# Weakly Supervised Semantic Segmentation Using Visual Explainability

**A DISSERTATION**

*submitted in partial fulfillment of the requirements*

*for the award of the degree of*

**Master of Technology (Computer Science)**

by

**Niladri Das**

**ROLL NO. CS2219**

Under the supervision of

**Ujjwal Bhattacharya**

Computer Vision and Pattern Recognition Unit



**INDIAN STATISTICAL INSTITUTE**

# CERTIFICATE

I hereby certify that the work which is being presented in the M.Tech. Dissertation entitled **Weakly Supervised Semantic Segmentation Using Visual Explainability** in partial fulfillment of the requirements for the award of the **Master of Technology (Computer Science)** is an authentic record of my own work carried out during a period from June, 2023 to June,2024 under the supervision of **Ujjwal Bhattacharya, Computer Vision and Pattern Recognition Unit, Indian Statistical Institute**.

The matter presented in this thesis has not been submitted for the award of any other degree elsewhere.

*Signature of Candidate*

**Niladri Das**

**Roll No. CS2219**

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

*Signature of Supervisor*

Date: 12 June 2024

**Ujjwal Bhattacharya**

**Computer Vision and Pattern Recognition Unit**

**Indian Statistical Institute**

# ACKNOWLEDGEMENT

I am profoundly grateful for the divine inspiration that fueled my passion and dedication, culminating in the successful completion of this research. I owe a debt of gratitude to my mentor, Ujjwal Bhattacharya from the Computer Vision and Pattern Recognition Unit at the Indian Statistical Institute. His unwavering support and insightful guidance were instrumental in deepening my understanding of the subject matter and contributed significantly to the progression and excellence of my M. Tech. dissertation.

The creation of this thesis was made possible by the unwavering encouragement of my peers. I extend my heartfelt appreciation to my parents for their endless love, support, and blessings. Finally, I offer my sincerest thanks to the divine for endowing me with the determination, fortitude, and insight necessary to fulfill this academic endeavor.

(Niladri Das)

# ABSTRACT

In the realm of remote sensing, the task of semantically segmenting landslide images is traditionally reliant on supervised learning techniques. These methods necessitate extensive training datasets and meticulous pixel-level annotations—a process that demands considerable human labor and incurs high costs. To mitigate these challenges, we introduce an innovative approach that employs weakly supervised learning, integrating Class Activation Maps (CAMs) with a Cycle Generative Adversarial Network (CycleGAN). This novel methodology leverages image-level labels in lieu of pixel-level annotations. Initially, CAMs are utilized to locate the landslide's rough area. Subsequently, CycleGAN generates a synthetic image devoid of landslides, which, when contrasted with the original, yields precise segmentation results. The efficacy of our approach is quantified using the mean Intersection-over-Union (mIOU) metric, demonstrating a superior performance with an mIOU of 0.228. Additionally, when juxtaposed with a U-Net network's supervised learning technique, which scored an mIOU of 0.408, our results affirm the viability of weakly supervised learning for accurate landslide semantic segmentation in remote sensing imagery. This method significantly alleviates the burden of data annotation.

Incorporating the advancements of Score-CAM [6], which surpasses Grad-CAM in object localization accuracy, we further refine our model. Score-CAM's enhanced precision in identifying relevant features contributes to the improved segmentation of landslide areas, promising a new frontier in remote sensing image analysis.

# Contents

# List of Figures

# Chapter 1

# INTRODUCTION

## 1.1   Introduction

In the evolving landscape of remote sensing image analysis, the precise segmentation of landslide areas remains a critical challenge. The comprehensive, high-quality data annotations required for traditional methods, which mostly rely on supervised learning, are expensive and labor-intensive. Recognizing these limitations, pioneering researchers Zhou, Y., Wang, H., Yang, R., Yao, G., Xu, Q., and Zhang, X. [1], have explored the potential of combining Cycle Generative Adversarial Networks (Cycle-GAN) with Gradient-weighted Class Activation Mapping (Grad-CAM [5]) to enhance segmentation accuracy. However, Grad-CAM's susceptibility to saturation and false confidence issues necessitates a more robust solution.

Enter Score-CAM [6], an innovative alternative proposed by Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu, designed to overcome the inherent drawbacks of gradient-based methods. Building upon this foundation, my thesis introduces a novel integration of Cycle-GAN with Score-CAM, aiming to harness the strengths of both techniques for the semantic segmentation of landslides in remote sensing images. This approach not only addresses the challenges posed by gradient-based methods but also significantly reduces the dependency on pixel-level annotations, thereby curtailing the annotation workload.

This thesis delves into the theoretical underpinnings of Cycle-GAN and Score-CAM, elucidating their respective roles in the proposed segmentation framework. Through a series of experiments and evaluations, it demonstrates the enhanced performance and reliability of the combined approach, offering a promising direction for future research in the field. The introduction of Score-CAM, in particular, marks a significant advancement in object localization, ensuring that the segmentation process is both accurate and efficient.

As we stand on the cusp of a new era in remote sensing technology, this thesis contributes to the ongoing discourse on machine learning methodologies, presenting a compelling case for the integration of weakly supervised learning techniques in environmental monitoring and disaster management applications. The insights gleaned from this research endeavor not only pave the way for more sophisticated analytical tools but also underscore the importance of innovation in addressing the pressing challenges of our time.

## 1.2    Limitations of Grad-CAM

1. **Gradient Saturation [6]:** The gradients in deep neural networks can become noisy and may vanish, especially under the influence of sigmoid saturation or within the zero-gradient regions of ReLU functions. This leads to gradients that are visually noisy, complicating the interpretation of Saliency Maps which rely on clear gradient signals to identify salient features.

2. **False Confidence [6]:** $A_i^l$ and $A_j^l$ are examples of linear combinations of activation maps created by Grad-CAM. It is expected that if $\alpha_c^i \geq \alpha_c^j$, then the region corresponding to $A_i^l$ for the target class $c$ is at least as significant as the one corresponding to $A_j^l$. This is the case when the weights $\alpha_c^i$ and $\alpha_c^j$ are assigned to these maps. On the other hand, when activation maps with larger weights contribute less to the network's output in comparison to a zero baseline, this can result in erroneous confidence. The gradient vanishing problem within the network and global pooling on gradients may make this problem worse.

## 1.3  Semantic Segmentation of Landslides

### 1.3.1  Score-CAM

Based on perturbation, Score-CAM [6] assesses the increase in confidence to determine the significance of activation maps. Below is a detailed description of the Score-CAM method.



Figure 1.1: Score-CAM Pipeline

**Definition 1 (Increase of Confidence):** Given a function $Y = f(X)$ that takes an input vector $X = [x_0, x_1, \ldots, x_n]^T$ and outputs a scalar $Y$. For a known baseline input $X_b$, the contribution $c_i$ of $x_i$ towards $Y$ is the change in output when the $i$-th entry in $X_b$ is replaced with $x_i$. Formally,

$$c_i = f(X_b \circ H_i) - f(X_b) \tag{1.1}$$

where $H_i$ is a vector with the same shape as $X_b$, but each entry $h_j$ in $H_i$ is defined as $h_j = I[i = j]$ and $\circ$ denotes the Hadamard product.

**Definition 2 (Channel-wise Increase of Confidence - CIC):** Given a CNN model $Y = f(X)$ that takes an input $X$ and outputs a scalar $Y$. For an internal convolutional layer $l$ with corresponding activation $A$, the contribution of the $k$-th channel $A_l^k$ towards $Y$ is defined as

$$C(A_l^k) = f(X \circ H_l^k) - f(X_b) \tag{1.2}$$

3

where

$$H_l^k = s(\text{Up}(A_l^k)) \tag{1.3}$$

Here, $\text{Up}(\cdot)$ denotes the operation that upsamples $A_l^k$ to the input size, and $s(\cdot)$ is a normalization function that maps each element in the input matrix to $[0, 1]$.

To generate a smoother mask $H_l^k$, the raw activation values in each activation map are normalized using:

$$s(A_l^k) = \frac{A_l^k - \min A_l^k}{\max A_l^k - \min A_l^k} \tag{1.4}$$

**Definition 3 (Score-CAM):** For a convolutional layer $l$ in a model $f$, given a class of interest $c$, the Score-CAM $L_{\text{Score-CAM}}^c$ is defined as:

$$L_{\text{Score-CAM}}^c = \text{ReLU}\left(\sum_k \alpha_k^c A_l^k\right) \tag{1.5}$$

where

$$\alpha_k^c = C(A_l^k) \tag{1.6}$$

The ReLU function is applied to the linear combination of activation maps because only the features that positively influence the class of interest are considered. The weights are derived from the CIC score for the corresponding activation maps on the target class, thus removing the dependence on gradients. Although the last convolutional layer is typically preferred, any intermediate convolutional layer can be utilized in this framework.

## 1.3.2 Cycle-GAN

Generative Adversarial Networks (GANs) [7] are a class of machine learning frameworks where two neural networks, a generator and a discriminator, are trained simultaneously through adversarial processes. Originally, GANs were utilized to generate images that mimic a specific style or distribution.

Cycle-Consistent Generative Adversarial Networks (CycleGANs) [10] expand on the capabilities of traditional GANs by enabling the translation of images from one domain to another without requiring paired examples. This means that CycleGANs can learn to transform images from one style to another while preserving the essential content of the original image.

Figure 1.2: Flow-Chart of Cycle-GAN

In our work, we leverage CycleGANs for style transfer between two distinct image domains: landslide images and non-landslide images. The fundamental idea is to treat landslide-affected images as one domain and images without landslides as another. By training a generative network, we can convert landslide images to their non-landslide counterparts, effectively generating virtual non-landslide images from the landslide images.

CycleGAN has a number of advantages over pix2pix GAN, including the ability to train without the need for paired pictures from both domains. This is particularly beneficial in scenarios where obtaining such pairs is challenging or impractical. In our case, we can utilize separate datasets of landslide and non-landslide images, eliminating the need for corresponding pairs of images captured before and after a landslide event.

The primary components of CycleGAN include two generators and two discriminators:

- **Generator** $G$: This generator learns to convert images from domain $X$ (landslide images) to domain $Y$ (non-landslide images).

- **Generator** $F$: This generator learns to convert images from domain $Y$ (non-

landslide images) to domain $X$ (landslide images).

- **Discriminator $D_X$**: This discriminator's role is to differentiate between real images from domain $X$ and fake images generated by $F$ (i.e., images that are translated from domain $Y$ to domain $X$).

- **Discriminator $D_Y$**: This discriminator's role is to distinguish between real images from domain $Y$ and fake images generated by $G$ (i.e., images that are translated from domain $X$ to domain $Y$).

CycleGAN utilizes a technique known as cycle consistency to ensure that the translations are meaningful and the core content of the images is preserved. An picture from one domain must closely resemble the original image when translated to the other domain and back again due to the cycle consistency loss. Formally, for an image $x$ in domain $X$ and an image $y$ in domain $Y$:

- Forward cycle consistency loss: $\|F(G(x)) - x\|$

- Backward cycle consistency loss: $\|G(F(y)) - y\|$

By minimizing these cycle consistency losses along with the adversarial losses, CycleGAN ensures that the generated images are not only stylistically accurate but also retain the essential features of the original images. This capability is crucial for our application, as it allows us to generate realistic non-landslide images from landslide images, thereby facilitating effective analysis and interpretation.

### 1.3.3 Proposed Approach

Our approach involves utilizing the Score-CAM method, which is a perturbation-based technique for visual explanations of CNN models. The steps involved in our approach are as follows:

1. Generate a non-landslide image from a given landslide image using CycleGAN.

2. Determine the difference between the photographs with and without landslides.

3. Map the difference to grayscale.

4. Generate the Score-CAM heatmap.

5. Apply thresholding to the heatmap to identify regions of interest.

6. Perform segmentation by taking the intersection of the thresholded heatmap with the original image.

## 1.4 Motivation

The primary motivation for this work stems from the need to enhance accuracy in image segmentation tasks using the Score-CAM method. Image segmentation, particularly in domains such as remote sensing and medical imaging, often relies on pixel-wise annotated datasets for training deep learning models. However, creating these annotated datasets is an extremely labor-intensive and costly process, requiring significant manual effort from experts.

Our proposed method addresses this challenge by utilizing Score-CAM, a perturbation -based technique that generates visual explanations for CNN models. Score-CAM enables us to identify and highlight regions of interest in images, facilitating more precise segmentation without the need for exhaustive manual annotations.

In the context of landslide detection, we leverage CycleGAN to transform landslide images into non-landslide images. This transformation allows us to create virtual non-landslide images, which serve as a reference to identify changes and segment the landslide areas effectively. By computing the difference between the landslide and non-landslide images and mapping these differences to a grayscale format, we can generate Score-CAM heatmaps that pinpoint the regions affected by landslides.

The intersection of these heatmaps with appropriate thresholds yields accurate segmentation results. This method not only enhances segmentation accuracy but also significantly reduces the dependency on manually annotated datasets, making it a cost-effective solution.

Beyond landslide detection, the versatility of this approach makes it applicable to various other fields. For instance, in medical imaging, precise segmentation of anatomical structures or pathological regions is crucial for diagnosis and treatment planning. Our method can be adapted to segment medical images, thereby

7

improving diagnostic accuracy and patient outcomes.

Furthermore, the cost-effectiveness of our approach is particularly beneficial for large-scale applications where obtaining annotated datasets is not feasible. By minimizing the need for manual annotations, we can streamline the segmentation process and make it more accessible for different research and industrial applications.

The potential to apply this method across different domains highlights its robustness and adaptability. Whether it's identifying tumors in medical scans, detecting changes in satellite imagery for environmental monitoring, or segmenting objects in various industrial applications, the underlying principles of our method remain applicable.

In summary, the dual goals of increasing segmentation accuracy and lowering the expense of producing annotated datasets are what motivate this effort. By harnessing the power of Score-CAM and CycleGAN, we provide a scalable and efficient solution that can be extended to a wide range of applications, thereby addressing a critical bottleneck in the field of image analysis and segmentation.

**Chapter 2** delves into the Literature Survey, **Chapter 3** outlines the Thesis and Methodology, **Chapter 4** showcases the Results, **Chapter 5** provides the Conclusion, and the **final chapter** explores Future Work.

# Chapter 2

# LITERATURE SURVEY

## 2.1 Introduction

This chapter provides a detailed review of the significant works in the field of image segmentation using deep learning techniques, particularly focusing on CycleGAN and Score-CAM methods. The advancements in landslide detection and segmentation, as well as the development of visual explanation methods for deep neural networks, are thoroughly examined. The motivation behind integrating CycleGAN and Score-CAM for enhanced segmentation accuracy is also discussed.

## 2.2 Landslide Detection and Segmentation

Since the introduction of the ImageNet dataset in 2009, the field of computer vision has witnessed rapid advancements driven by deep learning techniques. These advancements have significantly impacted image classification, detection, and segmentation tasks. As a result, in order to identify and categorize landslide regions in remote sensing photos, researchers have started investigating deep learning techniques. The application of deep learning for landslide identification and segmentation has received less attention than it has for other remote sensing image recognition problems.

Sameen et al. employed a deep residual detection approach based on a feature

fusion network to detect landslides in remote sensing images from the Kinmallan plateau in Malaysia. Their method demonstrated a notable improvement, with an increase of 0.13 in the F1 score and 0.1296 in the mean Intersection over Union (mIOU) compared to conventional convolutional layer stacking techniques. Similarly, Chen et al. proposed a change detection method leveraging a deep convolutional network, achieving a false recognition rate of 0.176. Cheng et al. developed the YOLO-SA model for landslide detection in Qiaojia and Ludian counties in Yunnan Province, China, achieving a recognition accuracy of 0.9408.

Segmentation of landslide areas in remote sensing images is crucial for precisely delineating boundaries, studying changes in landslide areas, and calculating the affected regions. However, the scarcity of datasets has limited research in remote sensing image segmentation tasks. Soares et al. utilized Digital Elevation Model (DEM) information as training data and employed the U-Net model to automatically segment landslides in Novo Fribourg, Brazil, achieving F1 scores of 0.55 and 0.58 on two different test sets. Du et al. compared six deep learning semantic segmentation models using a self-built Yangtze River coastal landslide dataset, with the GCN and DeepLabV3 models achieving mIOU accuracies of 0.542 and 0.740, respectively. Prakash et al. experimented with an improved U-Net network on a statewide landslide dataset in Oregon, achieving a detection rate of 0.72, outperforming traditional methods. Bo et al. applied deep learning semantic segmentation techniques to accurately detect landslide areas in remote sensing images of Nepal, achieving a recall rate of 0.65 and an accuracy rate of 0.55.

## 2.3   Visual Explanation Methods

Enhancing the transparency of deep neural networks (DNNs) by making some aspects of their inference interpretable by humans has become an important research focus. Among the various explanation methods, visualizing the importance of input features or learned weights is the most straightforward approach. Many techniques have been created to enhance the explanations of convolutions and Convolutional Neural Networks (CNNs), as spatial convolution is a common element in state-of-the-art models for both language and image processing.
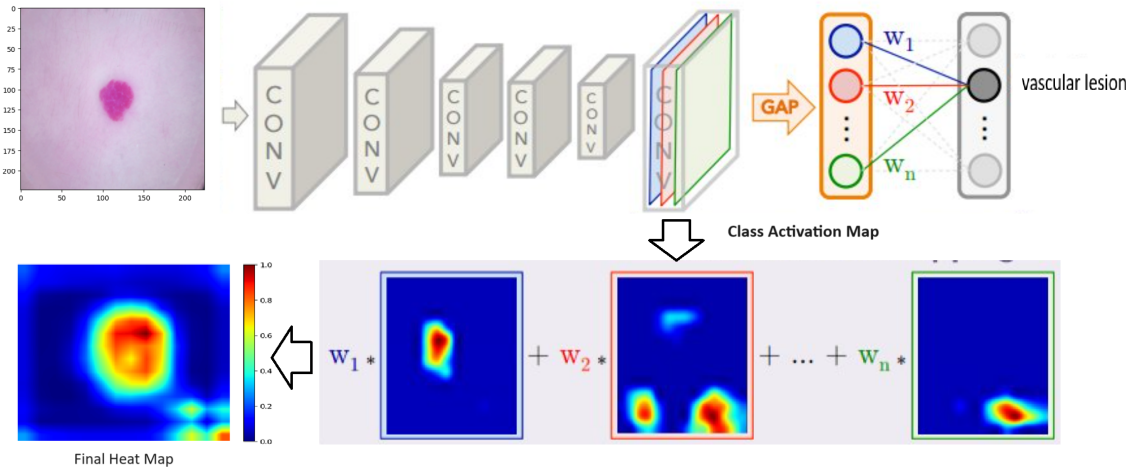
Figure 2.1: Class Activation Map Pipeline

Gradient-based methods [5] highlight image regions that influence predictions by backpropagating the gradient of a target class to the input. Saliency Map uses the derivative of the target class score with respect to the input image as an explanation. Other works manipulate this gradient to enhance the visual sharpness of the results. However, these gradient-based maps often suffer from low quality and noise.

Perturbation-based approaches perturb the original input and observe changes in model predictions to identify minimal regions influencing the output. These methods generally require additional regularization to improve results and can be computationally intensive.

Class Activation Map (CAM)-based [5] explanations provide visual explanations by creating a linear weighted combination of activation maps from convolutional layers. CAM generates localized visual explanations but requires a global pooling layer, making it architecture-sensitive. Grad-CAM and its variations, such as Grad-CAM++, generalize CAM to models without global pooling layers.

To address the limitations of gradient-based CAM variations, Score-CAM introduces a novel gradient-free visual explanation method.Rather than depending on gradient information, Score-CAM determines the significance of activation maps based on the contribution of their highlighted input features to the model output. This method provides a more accurate and intuitive explanation by bridging the gap between perturbation-based and CAM-based approaches.

## 2.4 Combining CAM and CycleGAN

An important study that integrates CAM and CycleGAN is the work by Zhou et al. [1]. They proposed a weakly supervised method for semantic segmentation of landslides in remote sensing images. In their approach, CycleGAN was employed for style transfer to generate non-landslide images from landslide images, while Grad-CAM was used for visual explanations to highlight the regions of interest. However, the segmentations produced by Grad-CAM were often noisy, resulting in less accurate delineation of landslide areas. They constructed masks by applying thresholding to the Grad-CAM heatmaps and then took the intersection with the CycleGAN outputs to improve them in order to address this.



Figure 2.2: Fully Combined Work Flow

In contrast, our approach utilizes Score-CAM instead of Grad-CAM. Score-CAM provides a more refined and precise segmentation due to its gradient-free mechanism, which reduces noise and enhances the clarity of the highlighted regions. By integrating CycleGAN for style transfer and Score-CAM for generating accurate heatmaps, our method achieves cleaner and more reliable segmentation results. This improvement is crucial for applications requiring high precision, such as monitoring and analyzing landslide-prone areas.

# Chapter 3

# Methodology

## 3.1 Dataset Description

The dataset for this study is derived from the Bijie landslide dataset [1], covering the entire Bijie city, which spans an area of 26,853 km² in the northwest of Guizhou province, China (Fig. 1). This region is situated in a transitional slope zone from the Tibet Plateau to the eastern hills, with altitudes ranging from 457 to 2900 meters. This region is one of the most landslide-prone in China because to its unstable geological features, many steep hillsides, significant yearly rainfall (between 849 and 1399 mm), and delicate biological environment.

The Bijie city area is prone to various types of landslides, including rock falls, rock slides, and, to a lesser extent, debris slides. Each year, numerous new landslides occur, causing significant damage to human settlements, infrastructure such as roads, bridges, and transmission lines, and agricultural lands. Traditionally, landslides in this region are detected through a combination of methods. One common approach involves indoor manual interpretation of satellite and aerial optical images alongside Digital Elevation Models (DEM), often followed by detailed field surveys to verify the findings.

This dataset provides a comprehensive basis for studying landslide occurrences, supporting the development and testing of deep learning models for remote sensing

---

[1]http://gpcv.whu.edu.cn/data/Bijie$_p$ages.html

landslide detection and segmentation.

## 3.2 Model and Architecture

### 3.2.1 CycleGAN Model

The CycleGAN model used in this study consists of two generators and two discriminators. The architecture for the generators and discriminators is detailed below.



Figure 3.1: Cycle GAN

**Generator Model**

The generator model employs a series of convolutional layers, instance normalization, ReLU activations, and residual blocks. The convolutional layers are used to extract features from the input images, and the residual blocks help in retaining important details while transforming the images from one domain to another. The model starts with a few convolutional layers to downsample the image, followed by multiple residual blocks to process the image in a deeper latent space, and finally, it uses transpose convolutional layers to upsample the image back to its original size. The detailed architecture is as follows:

------------------------------------------------------------------

| Layer (type) | Output Shape | Param # |
|---|---|---|
| Conv2d-1 | [-1, 64, 512, 512] | 1,792 |
| InstanceNorm2d-2 | [-1, 64, 512, 512] | 0 |
| ReLU-3 | [-1, 64, 512, 512] | 0 |
| Conv2d-4 | [-1, 128, 256, 256] | 73,856 |
| InstanceNorm2d-5 | [-1, 128, 256, 256] | 0 |
| ReLU-6 | [-1, 128, 256, 256] | 0 |
| Conv2d-7 | [-1, 256, 128, 128] | 295,168 |
| InstanceNorm2d-8 | [-1, 256, 128, 128] | 0 |
| ReLU-9 | [-1, 256, 128, 128] | 0 |
| Conv2d-10 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-11 | [-1, 256, 128, 128] | 0 |
| ReLU-12 | [-1, 256, 128, 128] | 0 |
| Conv2d-13 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-14 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-15 | [-1, 256, 128, 128] | 0 |
| Conv2d-16 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-17 | [-1, 256, 128, 128] | 0 |
| ReLU-18 | [-1, 256, 128, 128] | 0 |
| Conv2d-19 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-20 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-21 | [-1, 256, 128, 128] | 0 |
| Conv2d-22 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-23 | [-1, 256, 128, 128] | 0 |
| ReLU-24 | [-1, 256, 128, 128] | 0 |
| Conv2d-25 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-26 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-27 | [-1, 256, 128, 128] | 0 |
| Conv2d-28 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-29 | [-1, 256, 128, 128] | 0 |
| ReLU-30 | [-1, 256, 128, 128] | 0 |

| | | |
|---|---|---|
| Conv2d-31 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-32 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-33 | [-1, 256, 128, 128] | 0 |
| Conv2d-34 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-35 | [-1, 256, 128, 128] | 0 |
| ReLU-36 | [-1, 256, 128, 128] | 0 |
| Conv2d-37 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-38 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-39 | [-1, 256, 128, 128] | 0 |
| Conv2d-40 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-41 | [-1, 256, 128, 128] | 0 |
| ReLU-42 | [-1, 256, 128, 128] | 0 |
| Conv2d-43 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-44 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-45 | [-1, 256, 128, 128] | 0 |
| Conv2d-46 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-47 | [-1, 256, 128, 128] | 0 |
| ReLU-48 | [-1, 256, 128, 128] | 0 |
| Conv2d-49 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-50 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-51 | [-1, 256, 128, 128] | 0 |
| Conv2d-52 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-53 | [-1, 256, 128, 128] | 0 |
| ReLU-54 | [-1, 256, 128, 128] | 0 |
| Conv2d-55 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-56 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-57 | [-1, 256, 128, 128] | 0 |
| Conv2d-58 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-59 | [-1, 256, 128, 128] | 0 |
| ReLU-60 | [-1, 256, 128, 128] | 0 |
| Conv2d-61 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-62 | [-1, 256, 128, 128] | 0 |

| Layer (type) | Output Shape | Param # |
| --- | --- | --- |
| ResidualBlock-63 | [-1, 256, 128, 128] | 0 |
| Conv2d-64 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-65 | [-1, 256, 128, 128] | 0 |
| ReLU-66 | [-1, 256, 128, 128] | 0 |
| Conv2d-67 | [-1, 256, 128, 128] | 590,080 |
| InstanceNorm2d-68 | [-1, 256, 128, 128] | 0 |
| ResidualBlock-69 | [-1, 256, 128, 128] | 0 |
| ConvTranspose2d-70 | [-1, 128, 256, 256] | 295,040 |
| InstanceNorm2d-71 | [-1, 128, 256, 256] | 0 |
| ReLU-72 | [-1, 128, 256, 256] | 0 |
| ConvTranspose2d-73 | [-1, 64, 512, 512] | 73,792 |
| InstanceNorm2d-74 | [-1, 64, 512, 512] | 0 |
| ReLU-75 | [-1, 64, 512, 512] | 0 |
| Conv2d-76 | [-1, 3, 512, 512] | 1,731 |
| Tanh-77 | [-1, 3, 512, 512] | 0 |

==================================================================

Total params: 12,542,979

Trainable params: 12,542,979

Non-trainable params: 0

------------------------------------------------------------------

**Discriminator Model**

The discriminator model is designed to differentiate between real and generated images. It uses convolutional layers followed by instance normalization and leaky ReLU activations. The layers progressively reduce the spatial dimensions of the input image, ultimately leading to a single output that indicates whether the image is real or fake. This model helps in training the generator by providing feedback on how realistic the generated images are. The detailed architecture is as follows:

------------------------------------------------------------------

| Layer (type) | Output Shape | Param # |
| --- | --- | --- |
| Conv2d-1 | [-1, 64, 256, 256] | 3,136 |

```
        LeakyReLU-2           [-1, 64, 256, 256]              0
          Conv2d-3           [-1, 128, 128, 128]        131,200
   InstanceNorm2d-4          [-1, 128, 128, 128]              0
        LeakyReLU-5          [-1, 128, 128, 128]              0
          Conv2d-6            [-1, 256, 64, 64]         524,544
   InstanceNorm2d-7           [-1, 256, 64, 64]              0
        LeakyReLU-8           [-1, 256, 64, 64]              0
          Conv2d-9             [-1, 1, 63, 63]           4,097
================================================================
Total params: 662,977
Trainable params: 662,977
Non-trainable params: 0
----------------------------------------------------------------
```

### 3.2.2  VGG16 with Score-CAM

A deep convolutional neural network that was first created for picture categorization is the VGG16 model. Convolutional layers, ReLU activations, max-pooling layers, and completely connected layers are among its 16 layers. The fully connected layers carry out the final classification, while the convolutional layers are in charge of feature extraction. The detailed architecture is as follows:
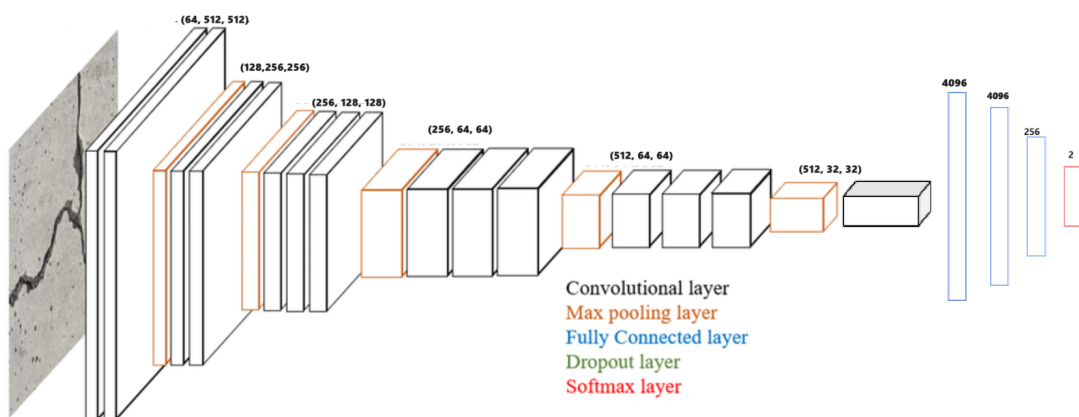


Figure 3.2: Classification Architecture

```
================================================================
Total params: 135,309,890
Trainable params: 1,049,346
Non-trainable params: 134,260,544
----------------------------------------------------------------
```

## 3.3  Training

**Data Preparation**

The first step in training the Cycle-GAN model was to prepare the dataset. The images from the Bijie landslide dataset were resized to dimensions of $3 \times 512 \times 512$ pixels. This resizing ensured uniformity in input size for the models, facilitating efficient processing during training.

**Batch Size and Optimization**

A batch size of 6 was chosen for training the Cycle-GAN model. This batch size strikes a balance between utilizing GPU memory efficiently and maintaining stability during training. With a smaller batch size, the model can update its weights more frequently, potentially aiding convergence.

The Adam optimizer was employed for training both the generator and discriminator networks in the Cycle-GAN model. The optimizer was configured with a learning rate of 0.0002 and betas set to $(0.5, 0.999)$. This choice of optimizer and learning rate was based on empirical observations and prior experimentation, aiming to achieve stable and efficient training.

**Loss Functions**

Several loss functions were used to train the Cycle-GAN model effectively. These included:

- **Adversarial Loss (GAN Loss):** Measured the ability of the generator to produce realistic images by fooling the discriminator.

- **Cycle Consistency Loss:** Ensured that the translated images from one domain to another and back remained close to the original input images.

- **Identity Loss:** Maintained the identity mapping between the input and output images, encouraging the generator to preserve essential features.

Combining the adversarial loss, cycle consistency loss, and identity loss—each weighted by its corresponding hyperparameters—was how the overall generator loss was calculated. This comprehensive loss formulation guided the training process towards generating high-quality translations between image domains.

**Training Process**

The training process involved iterating over the dataset for multiple epochs, with each epoch comprising forward and backward passes through the network. The models were trained on an RX400 GPU, with an average epoch duration of 764 seconds.

During each epoch, the generator and discriminator networks were updated iteratively. The generator attempted to minimize the total generator loss, while the discriminators aimed to correctly classify real and generated images.

The training progress was monitored closely, and model checkpoints were saved periodically, typically every 50 epochs. This checkpointing strategy allowed for the recovery of trained models in case of interruptions and facilitated model evaluation at different stages of training.

After 550 epochs of training, the model's performance was evaluated, with the best-performing model observed at epoch 450. This model was selected based on its ability to produce high-quality image translations and maintain consistency across different domains.

**Fine-Tuning of VGG16**

In addition to training the Cycle-GAN model, the VGG16 network was fine-tuned for image classification tasks. The images were resized to $3 \times 512 \times 512$ pixels, and a batch size of 32 was utilized for training on a Colab T4 GPU.

The fine-tuning process involved optimizing the VGG16 model parameters using the stochastic gradient descent (SGD) optimizer with a learning rate of 0.001 and

momentum of 0.9. A cross-entropy loss function was employed to compute the classification loss during training.

The model was trained for 100 epochs, with checkpoints saved every 10 epochs. The best-performing model was identified at epoch 20 based on its classification accuracy and overall performance on validation data.

## 3.4 Combining Generated Masks

**Cycle-GAN Generated Mask**

The Cycle-GAN model was trained to generate non-landslide images from landslide images. To obtain a mask highlighting the landslide regions, we adopted a simple differencing approach.

Given a landslide image $L$ and its corresponding non-landslide image $N$ generated by the Cycle-GAN, we computed the absolute difference $D = |N - L|$. This difference image $D$ contains intensity variations that correspond to the regions where landslides are present.

To convert the difference image into a binary mask, we first converted it to grayscale using the formula $0.299 \cdot R + 0.587 \cdot G + 0.114 \cdot B$, where $R$, $G$, and $B$ represent the red, green, and blue channels, respectively. This grayscale conversion ensured that the intensity of each pixel represented the magnitude of change between the landslide and non-landslide images.

Next, we normalized the pixel values in the grayscale image to the range [0, 1] to facilitate thresholding. We empirically determined a threshold value of 0.08 through experimentation, above which pixel intensities were considered indicative of landslide regions.

By applying this threshold, we obtained a binary mask representing the predicted landslide areas generated by the Cycle-GAN model.

**Score-CAM Generated Mask**

For the Score-CAM model applied to the VGG16 network, the goal was to generate a heatmap highlighting the regions of the image that contributed most to the landslide prediction.

Score-CAM operates by computing the importance scores for each neuron in the last convolutional layer of the network. These scores are then weighted by the corresponding activation maps to produce the final heatmap.

We experimented with various threshold values ranging from 0.1 to 0.6 to determine the optimal threshold for generating the Score-CAM mask. Through a grid search, we observed that a threshold of 0.4 yielded the highest precision in identifying landslide regions while minimizing false positives.

This thresholding process resulted in a binary heatmap where pixels above the threshold were considered significant contributors to the landslide prediction, while those below were disregarded.

## Combining Masks

Once the individual masks from the Cycle-GAN and Score-CAM models were obtained, we combined them to produce the final landslide mask.

We applied a logical AND operation (Intersection of 2 Masks) between the two binary masks generated by the Cycle-GAN and Score-CAM models. This operation resulted in a merged mask where any pixel identified as a landslide region by either model contributed to the final prediction.

The combined mask effectively leveraged the strengths of both models, utilizing the Cycle-GAN's ability to generate realistic images and the Score-CAM's capability to identify relevant image features. By aggregating the predictions from both models, we aimed to enhance the accuracy and robustness of the final landslide detection.

## Threshold Selection

I initially experimented with Niblack and Sauvola Thresholding methods, but they produced unsatisfactory results. Instead, Otsu Thresholding yielded superior performance.

In addition to combining the masks, we utilized Otsu Thresholding to optimize the threshold values for generating the masks. Otsu's method is an automatic threshold selection technique that determines the optimal threshold by maximizing the variance between foreground and background pixels, thereby minimizing the

within-class variance.

Applying Otsu Thresholding allowed us to achieve an optimal balance between sensitivity and specificity in identifying landslide regions. This method ensured that the threshold values were selected based on the image data distribution, leading to more accurate segmentation results.

We determined the final threshold values using Otsu's method on the validation dataset, comparing the combined mask's predictions against ground truth annotations. Through this automatic threshold optimization process, we aimed to produce a reliable and accurate landslide detection system capable of identifying hazardous areas with high precision.

By leveraging Otsu Thresholding, we ensured that the combined masks accurately captured true positives while minimizing false positives and false negatives. This methodical approach to threshold optimization contributed to the robustness and reliability of our landslide detection system.

# Chapter 4

# Results

## 4.1 Introduction

In this chapter, we present the evaluation of the VGG16 classification model along with the masks generated by Score-CAM, CycleGAN, and their combined version against the ground truth. The performance of these models and methods is assessed using four key metrics: Precision, Recall, and Mean Intersection over Union (mIoU). These metrics provide a comprehensive understanding of how well the models perform in the context of landslide remote sensing image semantic segmentation.

## 4.2 Methodology

### 4.2.1 Model Evaluation Method

In the domain of landslide remote sensing image semantic segmentation tasks, the evaluation of test results is commonly conducted using metrics such as precision, recall, and mIoU. These metrics are derived from the fundamental calculations involving True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN). Below, we outline the definitions and calculations of these metrics:

**Precision**

Precision measures the proportion of positive samples correctly predicted by the model among all samples predicted as positive. It reflects the accuracy of the model in identifying landslide areas from the images. The precision is calculated using the following formula:

$$\text{Precision} = \frac{TP}{TP + FP} \tag{4.1}$$

**Recall**

Recall indicates the proportion of actual positive samples that were correctly predicted by the model. It represents the model's ability to detect all the landslide areas present in the images. The recall is calculated as follows:

$$\text{Recall} = \frac{TP}{TP + FN} \tag{4.2}$$

**Mean Intersection over Union (mIoU)**

The mIoU metric is used to comprehensively evaluate the performance of the segmentation model. The Intersection over Union (IoU) is calculated by dividing the intersection of the predicted segmentation and the ground truth by their union. The mIoU is the average of the IoU values for each category. The formula for IoU is:

$$\text{mIoU} = \frac{A \cap B}{A \cup B} = \frac{TP}{TP + FP + FN} \tag{4.3}$$

## 4.2.2 Evaluation Metrics Calculation

To thoroughly evaluate the VGG16 classification model and the masks generated by Score-CAM, CycleGAN, and their combined version, we calculated the above metrics for each method. The ground truth masks were used as a benchmark to compare the predicted masks. The following steps outline the process:

1. **True Positives (TP):** Count the number of pixels correctly identified as part of the landslide. 2. **False Positives (FP):** Count the number of pixels incorrectly identified as part of the landslide. 3. **True Negatives (TN):** Count

the number of pixels correctly identified as not part of the landslide. 4. **False Negatives (FN):** Count the number of pixels incorrectly identified as not part of the landslide.

Using these counts, the precision, recall and mIoU were computed to assess the performance of each method.
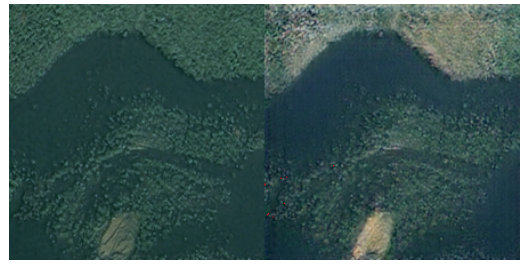
### 4.2.3    10 Fold Cross Validation with Modified Dataset

Due to the relatively small size of our dataset, we employed 10 Fold Cross-Validation to ensure robust and reliable evaluation of our models. Our original dataset comprised 770 landslide images and 2003 non-landslide images. To enhance the dataset and improve the training process, we performed the following modifications:

- **Data Augmentation:** Using a cyclic GAN, we generated an additional 400 landslide images. During this process, 10 images were manually identified as unsuitable and subsequently deleted. This brought the total number of landslide images to 1160. Some generated images with original images are shown below.

- **Data Cleaning:** For the non-landslide images, we manually reviewed the dataset to remove images containing buildings and highways, which were deemed irrelevant for the task. After this cleaning process, we were left with 1558 non-landslide images.
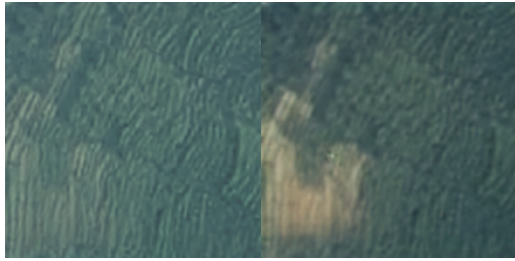
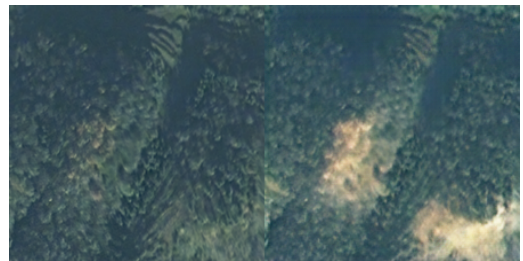Figure 4.1: A set of 10 images showing the transformation of original non-landslide images into landslide images using Cyclic-GAN.
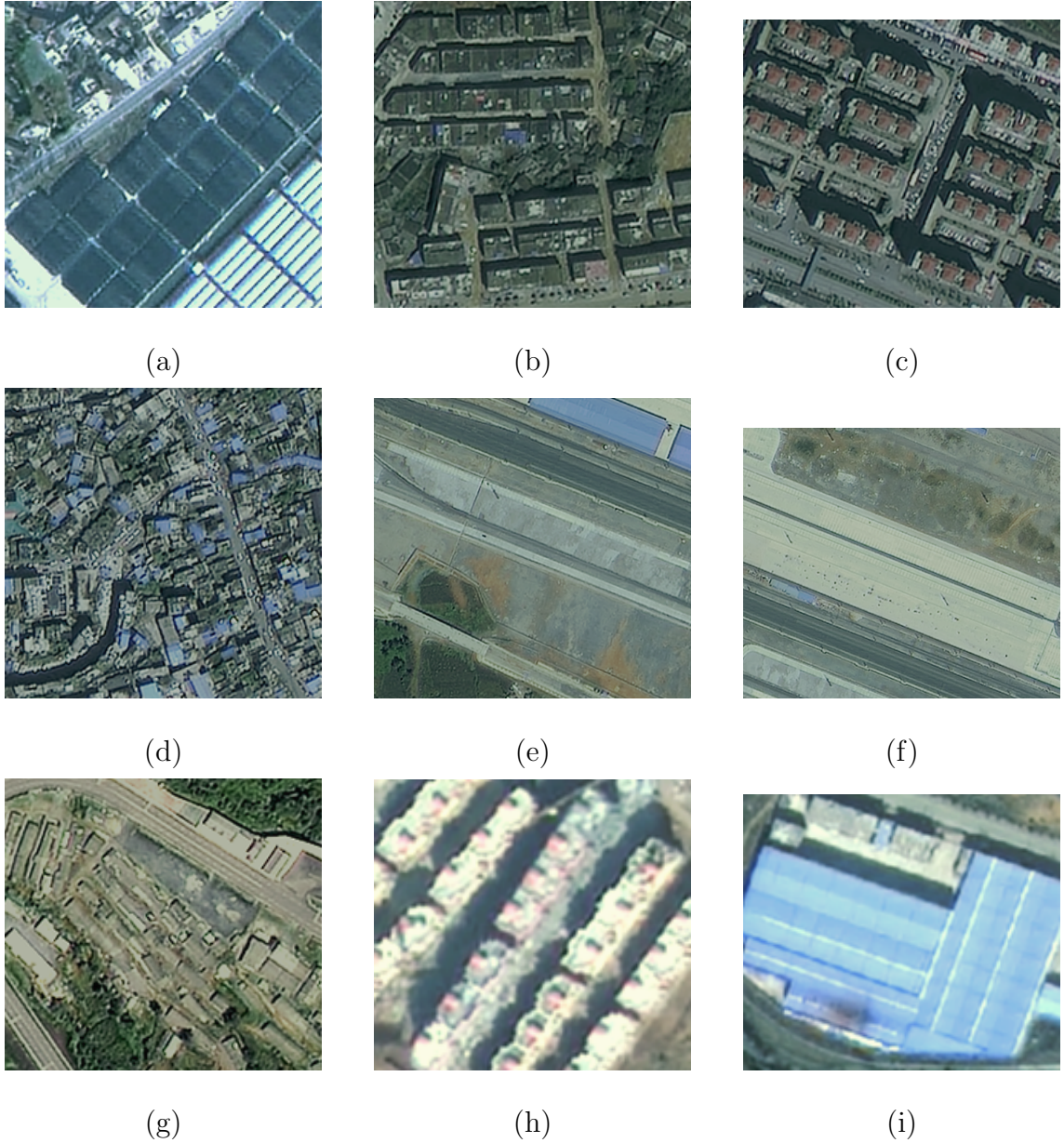
(a) (b) (c)







(d) (e) (f)







(g) (h) (i)

Figure 4.2: A set of 9 images depicting the deleted samples from the non-landslide class, which were removed due to the presence of highways and buildings.

With these modifications, our final dataset comprised 1160 landslide images and 1558 non-landslide images. This dataset was then divided into 10 equal parts to perform 10 Fold Cross-Validation. The procedure was as follows:

1. **Dataset Partitioning:** The entire dataset was randomly divided into 10 equal subsets. Each subset contained both landslide and non-landslide images in proportionate amounts.

2. **Model Training:** For each fold, the GAN and VGG16 models were trained

on 9 of the 10 subsets. This training process aimed to learn and optimize the models using the majority of the dataset while reserving one subset for validation.

3. **Mask Generation and Evaluation:** After training, the models were used to generate masks on the remaining subset. These generated masks were then compared to the ground truth masks to evaluate the model's performance on unseen data.

This cross-validation approach ensured that each image in the dataset was used for both training and validation, providing a comprehensive assessment of the model's performance across different data splits. By averaging the results across all 10 folds, we obtained a more reliable estimate of the models' generalization capabilities and robustness.

The use of 10 Fold Cross-Validation with our modified dataset allowed us to make the most of the available data and achieve a thorough evaluation of our proposed methods. This approach mitigates the risk of overfitting and ensures that the reported performance metrics are representative of the model's true capability in practical applications.

Table 4.1: Test results comparing fully supervised and weakly supervised learning methods

| Method | Precision | Recall | mIOU |
|---|---|---|---|
| **Weakly Supervised Learning** | | | |
| Grad-CAM | 0.692 | 0.593 | 0.159 |
| CycleGAN | 0.845 | 0.404 | 0.184 |
| Grad-CAM + CycleGAN | 0.924 | 0.383 | 0.237 |
| **Score-CAM + CycleGAN + Manual Threshold** | **0.918105** | **0.375309** | **0.228908** |
| **Score-CAM + CycleGAN + Otsu Threshold** | **0.917362** | **0.382849** | **0.227165** |
| **Supervised Learning** | | | |
| U-Net | 0.955 | 0.555 | 0.408 |

a. Original Image    b. Original Mask    c. GAN Mask    d. CAM Mask    e. Combined Mask

Figure 4.3



a. Original Image    b. Original Mask    c. GAN Mask    d. CAM Mask    e. Combined Mask

Figure 4.4



a. Original Image    b. Original Mask    c. GAN Mask    d. CAM Mask    e. Combined Mask

Figure 4.5



a. Original Image    b. Original Mask    c. GAN Mask    d. CAM Mask    e. Combined Mask

Figure 4.6

a. Original Image　　b. Original Mask　　c. GAN Mask　　d. CAM Mask　　e. Combined Mask

Figure 4.7



a. Original Image　　b. Original Mask　　c. GAN Mask　　d. CAM Mask　　e. Combined Mask

Figure 4.8



a. Original Image　　b. Original Mask　　c. GAN Mask　　d. CAM Mask　　e. Combined Mask

Figure 4.9
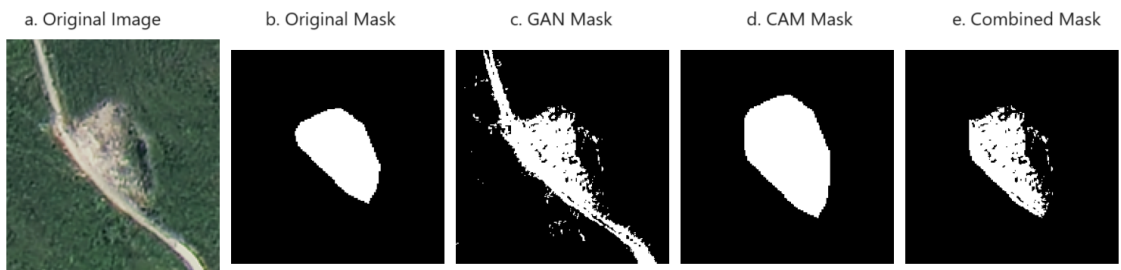


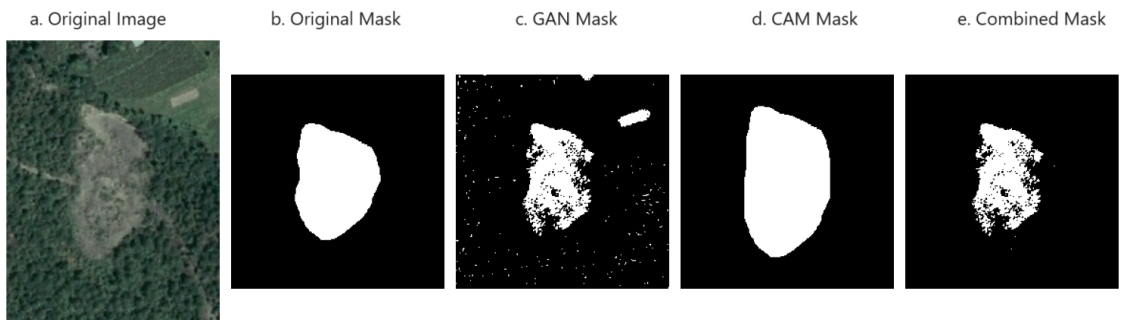a. Original Image　　b. Original Mask　　c. GAN Mask　　d. CAM Mask　　e. Combined Mask
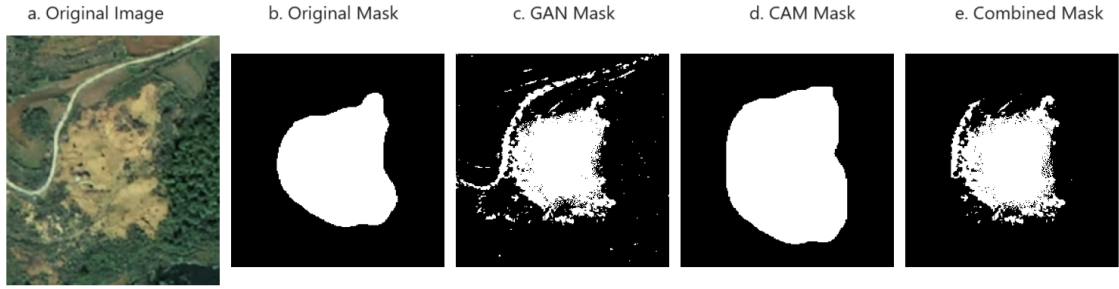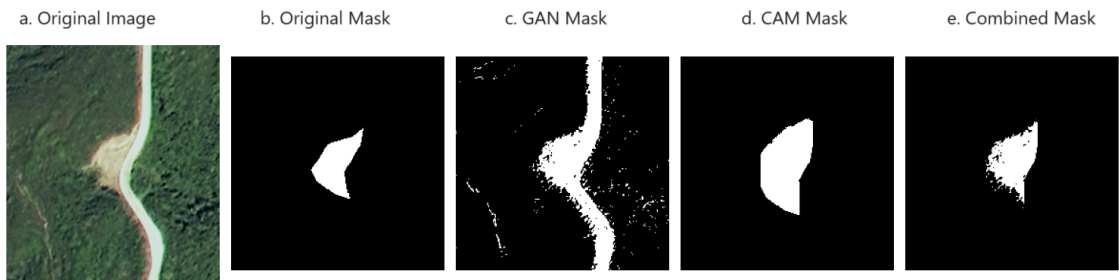
Figure 4.10

Figure 4.11



Figure 4.12

## 4.3   Some Important Observations

Based on our experimental assessments, we found that the GAN generated more precisely located masks for areas affected by landslides. The masks generated by the GAN were effective in capturing the broader areas of landslide occurrences. However, one significant drawback was that these masks often contained a considerable amount of noise, which could lead to false positive detections in certain scenarios.

On the other hand, the Score-CAM applied on the VGG16 model demonstrated a remarkable ability to capture the exact location of landslide regions with higher precision. The Score-CAM-generated masks were less noisy and more focused on the actual landslide areas, which significantly improved the precision of the detection.

Given these complementary strengths, we combined the masks generated by both the GAN and the Score-CAM. This combined approach leveraged the GAN's capability for broader localization and the Score-CAM's precision in identifying specific regions. As a result, the fusion of these two methods produced superior detec-

tion outcomes, balancing localization and precision effectively. This combination provided a more accurate and reliable mask for landslide detection, significantly improving the overall performance of our system.

# Chapter 5

# Conclusion

In this research, we explored and evaluated the efficacy of combining Score-CAM and Cycle-GAN for the task of landslide detection and segmentation from remote sensing images. The integration of these two methodologies aimed to leverage the strengths of both approaches, ultimately enhancing the performance metrics of our landslide detection system.

## 5.1 Summary of Findings

The primary objective of our study was to improve the precision, recall, and mean Intersection over Union (mIoU) scores of landslide segmentation models. Through extensive experimentation, we found that combining Score-CAM and Cycle-GAN significantly improved the recall and mIoU scores compared to using each method individually.

- **Score-CAM**: Score-CAM was effective in pinpointing the exact location of landslides within the images. This method capitalized on the interpretability of convolutional neural networks (CNNs), specifically focusing on the activation maps to generate precise heatmaps for landslide regions.

- **Cycle-GAN**: Cycle-GAN excelled in generating localized masks for landslide areas. However, the generated masks often contained noise, which could potentially reduce the overall accuracy and increase the false positive rate.

- **Combined Approach**: By integrating the detailed localization capabilities

of Score-CAM with the robust mask generation of Cycle-GAN, we achieved a more reliable segmentation model. This combined approach not only improved the recall but also enhanced the mIoU score, demonstrating a better balance between sensitivity and specificity in identifying landslide regions.

## 5.2 Performance Evaluation

Our evaluation metrics, which included precision, recall, mIoU, and false positive rate (FPR), provided a comprehensive assessment of the model's performance. The combined Score-CAM and Cycle-GAN approach yielded the following improvements:

- **Recall**: The recall score increased, indicating a higher percentage of actual landslide areas correctly identified by the model. This is crucial for applications where missing a landslide area could have severe consequences.

- **mIoU**: The mean Intersection over Union score saw a notable improvement, reflecting a better overlap between the predicted landslide areas and the ground truth annotations.

## 5.3 Key Observations

- **Localization vs. Noise**: While Cycle-GAN produced well-localized masks, the presence of noise was a significant drawback. In contrast, Score-CAM provided precise localization but lacked the robustness of Cycle-GAN in some instances. The combination of both methods successfully mitigated these individual shortcomings.

- **Threshold Optimization**: The use of Otsu Thresholding played a critical role in fine-tuning the segmentation masks. By optimizing the threshold values, we were able to achieve a better balance between detecting true positives and minimizing false positives and negatives.

- **Dataset Augmentation and Modification**: The augmentation of the dataset with additional landslide images generated by Cycle-GAN, and the careful

curation of non-landslide images, further contributed to the improved performance of the combined model.

## 5.4 Future Work

The promising results of this study pave the way for several avenues of future research:

- **Real-World Applications**: Further testing and validation on diverse real-world datasets would help in assessing the generalizability and robustness of the combined model in different geographic and climatic conditions.

- **Algorithm Enhancements**: Exploring advanced GAN architectures and more sophisticated activation mapping techniques could potentially yield even better results. Additionally, integrating other weakly supervised and unsupervised learning methods could further enhance the model's performance.

- **Automated Thresholding**: Developing automated methods for threshold optimization, possibly through machine learning techniques, could streamline the process and adapt to varying data conditions dynamically.

## 5.5 Conclusion

The integration of Score-CAM and Cycle-GAN has proven to be an effective strategy for improving the performance of landslide detection and segmentation models. By leveraging the strengths of both methodologies, we achieved higher recall and mIoU scores, making our approach a valuable tool for remote sensing image analysis. This research underscores the importance of combining complementary techniques to tackle complex image segmentation tasks, ultimately contributing to more accurate and reliable landslide detection systems.

# Bibliography

[1] Zhou, Y.; Wang, H.; Yang,R.; Yao, G.; Xu, Q.; Zhang, X. A Novel Weakly Supervised Remote Sensing Landslide Semantic Segmentation Method: Combining CAM and cycleGAN Algorithms. *Remote Sensing* 2022, 14, 3650. https://doi.org/10.3390/rs14153650

[2] Johnson, J., Alahi, A., Fei-Fei, L. (2016). Perceptual losses for real-time style transfer and super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 694-711).

[3] He, K., Zhang, X., Ren, S., Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 770-778).

[4] Ronneberger, O., Fischer, P., Brox, T. (2015). U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention* (pp. 234-241). Springer, Cham.

[5] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision (ICCV)* (pp. 618-626).

[6] Zhang, X., Wang, Y., Su, H. (2020). Score-CAM: Score-weighted visual explanations for convolutional neural networks. *CVPR'20 Workshop*.

[7] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Bengio, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems* (pp. 2672-2680).

[8] Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K. Q. (2018). Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 172-189).

[9] Kingma, D. P., Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980.*

[10] Zhu, J.-Y., Park, T., Isola, P., Efros, A. A. (2017). Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. *Proceedings of the IEEE International Conference on Computer Vision (ICCV).* https://arxiv.org/abs/1703.10593.