- A **public repository containing all the code**<sup>1</sup> used to produce (and therefore to replicate) this work and the pre-trained models used to generate the images in this paper.

---

<sup>1</sup> <https://gitlab.com/4ndr3aR/cultural-artifacts-crack-segmentation>

## 2 Related Work

Automated crack segmentation is a critical component in structural health monitoring and integrity systems, with applications spanning road infrastructure, buildings, and cultural heritage. Research in this field has evolved from traditional image processing techniques to sophisticated deep learning paradigms.

### 2.1 Traditional vs. Deep Learning Approaches for Crack Segmentation

While these methods are computationally less intensive, traditional computer vision approaches often struggle with complex crack image backgrounds, varying illumination, noise, and the intricate spatial details of cracks, leading to limited effectiveness and reliability in heterogeneous or noisy environments. Early methods for crack detection relied on traditional image processing techniques such as edge detection (e.g., Canny [9,10], Sobel [11]), thresholding [12], morphological operations [13], and statistical analysis.

The advent of deep learning has significantly advanced crack segmentation. Deep learning approaches, particularly Convolutional Neural Networks (CNNs) can directly learn complex spatial patterns and feature representations from labeled data. Their ability to extract and match the most relevant features significantly improves the segmentation accuracy.

### 2.2 CNN and UNet-based Architectures

In the domain of cultural heritage, the scientific community has primarily focused on the classification and the object detection of cracks. In [14], the authors construct a dataset of 6002 images for binary crack classification on images of temples captured by an AUV in the Ayutthaya region of Thailand. They then train a three-layer CNN by comparing three different types of classifiers for CNN features: the fully connected layer of the CNN itself, a Support Vector Machine (SVM), and a Random Forest (RF). In [15], the researchers propose a dataset of 4,374 images for object detection of cracks on various masonry materials (cob, brick, stone, and tile). The authors share their dataset publicly on Kaggle and train a YOLOv5 model on it, achieving mAP50 values between 94.4% and 70.3%, depending on the material.

U-Net, introduced by Ronneberger et al. [16], is a seminal deep learning architecture widely adopted for image segmentation. Its "U"-shaped architecture comprises a contracting path (encoder) that progressively down-samples the input image through convolutional layers and pooling operations, capturing contextual information and reducing spatial resolution. The expansive path (decoder) then up-samples these features, gradually restoring spatial resolution and generating a segmentation map. A key innovation of U-Net is to directly concatenate feature maps from the contracting path to the corresponding decoder

layers. This mechanism is crucial for preserving fine-grained spatial information that might otherwise be lost during downsampling, enabling the network to produce high-resolution predictions with precise boundaries.

For the U-Net models, we utilize four prominent CNN backbones as encoders:

- **ResNet-50:** ResNet-50 is a 50-layer deep convolutional neural network known for addressing the vanishing gradient problem through residual blocks and shortcut connections. We use a standard *Torchvision* model with 25.6 M parameters, fine-tuned on images resized to 270x270 pixels.
- **ResNet-101:** ResNet-101 is a deeper variant of ResNet-50, comprising 101 layers (44.5 M parameters). This model is trained on images resized to 540x540 pixels, leveraging its greater capacity for capturing finer details at a larger scale.
- **ConvNeXt V2 Base:** ConvNeXt V2 is a purely convolutional architecture that significantly improves the performance of ConvNets, particularly when co-designed with masked autoencoders for self-supervised learning. It is pre-trained on the popular *ImageNet-22k* dataset and is known for its effective design, balancing model complexity (88.7 M parameters) and performance. We fine-tune this model on images resized to 384x384 pixels.
- **ConvNeXt V2 Huge:** This is the largest variant of the ConvNeXt V2 family (660 M parameters). Like its Base counterpart, it is pre-trained on the *ImageNet-22k* dataset and has been fine-tuned using images resized to 512x512 pixels to maximize detail capture and segmentation precision.

The decoder for all U-Net variants follows the standard U-Net design, incorporating upsampling layers and concatenating features from the respective CNN encoder via skip connections.

### 2.3 Crack Detection in Cultural Heritage and Dataset Challenges

While crack detection is widely studied for civil infrastructures [3,5], its application to CH poses unique challenges due to its peculiarities. Existing datasets for cultural heritage are scarce and often small, making large-scale supervised training difficult. For instance, the Historical-Crack18-19 dataset [4] contains 3886 annotated images from an ancient mosque, but only 757 are crack images, highlighting the imbalance and limited scope.

The scarcity of high-quality, labeled vision data for damaged artifacts is a significant impediment to developing robust deep learning models. This often necessitates data augmentation strategies, including traditional methods like rotations, resizing, and noise addition, or more advanced techniques like synthetic data generation. Synthetic data can augment real-world datasets, improving model performance and generalization, especially for rare defect occurrences. Tools like *UnrealROX* [17] and *EasySynth* [18] can generate realistic synthetic images to overcome training data limitations.

### 3 Methodology

#### 3.1 Problem Definition and Scope

Our objective is to perform semantic segmentation of cracks on the surfaces of statues and monuments. This involves assigning a binary label (crack or non-crack) to each pixel in an input image, thereby precisely delineating the extent and morphology of the damage. The inherent challenges in this task include the fine-grained nature of cracks, their variable morphology (e.g., hairline or structural), the complex and often textured backgrounds of cultural heritage artifacts, and the significant scarcity of annotated ground truth data for this specific domain. We also conduct a comparative study of the ability of the U-Net architectures with various CNN backbones to address these challenges.

#### 3.2 Dataset for Training and Quantitative Evaluation

The primary dataset for our study is OmniCrack30k. OmniCrack30k is a large-scale, systematic, and thorough benchmark dataset specifically designed for universal crack segmentation. It comprises 30,000 samples compiled from over 20 diverse datasets, for a total of 9 billion pixels. This compilation features images of cracks on a wide array of materials, including asphalt, ceramic, concrete, masonry, and steel. The dataset's comprehensive nature aims to reduce biases and enable robust benchmarking for general crack segmentation tasks.

While OmniCrack30k provides extensive data on various cracked materials, it does not specifically include images of statues or monuments. This inherent diversity, however, makes it an excellent foundation for training models that can generalize to different surface textures and crack patterns. To enhance model robustness and generalization, standard data augmentation techniques are also applied during training. The dataset is split into training, validation, and test sets, with the test set used only for the quantitative evaluation of our fine-tuned models.

#### 3.3 Experimental Setup

We fine-tune and evaluate the U-Net architectures with four distinct CNN backbones, representing a range of model capacities as introduced in Section 2.2.

The training took place within a Jupyter Notebook environment running Python 3.10 and using the popular *Fast.ai* library now at its second version [19]. *Fast.ai* adds an additional layer of abstraction above *Pytorch* [20], therefore it is very convenient to use for speeding up the "standard" and repetitive tasks of training a neural network.

**Data augmentation** We train the four architectures both with and without data augmentation. For the data augmentation pipeline, we employ the popular *Albumentations* library. We apply a stochastic augmentation pipeline with probability  $p = 1$ , that combines the following data-augmentation branches:

- geometric transformations (*HorizontalFlip*, *RandomRotate90*, *Transpose*, *ShiftScaleRotate*) at  $p = 0.25$  per operation;
- moderate distortions (*Blur*, *ElasticTransform*, *GridDistortion*, *OpticalDistortion*) are introduced with probability  $p = 0.1$  with limited displacement and distortion intensities;
- photometric variations via *HueSaturationValue* and *CLAHE* ( $p = 0.1$  each) simulate illumination and contrast shifts.

**Optimization and Hyperparameters** All models are trained using the Adam optimizer, known for its adaptive learning rate capabilities, initialized with the default *Fast.ai* parameters. A cosine annealing learning rate schedule is employed, starting with an initial learning rate that ranges between  $1e^{-3}$  and  $1e^{-4}$ . The batch sizes are set to 12 and 8 for the two U-Nets with ResNet backbones and to 24 and 5 for the two U-Nets with ConvNext V2 backbones. For the loss function, we employ the standard **Binary Cross-Entropy (BCE) Loss** that calculates probabilities and compares each actual class output with predicted ones, making it suitable for pixel-level binary classification (crack vs. background).

**Training details** To maximize the level of automation during the training of the network, some *Fast.ai* callbacks have been used to perform the early stopping of the training (with *patience* = 2, i.e. the training stops when the validation loss of the network does not improve for 2 consecutive epochs) and to automatically save the best model of the current training round (according to both validation loss, Dice and Jaccard metrics). Later, that model is reloaded for the final validation phase and to show predicted images on the validation and test set. Experiments were conducted on a workstation equipped with three Nvidia RTX A6000 GPUs (48 Gb of VRAM each), an AMD Ryzen Threadripper Pro 7965WX CPU (with 24 cores/48 threads) and 128 Gb of DDR5 RAM. Each each model was trained on only one GPU at a time, so training times reported in 1 refer to single-GPU training runs.

### 3.4 Evaluation Protocol

Our evaluation protocol consists of two distinct phases to thoroughly assess model performance.

**Quantitative Evaluation on OmniCrack30k Test Set** Upon completion of training, a comparative quantitative evaluation of the different models is performed on the dedicated test set of the OmniCrack30k dataset. The following popular segmentation metrics are calculated: **i) Mean Intersection over Union (mIoU); ii) Dice Coefficient, and iii) Jaccard Index:** [21].

These metrics provide an objective measure of how well each model performs pixel-wise crack segmentation on a diverse, large-scale dataset.

**Out-of-Distribution Qualitative Evaluation** Given the lack of publicly available, annotated datasets specifically for cracked statues and monuments for semantic segmentation, an out-of-distribution qualitative evaluation is performed. This involves applying the trained models to an unlabeled test set comprising real damaged and cracked statues and monuments downloaded from the web. The results are then qualitatively rated by the authors, focusing on: **i) Crack Continuity and Completeness**, that is, how well the models detect entire crack networks; **ii) Boundary Precision** accuracy of crack outlines; **iii) False Positives/Negatives:** instances of over-segmentation (e.g., misclassifying textures as cracks) or under-segmentation (missing actual cracks); **iv) Generalization:** the models’ ability to perform on visually distinct materials and environments not seen during training.

This qualitative assessment is crucial for understanding the practical utility and generalization capabilities of the models in a real-world cultural heritage context.

## 4 Experiments

U-Net Architecture	Train Loss	Val Loss	mIoU	Dice	Jaccard	Time per Epoch (h:mm)
ResNet-50, 270px	0.026	0.027	0.634	0.840	0.755	0:19
ResNet-101, 540px	0.031	0.030	0.624	0.821	0.735	1:06
ConvNeXt V2 Base, 384px	0.022	0.026	0.638	0.848	0.765	0:24
ConvNeXt V2 Huge, 512px	0.021	0.025	0.641	0.859	0.778	5:33

**Table 1.** Quantitative Metrics (best epoch) on OmniCrack30k Training and Validation Sets (no data augmentation regime)

U-Net Architecture	Train Loss	Val Loss	mIoU	Dice	Jaccard	Time per Epoch (h:mm)
ResNet-50, 270px	0.027	0.026	0.626	0.828	0.742	0:20
ResNet-101, 540px	0.034	0.035	0.619	0.808	0.720	1:06
ConvNeXt V2 Base, 384px	0.033	0.030	0.626	0.822	0.736	0:24
ConvNeXt V2 Huge, 512px	0.028	0.026	0.636	0.851	0.768	5:33

**Table 2.** Quantitative Metrics (best epoch) on OmniCrack30k Training and Validation Sets (with data augmentation regime)

### 4.1 Quantitative Evaluation on OmniCrack30k Test Set

The quantitative results for the four U-Net architectures, trained and evaluated on the OmniCrack30k dataset, are summarized in Table 1 to Table 4. These metrics provide insights into the models’ learning efficiency and segmentation accuracy.

<b>U-Net Architecture</b>	<b>Test Loss</b>	<b>mIoU</b>	<b>Dice</b>	<b>Jaccard</b>
ResNet-50, 270px	0.024	0.652	0.853	0.771
ResNet-101, 540px	0.029	0.645	0.831	0.745
ConvNeXt V2 Base, 384px	0.026	0.662	0.852	0.770
ConvNeXt V2 Huge, 512px	0.024	0.666	0.865	0.786

**Table 3.** Quantitative Metrics (best epoch) on OmniCrack30k Test Set (no data augmentation regime)

<b>U-Net Architecture</b>	<b>Test Loss</b>	<b>mIoU</b>	<b>Dice</b>	<b>Jaccard</b>
ResNet-50, 270px	0.025	0.651	0.841	0.757
ResNet-101, 540px	0.034	0.645	0.822	0.734
ConvNeXt V2 Base, 384px	0.030	0.647	0.830	0.744
ConvNeXt V2 Huge, 512px	0.024	0.659	0.862	0.782

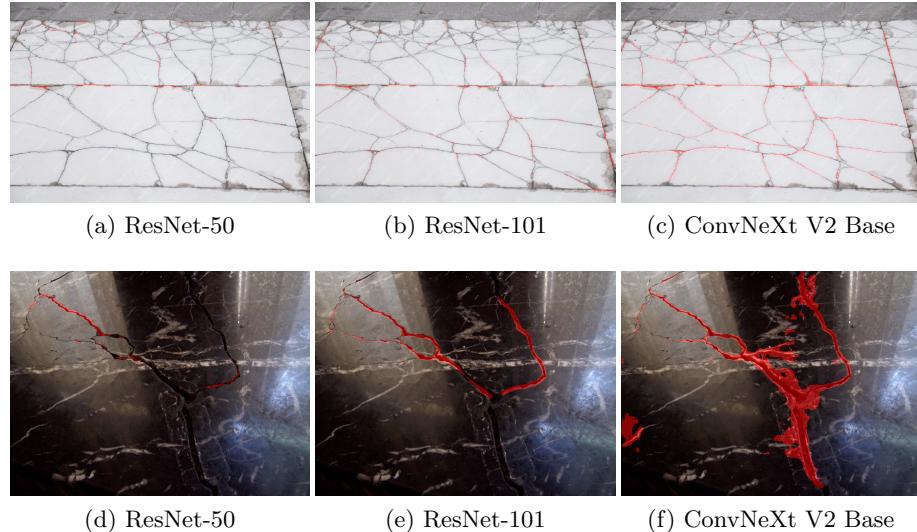
**Table 4.** Quantitative Metrics (best epoch) on OmniCrack30k Test Set (with data augmentation regime)

### Analysis of Results

The quantitative metrics reveal several key patterns. First, ConvNeXt V2 Huge consistently achieves the highest segmentation accuracy across all configurations (e.g., test mIoU of 0.666 without augmentation), validating its architectural superiority for artifacts crack detection. However, this comes at substantial computational cost, with 5× longer training times vs. ResNet-101. Notably, the data augmentation pipeline, as it has been conceived, *reduces* the performance for most models – e.g. ConvNeXt V2 Base test mIoU drops from 0.662 to 0.647 – suggesting possible over-regularization on this dataset, eventually causing the models to “*hallucinate*” cracks (false positives).



**Fig. 1.** Out-of-distribution predictions obtained with the *ConvNeXt V2 Huge* U-Net model (no data augmentation regime) on images of black and white marble (therefore closer to images in the training set).



**Fig. 2.** The same images as in Fig. 1 processed with the other three fine-tuned models (no data augmentation regime as before).

#### 4.2 Out-of-Distribution Qualitative Evaluation

The qualitative evaluation on unlabeled images of real damaged and cracked statues and monuments provides crucial insights into the models’ generalization capabilities to a domain not explicitly seen during training. Despite the training data (OmniCrack30k) consisting of cracks on materials like asphalt, ceramic, concrete, masonry, and steel, and not actual statues or monuments, the models demonstrate a remarkable generalization ability and are able to segment cracks in this new context.

Specifically, as shown in Fig. 1, the U-Net utilizing ConvNeXt V2 Huge exhibits superior performance in accurately segmenting fractures across materials with different colors and textures. In contrast, other architectures - as evidenced by Fig. 2 - fail to achieve comparable results. This limitation is particularly pronounced in models employing ResNet backbones. ConvNeXt V2 Base, instead, is prone to false positives segmenting black marble images, probably due to its lower image resolution and the limited exposure to such samples during training.

As illustrated in Fig. 3, ConvNeXt V2 Huge effectively identifies cracks in images of statues. Nevertheless, it occasionally produces false positives in the form of overflowing and imprecise boundaries.

In contrast, ResNet-based backbones exhibit some limitations, as evidenced in Fig. 4. These architectures struggle to reliably detect cracks, frequently yielding both false positives (e.g. Fig. 4h, both the eye and the mouth are incorrectly highlighted) and false negatives (e.g. Fig. 4k). Notably, ResNet-50 fails to detect entire cracks in several cases. While ConvNeXt V2 Base achieves reasonable seg-



**Fig. 3.** Out-of-distribution predictions obtained with the *ConvNeXt V2 Huge* U-Net model (no data augmentation regime) on images depicting statues (therefore quite distant from images in the training set).

mentation performance across most statues, its predictions degrade substantially on dark-colored materials, showing pronounced false-positive artifacts.

Overall, the qualitative evaluation confirms the quantitative findings: models with more advanced and larger CNN backbones, particularly the ConvNeXt V2 variants, show better generalization capabilities for crack segmentation on real-world statues and artifacts. This is a non-trivial and non-obvious finding, as it

suggests that large convolutional architectures trained on diverse crack datasets can effectively transfer their learned knowledge to new, visually distinct domains without any explicit domain-specific training, incurring in minimal false-positive rates.

## 5 Conclusion

This study presented a comparative analysis of U-Net architectures with various CNN encoders for the semantic segmentation of cracks, a critical task for cultural heritage preservation. The research leveraged the OmniCrack30k dataset for training and quantitative evaluation, acknowledging the significant challenge posed by the lack of large-scale publicly available datasets specifically annotated for the semantic segmentation of cracked statues and artifacts.

Our methodology involved training U-Net models with ResNet-50, ResNet-101, ConvNeXt V2 Base and ConvNeXt V2 Huge as encoders, at resolutions of 270px, 540px, 384px and 512px respectively. The quantitative evaluation on the OmniCrack30k test set demonstrated a clear performance hierarchy, with the ConvNeXt V2 Huge backbone achieving the highest mIoU, Dice, and Jaccard scores, followed by ConvNeXt V2 Base, ResNet-101, and ResNet-50. This highlights the benefits of increased model capacity, higher input resolution, and advanced CNN architectures like ConvNeXt V2, which benefit from masked autoencoder pre-training and for being trained on *ImageNet-22k* (thus on all 14 million images).

Moreover, the out-of-distribution qualitative evaluation on images of real damaged and cracked statues and artifacts downloaded from the web revealed a promising generalization capability. The ConvNeXt V2-based U-Nets, in particular, demonstrated superior performance in detecting fine hairline cracks, maintaining continuity, and exhibiting robustness to complex textures, materials and lighting conditions inherent in cultural heritage images. This suggests that the learned features from diverse crack patterns are highly transferable.

On the other hand, while the models already show strong generalization capabilities, fine-tuning on a dedicated, albeit small, dataset of cultural heritage cracks could further enhance their performance and reduce out-of-domain false positives.

Future research should focus primarily on the creation of large-scale, publicly available datasets for statue and monument segmentation. Beyond manual annotation, this objective can be achieved combining image segmentation with 3D reconstruction, with synthetic 3D data generation or leveraging diffusion models for domain adaptation (e.g. adapting human-centric datasets to the statue domain while preserving semantic boundaries). Moreover, future efforts should involve the technical expertise of conservators and archaeologists, potentially for labeling future datasets, providing qualitative assessment of model performance, and, most importantly, identifying underrepresented artifact types (e.g. beyond statues and monuments) critical to cultural heritage preservation. In conclusion, this research represents a significant step towards leveraging advanced

deep learning for the automated preservation of our invaluable global cultural heritage.

## Acknowledgements

This work was carried out within the framework of the Italian Ministry of Business and Made in Italy (MIMIT) project, House of Emerging Technologies Genoa (CTEGE): Digital factory for culture.

This work was also partially funded within the framework of the activities of the National Recovery and Resilience Plan (NRRP) M4C2 Inv. 1.4 – CN MOST – Sustainable Mobility Center, Spoke 7, whose financial support is gratefully acknowledged.

## References

- Christian Benz and Volker Rodehorst. Omnicrack30k: A benchmark for crack segmentation and the reasonable effectiveness of transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3876–3886, 2024.
- Moez Krichen. Convolutional neural networks: A survey. *Computers*, 12(8):151, 2023.
- Elia Moscoso Thompson, Andrea Ranieri, Silvia Biasotti, Miguel Chicchon, Ivan Sipiran, Minh-Khoi Pham, Thang-Long Nguyen-Ho, Hai-Dang Nguyen, and Minh-Triet Tran. Shrec 2022: Pothole and crack detection in the road pavement using images and rgb-d data. *Computers & Graphics*, 107:161–171, 2022.
- Esraa Elhariri, Nashwa El-Bendary, and Shereen A Taie. Historical-crack18-19: A dataset of annotated images for non-invasive surface crack detection in historical buildings. *Data in Brief*, 41:107865, 2022.
- Rui Fan, Mohammad Junaid Bocus, Yilong Zhu, Jianhao Jiao, Li Wang, Fulong Ma, Shanshan Cheng, and Ming Liu. Road crack detection using deep convolutional neural network and adaptive thresholding. In *2019 IEEE Intelligent Vehicles Symposium (IV)*, pages 474–479. IEEE, 2019.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16133–16142, 2023.
- Anders Landstrom and Matthew J Thurley. Morphology-based crack detection for steel slabs. *IEEE Journal of selected topics in signal processing*, 6(7):866–875, 2012.

10. Qiang Zhou, Zhong Qu, and Chong Cao. Mixed pooling and richer attention feature fusion for crack detection. *Pattern Recognition Letters*, 145:96–102, 2021.
11. Ankur Dixit and Hiroaki Wagatsuma. Investigating the effectiveness of the sobel operator in the mca-based automatic crack detection. In *2018 4th International Conference on Optimization and Applications (ICOA)*, pages 1–6. IEEE, 2018.
12. Yuslena Sari, Puguh Budi Prakoso, and Andreyan Rezky Baskara. Road crack detection using support vector machine (svm) and otsu algorithm. In *2019 6th International Conference on Electric Vehicular Technology (ICEVT)*, pages 349–354. IEEE, 2019.
13. Liu Zhen-Liang, Zhou An, Ran Xin-Ru, Wu Yun-Peng, Zhao Wei-Gang, and Zhang Hao. A crack detection and quantification method using matched filter and photograph reconstruction. *Scientific Reports*, 15(1):25266, 2025.
14. Krisada Chaiyasarn, Mayank Sharma, Luqman Ali, Wasif Khan, and Nakhon Poovarodom. Crack detection in historical structures based on convolutional neural network. *Geomatè Journal*, 15(51):240–251, 2018.
15. Narges Karimi, Mayank Mishra, and Paulo B Lourenço. Automated surface crack detection in historical constructions with various materials using deep learning-based yolo network. *International Journal of Architectural Heritage*, 19(5):581–597, 2025.
16. Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages 234–241, Cham, 2015. Springer International Publishing.
17. Pablo Martinez-Gonzalez, Sergiu Oprea, Alberto Garcia-Garcia, Alvaro Jover-Alvarez, Sergio Orts-Escalano, and Jose Garcia-Rodriguez. Unrealrox: an extremely photorealistic virtual reality environment for robotics simulations and synthetic data generation. *Virtual Reality*, 24:271–288, 2020.
18. Nikola Jovanović and Filip Panjević. Easysynth - a plugin for unreal engine. <https://github.com/ydrive/EasySynth>, 2025.
19. Jeremy Howard and Sylvain Gugger. Fastai: A layered API for deep learning. *Information*, 11(2):108, Feb 2020.
20. Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019.
21. Michel Marie Deza and Elena Deza. *Encyclopedia of Distances*. Springer Berlin Heidelberg, 2009.

## A Appendix

This appendix presents an ablation study evaluating the impact of individual data augmentation techniques on the performance of our smallest architecture, a *U-Net* with a *ResNet-50* encoder trained with images at *270px* resolution. Each model was trained on the OmniCrack30k dataset using a single augmentation transform, with validation and test metrics reported in Tables 5 and 6, respectively. Key observations are summarized below.

Augmentation	Valid Loss	mIoU	Dice	Jaccard
Transpose	0.025	0.643	0.854	0.772
CLAHE	0.026	0.643	0.851	0.769
GridDistortion	0.026	0.646	0.849	0.766
ElasticTransform	0.027	0.641	0.845	0.761
RandomRotate90	0.028	0.633	0.843	0.759
OpticalDistortion	0.027	0.633	0.841	0.756
Blur	0.028	0.635	0.839	0.755
HorizontalFlip	0.028	0.631	0.837	0.752
ShiftScaleRotate	0.028	0.628	0.834	0.749
HueSaturationValue	0.033	0.625	0.821	0.734

**Table 5.** Ablation study on data augmentation techniques: validation metrics sorted by Dice score for ten different *ResNet-50 @270px* models trained with a single *Albumentation* transform type each.

Augmentation	Test Loss	mIoU	Dice	Jaccard
Transpose	0.023	0.675	0.866	0.787
GridDistortion	0.024	0.672	0.863	0.783
CLAHE	0.025	0.668	0.861	0.780
ElasticTransform	0.025	0.669	0.856	0.775
RandomRotate90	0.026	0.660	0.854	0.772
Blur	0.026	0.661	0.854	0.772
OpticalDistortion	0.025	0.658	0.852	0.770
HorizontalFlip	0.025	0.655	0.850	0.768
ShiftScaleRotate	0.026	0.655	0.842	0.758
HueSaturationValue	0.034	0.650	0.834	0.749

**Table 6.** Ablation study on data augmentation techniques: test metrics sorted by Dice score for ten different *ResNet-50 @270px* models trained with a single *Albumentation* transform type each.

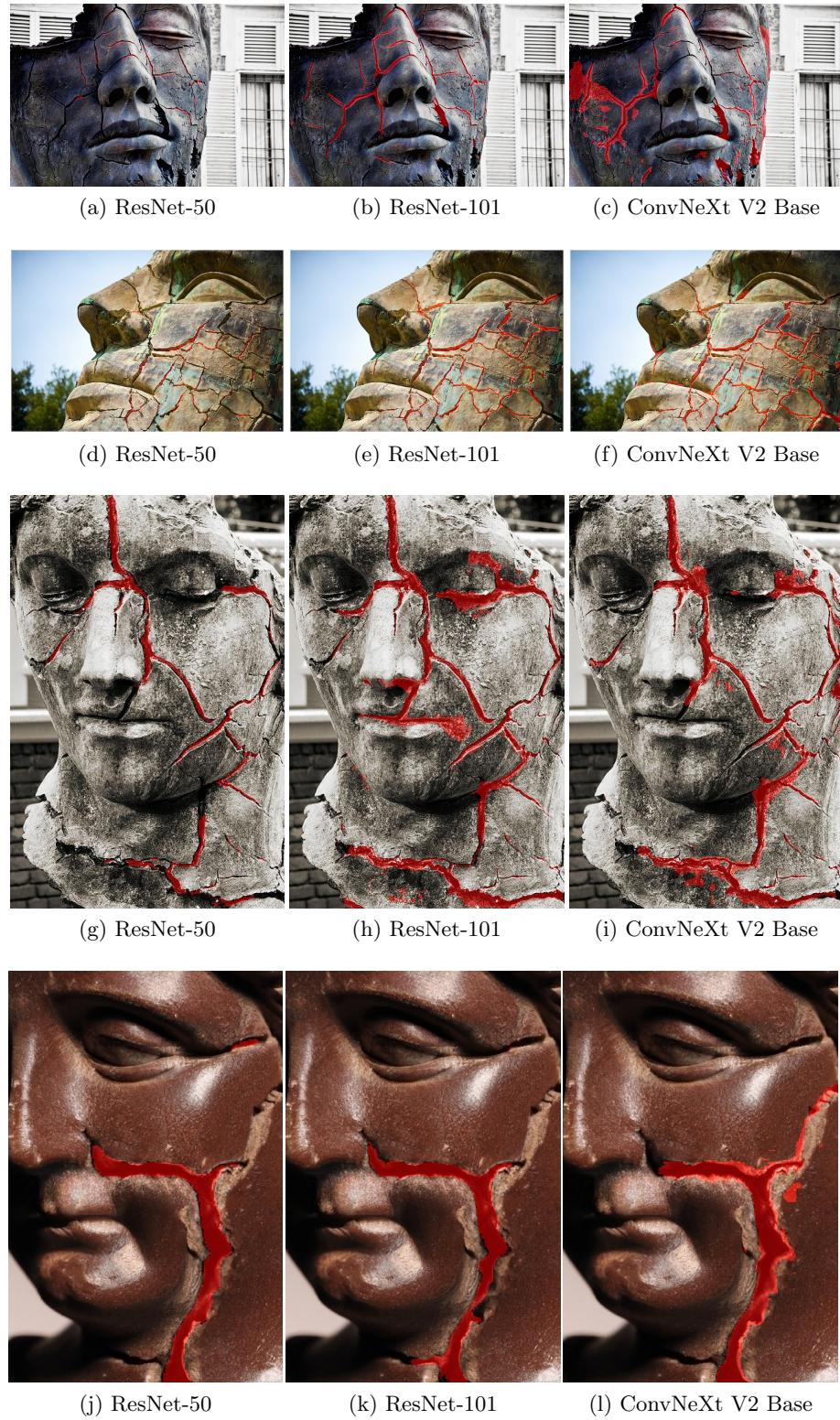
The ablation study shows that top-performing augmentations, selected by their *Dice* score, are *Transpose*, *CLAHE*, and *GridDistortion* with *Dice* scores ranging from 0.854 to 0.849 on the validation set and 0.866–0.861 on the test set<sup>2</sup>. These techniques likely enhance crack detection for the following reasons:

<sup>2</sup> Slightly higher *Dice* score values in the test set compared to those in the validation set are consistent with what was observed in Tables 1 and 3, probably due to the intrinsic conformation of the dataset.

- **Transpose** preserves structural patterns while altering orientation in  $90^\circ$  steps, improving robustness to spatial variations.
- **CLAHE** (Contrast Limited Adaptive Histogram Equalization) enhances local contrast in low-light regions, amplifying subtle crack features without introducing further noise.
- **GridDistortion** simulates natural surface deformations (e.g., material warping), aiding generalization to unseen images.

On the other hand, the most detrimental augmentations are *HueSaturationValue* and *ShiftScaleRotate* with Dice scores ranging from 0.821 to 0.834 on the validation set and 0.834-0.842 on the test set. Potential reasons include:

- **HueSaturationValue** distorts color channels, obscuring grayscale crack features critical for detection.
- **ShiftScaleRotate** performs translations, scaling and rotations, the combination of which can erase the finest cracks from the image or make them too blurred (due to rotation, in fact, *Blur*, understandably, has similar impacts on segmentation).



**Fig. 4.** The same images as in Fig. 3 processed with the other three fine-tuned models (no data augmentation regime as before).