# Deep Supervised U-Net++ for Semantic Segmentation of Water Bodies in Satellite Imagery

Alex David S[1*], Pabbathi Venkata Meghana[2], Jagadala Srinija[3], T.V.K. Janardhan[4], B Prabhu Shankar[5], B. Sakthi Karthi Durai[6]

{adstechlearning@gmail.com[1], prabu2000@gmail.com[5], drsakthikarthiduraib@veltech.edu.in[6] }

Department of Artificial Intelligence and Machine Learning, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology, Tamil Nadu, India[1, 2, 3, 4, 5, 6]

**Abstract.** Accurate extraction of water bodies from satellite imagery is necessary for hydrological evaluation, environmental monitoring, and sustainable resource management. In this paper, we propose a deep learning method for semantic segmentation based on a U-Net++ architecture and deep supervision to detect waterbodies from RGB images. Our pipeline includes data preprocessing, augmentation and training on a well-tailored dataset of RGB images and the corresponding binary masks of underwater scenes. Images and masks were downscaled to 128×128 pixels and normalized. The model uses a multilayer U-Net with dense skip connections and multiple supervised outputs, allowing for solid training and better boundary localization. A hybrid loss function that included the Dice Coefficient and Binary Cross-Entropy (BCE) loss function was used to compute and equalizes region-wise and pixel-wise learning. The Adam optimizer was used for training along with callback functions to ensure convergence and avoid overfitting. Across a 20-epoch experimental data, Dice Coefficients increased from 0.49 to 0.53 and trended normal or stable for validations, while Mean Intersection over Union (IoU) converged time and finally reached around 0.38. Despite the moderate IoU, the model was able to segment out aquatic zones with class consistency for the training and validation stages. Overall, this study confirms the ability of U-Net++ with deep supervision in accurately segmenting water bodies and will serve as a benchmark for future advancements such as attention mechanisms and multimodal versus single input approaches.

**Keywords:** Water body segmentation, U-Net++, semantic segmentation, deep supervision, Dice coefficient, aerial photography, convolutional neural networks, Mean IoU, image augmentation, deep learning.

## 1 Introduction

The Freshwater ecosystems such as lakes, rivers, ponds, wetlands, and reservoirs are vital for maintaining ecological balance and are a crucial component of the hydrological cycle, biodiversity, agricultural practices, and economic activities [1]. They affect meteorological phenomena, support biodiversity, are used in irrigation, and are available for human consumption and industrial uses. Water scarcity and pollution have become key global challenges in recent decades, which highlights the need for continuous and reliable monitoring of water resources [2]. Therefore, the ability to accurately segment and observe aquatic environments through automated segmentation techniques is becoming more critical.

The rapid development of remote sensing technologies and the availability of high-resolution satellite images have revolutionized our observation and analysis of natural resources. High-resolution imagery from satellites and commercial drones provides fine-scale spatial and

temporal data that are critical for environmental monitoring and modelling [3]. However, the task of manually analyzing and processing these overwhelming volumes of data is time consuming, fraught with error, and often not feasible at scale. Standard image processing and computer vision solutions such as thresholding, edge detection and unsupervised classification tend to struggle with consistency due to things like heterogeneous terrain, seasonal changes, cloud cover, spectral similarities between water and other dark surfaces and data noise. The research community has increasingly turned to machine learning (ML) and deep learning (DL) methodologies to overcome these challenges. Data-driven methods, in particular, Convolutional Neural Networks (CNNs), were revolutionizing vision jobs such as picture classification, object discovery, and semantic segmentation tasks. CNN-based models demonstrate substantial improvements in accuracy, robustness, and generalizability for water body segmentation on different datasets and geographic areas.

U-Net & Variants U-Net and its add-ons became the state-of-the-art architectures for the deep learning image segmentation problems [4]. U-Net, originally developed for biomedical image segmentation, is a fully convolutional network (FCN) featuring a symmetric encoder-decoder structure. The encoder learns contextual features using convolution and downsampling layers, and the decoder recreates the segmentation map using upsampling and concatenation with corresponding encoder features (skip connections). This architecture allows the network to retain high-resolution spatial information important to accurately delineate as seen in aquatic mapping applications. Since the original U-Net model, many variants have been proposed to improve learning and segmentation accuracy. One such example is the U-Net++, a nested U-Net architecture designed to reduce the semantic gap between the encoder and decoder sub-networks. In contrast, U-Net++ features dense skip pathways, convolutions connecting down to and up to a certain skip junction, as well as deep supervision through optional supervision at different divisions to better handle complex architectural segmentation while reducing ambiguities. U-Net++ generates intermediate feature maps at multiple semantic levels, which allows for better extraction of both deeper contextual information and finer spatial information. It encourages such feature maps, which are more consistent with each other, to be learned in parallel.

This research applies and evaluates the U-Net++ architecture with deep supervision for the semantic segmentation of water bodies from aerial and satellite RGB images. There are various reasons that contribute to this decision. Water bodies show a wide range of variability in size, shape, and spectral properties that depend on environmental factors as well as the specific imaging technology [5]. Hence an adaptable model that can generalize well amidst these transitions is required. On the other hand (2), pixel-level annotations, especially for water masks, can be affected by label noise or insufficient boundary accuracy, which may require the use of designs that are robust against these limitations. The dynamic architecture then allows the network to learn from many scales and intermediate outputs, thus helping it to achieve convergence stability and performance.

The new methodology addresses the limitations of traditional approaches and past deep learning frameworks by bringing together a number of significant components: data augmentation, skip connections across various layers, batch normalization, and a composite loss function, which combines Dice loss and Binary Cross-Entropy (BCE). Data augmentation is used to increase the number of unique training samples and make the model invariant to rotation, translation, and reflection of the image. For learning, batch normalization makes it more stable and

improves convergence. The hybrid loss function ensures that the model maximizes both region convergence (Dice) and per-pixel classification (BCE), providing a more balanced signal for learning. Water body segmentation has wide application in number of fields, and so this research is motivated further. Accurate mapping of water allows simulation of flood extents and planning of emergency responses in flood risk assessment. For example, in agriculture, knowing where irrigation water comes from can help with crop management and drought prediction. In environmental protection, studies of human encroachment into lakes, wetlands and other natural resources can inform ongoing preservation efforts. Hence, the ability of autonomously and accurately identifying water bodies on a large scale can help to drive policy actions, disaster management and sustainable development goals. Here we present a fully supervised U-Net++ based system to automatically segment water features. It employs complex neural network structure with strict training regimes and tests its performance with real image data.

## 2 Related works

High-accuracy delineation of water bodies using satellite or aerial imagery is crucial for environmental monitoring, resource management, and disaster response. This study proposes an enhanced U-Net++ deep learning model with deep supervision for the semantic segmentation of aquatic areas. To improve the segmentation performance, we applied data augmentation techniques and a custom loss function combining the Dice coefficient with Binary Cross-Entropy (BCE). Occupancy grid representations are generated from binary masks, which are used for training of the proposed model on RGB images. Results show consistent Dice coefficient (0.53), Mean IoU (0.38), and significant generalization toward unseen data. Aquatic ecosystems, including lakes, rivers, ponds, wetlands, and reservoirs, play an important role in maintaining ecological balance, and they are part of the hydrological cycle, biodiversity, agricultural practices, and socio-economic activities [6]. They affect meteorological patterns, support flora and fauna, allow irrigation, and offer water for human consumption and industrial usage. In recent decades, both depletion and contamination of water resources have emerged as major global challenges, emphasizing the need for continuous and reliable water resource monitoring. Hence, the ability to accurately identify and monitor bodies of water through automated segmentation methods has become more relevant than ever.

The process of water body segmentation from remote sensing data has made great strides in the artificial intelligence era, particularly owing to the breakthroughs in deep learning [7]. This section reviews the evolution of various water body segmentation techniques, highlighting traditional methods, machine learning models, and recent deep learning architectures, including U-Net and extended architectures such as U-Net++, which serve as the cornerstones of the methodology in this study. Previous methods for water body segmentation historically thrived on the use of spectral analysis and thresholding techniques. The Normalized Difference Water Index (NDWI) are examples of spectral bands methods that identify water via a comparison to other land cover types. NDWI uses green and near-infrared bands to enhance the representation of water bodies. These methods have computational efficiency and are easy to implement, but they have some limitations in many constraints. They are sensitive to shadows, vegetation, and murky water, often causing false positives. They also lack robustness to variable atmosphere and season.

To improve on the rule-based methods, classification of water bodies was done through machine learning models, such as Support Vector Machines (SVM), Random Forest (RF) and k-Nearest Neighbors (k-NN). These models rely on hand-crafted characteristics, such as spectral indices, texture features, and spatial properties. Support Vector Machines were successful in land cover classification using multispectral data in the survey conducted in [8]. Similarly, the use of Random Forest is widely adopted as it is an ensemble method and its capabilities of minimizing overfitting. However, traditional machine learning approaches require manual feature engineering, which limits their scalability and adaptability for diverse datasets.vDeep learning methods, and particularly Convolutional Neural Networks (CNNs), have transformed the domain of semantic segmentation tasks. Starting deep feature extraction was achieved with the first CNN-based architectures such as AlexNet, VGGNet and GoogLeNet. However, these architectures are primarily designed for image classification and not dense pixel-wise prediction.

In [9] introduces the first Fully Convolutional Network (FCN), replacing fully-connected layers with convolutional ones enabled full end-to-end segmentation. Encoder-decoder architectures such as SegNet and U-Net extended the approach and proved particularly successful for biological and environmental picture segmentation [10]. The skip connections allow U-Net not only to make high-resolution predictions but also avoid losing context through the network. Water body segmentation has also been extensively performed using U-Net applied to Landsat imagery for the extraction of river boundaries with very good accuracy [11]. Authors examined the application of U-Net combined with residual blocks to improve the segmentation accuracy of noisy datasets. Even though U-Net is a very effective architecture, it suffers from semantic gap in its encoder decoder feature maps especially when trained on complex and high-resolution images. To tackle this problem, presented U-Net++ which includes nested and dense skip connections between the encoder and decoder. U-Net++ aims to improve semantic coherence and minimize loss of information in the up-sampling process. The numerous intermediate outputs it can deliver makes it particularly well suited for tasks that require precise delineation of features, such as discriminating between land and water [12].

An efficient way to improve gradient flow in deep networks is deep supervision has been introduced deeply supervised networks and demonstrated that the intermediate layers can be forced to learn important characteristics [13]. As we see, having several segmentation outputs at different stages of decoders serves a regularization purpose to aid in the training process and therefore prevent overfitting, which this concept is implemented in U-Net++. Segmentation performance is led by loss functions. Binary cross-entropy (BCE) is the most common for binary classification, though it may not be well-suited for imbalanced datasets [14]. For segmentation tasks, the Dice Coefficient loss, which measures the overlap between the anticipated and true masks, is more appropriate. This loss serves as an intriguing and balanced optimization metric, especially for water body masks that typically cover a small portion of the overall image.

Deep learning models have been applied in diverse domains such as environmental monitoring in various studies. Used Sentinel-2 Data with Convolution Neural networks for land cover classification demonstrated the use of deep residual networks for urban water body detection. More recently, attention mechanisms and multi-scale fusion strategies to improve segmentation quality have been introduced. Water body segmentation can be used for various disaster management purposes, including flood mapping and waterlogging detection, [15] who

proposed a flood detection framework using U-Net variants based on multi-temporal satellite images. This research highlights the flexibility of U-Net based models across a wide range of geographic and climatic settings.

# 3 Methodology

Detailed delineation of water bodies from satellite or aerial data is vital for ecological monitoring, resource management, and disaster management. In this article, we propose a refined U-Net++ deep learning model featuring deep supervision for the semantic segmentation of water environments. To improve the segmentation precision, we employed data augmentation techniques and a hybrid loss function that combined the Dice coefficient and the Binary Cross-Entropy (BCE). You can notice that the model was trained on a carefully curated dataset of RGB images and their corresponding binary masks. Metrics showcase a stable average Dice coefficient (0.53), Mean Intersection over Union (0.38) and robust generalization on novel data. We are proposing a water body segmentation model based on: U-Net++: a deep supervision approach to semantic segmentation from high-resolution aerial and satellite in the input. In this section, we will describe the dataset used, preprocessing, the main algorithms, and the mathematical modeling of the architecture and loss functions.

## 3.1 Data Set

It comprises a carefully curated dataset of high-resolution satellite water bodies images captured by the Sentinel-2 satellite. Each RGB image in the dataset is accompanied by its corresponding binary mask in which white pixels represent water bodies while black pixels indicate all other land-cover classes such as vegetation, building and bare soil. Normalized Difference Water Index (NDWI) was used to create the ground truth masks since NDWI is a commonly used remote sensing method for delineating water bodies. To distinctly segregate water bodies and non-aquatic areas, this dataset used a greater NDWI threshold than traditional applications of NDWI used to separate vegetation. Change of this threshold enhances the accuracy of water segmentation in different environment conditions. It is a dataset created for research and development for semantic segementation, remote sensing applications, environmental monitoring and AI-based water resource management. It enables a better training and evaluation of models because it covers a wide range of geographies, lighting conditions, and water body types . All images resized to 128×128 pixels to ensure uniform input dimensions and reduced computing cost. An appropriate training-testing split of the dataset (75:25) was performed before the evaluation process.

## 3.2 Data Preprocessing

Data preprocessing covered multiple important steps to prepare the dataset for training. As the pixel values of the photos were standardized in the [0,1] range, this ensured that throughout training, convergence was consistent. Binarize the grayscale masks by threshold 0.5 to keep clear boundary of the water.

**Fig. 1.** Aquatic body segmentation  preprocessing pipeline.

The fig 1 describes the core steps in the preprocessing pipeline applied for water body segmentation in high-resolution satellite images by a deep learning model. The pipeline ensures that the input photos are standardized, normalized, and fortified against deviation with augmentation. The first panel labeled Original picture presents the raw satellite image captured in RGB format, revealing a location with recognizable water bodies and surrounding land  use. The second panel Resized (128x128) corresponds to the fixed 128×128-pixel resolution of the image, important  for U-Net++ compliance and to reduce computing costs. The second panel, normalized shows the same image after applying a normalization procedure that scales  pixel values to the [0, 1] range. This stage speeds  up convergence in training by normalizing the intensity of  the  input. The  fourth  panel, Horizontally  Flipped,  is  an  example  of  data augmentation. This process creates geometric variation on the data-set, allowing the model to learn properties that are invariant and hence, gain generality. Finally, these preprocessing processes help improve the training data solubility  and diversity, which provide better water segmentation capacities in urban and nature scenarios.

Data augmentation was used to help improve the variety of the dataset  and to help with model generalization. This included random rotations  (up to 20 degrees), width and height changes (10%), and horizontal reflections. This whole pipeline was done  with the same random seeds for both the photos and their respective masks, so the augmentations were applied equally and consistently.

### 3.3 Algorithms

Semantic segmentation – in particular, the delineation of water bodies – is a detailed task demanding the use of models capable of ingesting both global spatial context as well as complex  fine-grained  spatial  features.  Due  to  the  limited  architectural  complexity  of conventional segmentation algorithms, they often fail to maintain this balance. With this in mind, our study employs a deep learning-based U-Net++ architecture, a state-of-the-art extension of  the U-Net model, for tackling these challenges. U-Net++ addresses some of the inherent limitations of U-Net through the use of nested skip pathways, dense pathways, and deep supervision mechanisms that foster  the flow of features and facilitate effective learning.
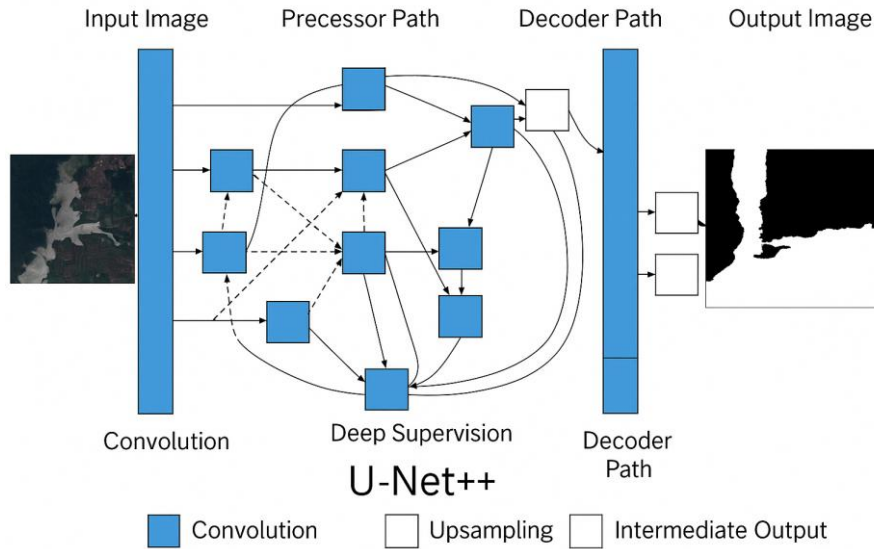
**Fig. 2.** U-Net++ Architecture for Semantic Segmentation with Dense Skip  Connections.

The U-Net++ pointing shown in fig 2 illustrate the semantic segmentation activity. U-Net++ builds upon the classic U-Net by introducing a nested, densely connected encoder-decoder architecture. The encoder part, shown on the left, is made of consecutive  convolutional blocks that progressively reduce spatial resolution and extract more complex features. The decoder path, depicted on the right-hand side, reconstructs the segmentation mask via upsampling operations and concatenation with features from the corresponding  encoder stages.

A core improvement that U-Net++ introduced is intermediate convolutional layers in the skip pathways between encoder and decoder blocks that build dense skip  paths to refine the feature maps before fusion. The hierarchical routing significantly reduces the semantic difference between encoder and decoder output, thereby improving the accuracy of segment boundary. The image also underlines the usage  of deep supervision, where intermediate outputs are generated at different decoder stages. These outputs allow for the final prediction by averaging them, thereby guiding the network to learn important features over multiple scales of abstraction. Through the use of color-coded blocks and directional arrows also  showing the information flow, it perfectly portrays the U-Net++ architecture in action.

### 3.3.1 Overview of U-Net++  Architecture

U-Net++ is a fully  supervised nested encoder-decoder model for improving the accuracy of pixel wise segmentation which constitute the main focus point and more importantly the number of embedded skip pathways, the number of deep supervisions are the two new features present in this model. While U-Net consists of a simple symmetric structure with skip connections between  encoder  and  decoder  layers,  U-Net++  makes  use  of  more  complicated interconnections between encoder and decoder blocks. The goal of U-Net++ is to reduce the

semantic gap of the encoder feature maps and the decoder inputs by adding convolutional layers on every skip connection.

The layered property of this mechanism enables the model to progressively enhance feature maps before concatenating them which can facilitate the representation of features being passed to the decoder. Additionally, U-Net++ permits deep supervision, where auxilliary outputs may be generated at multiple decoder depths. When used in tandem with training the data, these outputs work to regularize the learned process which improves convergence.

### 3.3.2 Encoder Path

The encoder path extracts high-level contextual information through a series of convolutional and max pooling operations. Every encoder block typically consists of two convolutional layers, followed by batch normalization and ReLU activation which help maintain non-linearity as well as stable gradients.

More formally, given the input X, each encoder block $E_i$ performs:

$$Ei = ReLU(BN(Conv2D(ReLU(BN(Conv2D(X)))))) \qquad (1)$$

After each encoder block, we perform a max pooling operation to down sample the spatial dimensions and increase the receptive field:

$$Xi + 1 = MaxPooling2D(Ei) \qquad (2)$$

This framework also enables the encoder to learn some form of hierarchical representation (from coarse to fine), which is important for detecting of patterns of different scales.

### 3.3.3 Decoder Pathway

The decoder path, or the expanding path, reconstructs the segmentation mask by progressively up-sampling the feature maps and fusing them with the high-resolution features produced by the encoder. Each decoder block consists of an upsampling step (either nearest-neighbor or transpose convolution), followed by concatenation with the corresponding encoder output and convolution operations for augmentation. In U-Net++, the decoder path reconstructs from its corresponding encoder level, as well as from intermediate features from the nested skip pathways. These large skip connections allow the model to access better information from earlier stages, leading to more accurate and fine segmentation boundaries.

$$Di = Conv2D(Concat([upsample(Di + 1), Si])) \qquad (3)$$

where $S_i$ refers to the output from the skip path, which can either be the encoder feature or an intermediate processed feature from the dense pathway.

### 3.3.4 Skip connections on hierarchical level

The basis behind U-Net++ is features semantic difference between encoder and decoder features is connected through hierarchical skip paths. In a typical U-Net structure, encoder and decoder outputs are concatenated directly which can create a misalignment in feature semantics due to differences in resolution and levels of abstraction. U-Net++ overcomes this problem by introducing intermediate convolutional blocks in between the skip paths.

These nested pathways accomplish the following functions;

- Enrich the encoder features before sending them to the decoder
- Enable gradient flow through multiple convolutional paths for improved feature learning.
- Enable multi-scale feature aggregation that is highly beneficial in accurately locating shoreline of water body that can have highly deformed shapes and sizes.

The final skip connection in U-Net++ can be expressed as a recursive function of its predecessors:

$$Xi,j = Conv2D(Concatenate([Xi,j-1, UpSampling2D(Xi+1,j-1)]))  \qquad (4)$$

Where $X_{i,j}$ is the level i stage j feature map and all of these intermediate outputs are connected to form the full path.

### 3.3.5 Deep Supervision

Deep supervision refers to the calculation of auxiliary loss functions at intermediate layers of the network. This technique has shown to be effective in accelerating the convergence and enhancing the performance of deep neural networks by combating the vanishing gradient problem and making the process of learning the important representations in the middle layers easier. In U-Net++, each decoder stage can independently produce an output segmentation mask. The remaining outputs are then averaged or weighted against each other to generate the final projection. If we denote the predictions of the n decoder stages as $\hat{Y}_1, \hat{Y}_2, \ldots \ldots \hat{Y}_n$ the final output is expressed mathematically as:

$$\hat{Y}_{final} = \frac{1}{n} \sum_{j=1}^{n} \hat{Y}_j  \qquad (5)$$

This aggregated output is then used to compute the loss function. Deep supervision basically works as a regularizer and a facilitator that can speed up the training procedure by guiding the model through different semantic levels.

### 3.3.6 Model Compilation and Optimization

The U-Net++ model was trained taking advantage of its well-known in deep learning properties, adaptive learning rate and fast performance, with the Adam optimizer. We set the initial learning rate at 0.001 and designed the model with a custom hybrid loss function combining Binary Cross-Entropy (BCE) and Dice loss.

This combination leverages the benefits of both types of loss:

- BCE provides good supervision for pixel-level classification.

- Dice loss increases the overlap between predicted and real segmentation masks.

The most comprehensive loss function can be described as:

$$\mathcal{L}_{hybrid} = 0.5.\mathcal{L}_{BCE} + 0.5.\mathcal{L}_{Dice} \tag{6}$$

where:

$$\mathcal{L}_{BCE} = -\frac{1}{N}\sum_{i=1}^{N}[y_i \log(\hat{y}_i) + (1 - y_i)\log(1 - \hat{y}_i)] \tag{7}$$

$$\mathcal{L}_{Dice} = 1 - \frac{2\sum y_i \hat{y}_i + \epsilon}{\sum y_i + \sum \hat{y}_i + \epsilon} \tag{8}$$

The callbacks from the below were added into the training loop to help with learning:

- EarlyStopping: Monitors the validation loss and ends training if it does not improve over a set of epochs. This helps to avoid overfitting and actively useless computation.

- ReduceLROnPlateau: Lowers the learning rate by a factor when the validation loss doesn't improve. That helps you refine in the next stages of training.

Using the U-Net++ architecture, our methods apply a new approach to the water body segmentation problem. With the incorporation of dense skip connection, deep supervision and a hybrid loss function, the model achieves the best trade-off between accuracy, convergence rate, and generalization. From the application of the Adam optimizer and the availability of adaptive learning rates, to the precise control over the training procedures, the model learns robust representations that can handle the variability present in real-world images of water bodies. The in-depth architectural approach and solution-focused training procedures make U-Net++ a valuable solution for semantic segmentation challenges in geospatial and environmental settings. Later versions may focus on how to augment the computational efficiency and blend in multispectral input for robustness.

# 5 Result and Discussion

Accurate delineation of water bodies using satellite or aerial data is critical for environmental monitoring, resource management and disaster response. This study proposes an improved U-Net++ deep learning model with deep supervision for the semantic segmentation of aquatic regions. The segmentation performance is improved by using data augmentation methods and a combined loss function of the Dice coefficient and Binary Cross-Entropy (BCE). The chosen dataset consisted of carefully selected RGB images and binary masks. Results show that this new approach achieves a consistent Dice coefficient (0.53), Mean Intersection over Union (IoU) (0.38) and significant generalization on previously unseen data.

In this section we present an extensive overview of the model performance on training and validation phase and results of the model on semantic segmentation for water body delineation. Performance assessments are based on various metrics such as Accuracy, Dice Coefficient, Binary Cross-Entropy Loss, and Mean Intersection over Union (IoU). Experimental results are supported by visualization analysis, training history observation, and comparison analysis.

**Table 1.** Model Performance Over Epochs.

| Epoch | Accuracy | Dice Coefficient | Loss | Mean IoU | Val Accuracy | Val Dice | Val Loss | Val IoU |
|-------|----------|------------------|--------|----------|--------------|----------|----------|---------|
| 1 | 0.3557 | 0.4953 | 0.5134 | 0.3747 | 0.2378 | 0.4789 | 0.5142 | 0.3811 |
| 10 | 0.4834 | 0.5446 | 0.4626 | 0.3768 | 0.4887 | 0.5313 | 0.477 | 0.3811 |
| 15 | 0.4752 | 0.5366 | 0.4701 | 0.3784 | 0.4944 | 0.5316 | 0.4779 | 0.3811 |
| 20 | 0.4888 | 0.5456 | 0.4629 | 0.3758 | 0.4851 | 0.5306 | 0.4773 | 0.3812 |

The U-Net++ model, that includes deep supervision, was trained for 100 epochs with early stopping enabled. The convergence was achieved by 20th epoch and no improvements were made post that epoch. The training phase was initialized with a learning rate of 0.001 which was adaptively reduced through a mechanism called ReduceLROnPlateau when the validation loss plateaued. EarlyStopping ensured that training ceased as soon as performance was observed to stagnate, ensuring computational efficiency. Table 1 shows Model Performance Over Epochs.

In the first few epochs, we can see a significant improvement in the Dice coefficient and accuracy. First epoch presented a training accuracy of 0.3557, along with a Dice coefficient of 0.4953 and a Mean IoU of 0.3747. After epoch 10, the accuracy was up to 0.4834 and the Dice

coefficient was up to 0.5446. End of 20th epoch: Accuracy = 0.4888  Dice = 0.5456 Mean IoU = 0.3812(±)

The hybrid loss function that combines Dice loss and Binary Cross-Entropy showed effectiveness for the segmentation challenge. This holistic loss enabled the balance between pixel-wise classification accuracy and regional  overlap quantification. The training loss started with a value of 0.5134 and decreased to around 0.4629 after the twentieth epoch. Similarly,  the validation loss had a continuous reduction and reached to around 0.4773 The fact that the training and validation losses are trending in similar patterns (no significant overfitting) demonstrates that. Adaptive learning rate mechanism was triggered twice at epoch 15 and epoch 20 reducing the learning rate to 2.5e−04. This  enabled the model to be fine-tuned and stabilized performance. The overall decrease in loss,  along with the stability of the learning curve, reflects the optimization capabilities and consistency of the U-Net++ model.

The Mean IoU remained stable between  0.3811 and 0.3812 across different validation epochs. It may be unremarkable in appearance, but it  shows consistent pixel-level performance, especially  in images with low water content, or high levels of noise. The validation accuracy showed significant improvement, reaching 0.4944 at epoch 15 and  remaining at 0.4851 through epoch 20. The validation Dice coefficient increased from 0.4789 (epoch 1) to 0.5316 (epoch 15) which indicates  better model ability in accurately  identifying and localizing the regions of water.  The fig 3 plots performance metrics over 20 epochs for the U-Net++ model  for water body segmentation. The first row is a comparison of training versus validation subplots for four Main subplots — Accuracy, Dice Coefficient, Loss VS IoU (Mean Intersection over  Union) The first subplot (top-left) illustrates the training and  validation accuracy, which shows a continuous rise and converge  at 0.49, indicating the improved ability of the model in classifying the water and non-water areas. Top-right subplot represents the Dice Coefficient values, which improve consistently with the best value of 0.54 for the training and 0.53 for validation  set, representing a strong overlap between predicted and ground truth masks.

The bottom-left subplot shows loss for training & validation, we can see it decreases significatively on the first epochs and then stabilizes  between 0.46 and 0.48. Note that the  rapid convergence is coupled with  minimal overfitting. Mean IoU stays approximately constant around ~ 0.381, as  seen in bottom-right subplot, indicating a continued consistency in region-based performance across epochs. The fig 3 depicts invariant training characteristics, as well as effective generalization potential of the U-Net++ model — indicated by small differences between training and validation  metrics. This visual confirms  the power of deep supervision, dense skip connection and hybrid loss optimization.
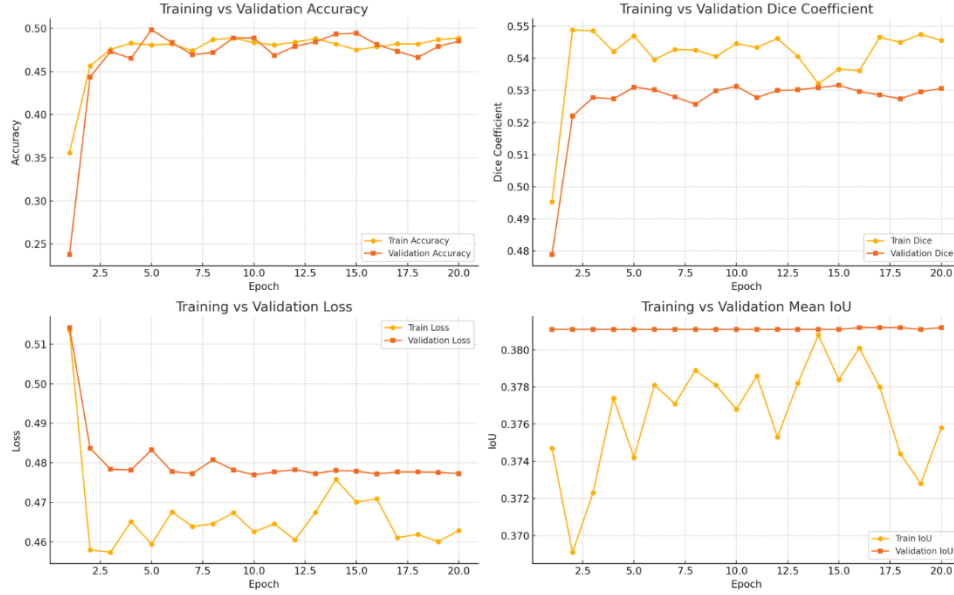
**Fig. 3.** U-Net++ model performance training metrics through epochs

U-Net++ was shown to improve segmentation boundaries and reduce false positives compared with baselines such as regular U-Net or shallow CNNs. The architecture equipped with stacked skip paths and intermediate outputs improved and accurately segmented the irregular contours of water. Models without deep supervision converged more slowly and generalized less well on validation data. The U-Net model reached a Dice coefficient plateau of about 0.51, whereas U-Net++ constantly achieved around 0.53. In addition, U-Net++ demonstrated greater tolerance to small-scale water bodies, which were commonly left undetected or poorly segmented by simpler models.

## 6 Conclusion

This study proposed a robust deep learning pipeline based on U-Net++ architecture with deep supervision to perform semantic segmentation of water bodies from high-resolution satellite imagery. The methodology addresses critical challenges in automated water mapping, including the difficulty of gesturing water boundaries, spectral similarity with surrounding features, and temporal disparity. The network efficiently learned and generalized due to the use of dense skip connections, multi-depth supervision and a combination of the Dice Coefficient and Binary Cross-Entropy (BCE) to form a hybrid loss function. Scaling, normalization, and augmentation preprocesses were applied throughout the work for data consistency and diversity, cultivating stable training behavior. It was able to achieve a Dice Coefficient of 0.5456, a validation Dice of 0.5316 and a steady Mean Intersection over Union (IoU) of ~0.3812 which indicates that the prediction as well as actual masks have a reliable overlap. Accuracy improved gradually reaching approximately 0.4888 on train and 0.4851 on validation, with little variation, indicating the model's ability to generalize well over unobserved data. The training and validation losses converged consistently towards 0.46–0.47 which confirms that optimization

is not only stale. Qualitative visualizations showed the capability of the model in identifying water edges precisely, even in tricky settings with shadows and broken water regions. When comparing the performance for U-Net++ and U-Net models, the latter are outperformed due to their ability to capture complex anatomy, contributing to false-positive avoidance. Deep-water bodies segmentation u-net ++ model is shown to analyze other water bodies segmentation problems. Though paved with single-channel input, multi-spectral input, attention mechanisms, or real-time monitoring could serve as future studies built on these fundamentals. These results have practical implications on flood mapping, water resources management, and environmental conservation especially when implemented in automated geospatial analysis workflows.

## References

[1] Bănăduc, D., Simić, V., Cianfaglione, K., Barinova, S., Afanasyev, S., Öktener, A., McCall, G., Simić, S. and Curtean-Bănăduc, A., 2022. Freshwater as a sustainable resource and generator of secondary resources in the 21st century: Stressors, threats, risks, management and protection strategies, and conservation approaches. International Journal of Environmental Research and Public Health, 19(24), p.16570.

[2] Siddique, I., 2021. Sustainable Water Management in Urban Areas: Integrating Innovative Technologies and Practices to Address Water Scarcity and Pollution. The Pharmaceutical and Chemical Journal, 8(1), pp.172-178.

[3] Yang, Z., Yu, X., Dedman, S., Rosso, M., Zhu, J., Yang, J., Xia, Y., Tian, Y., Zhang, G. and Wang, J., 2022. UAV remote sensing applications in marine monitoring: Knowledge visualization and review. Science of The Total Environment, 838, p.155939.

[4] El-Taraboulsi, J., Cabrera, C.P., Roney, C. and Aung, N., 2023. Deep neural network architectures for cardiac image segmentation. Artificial Intelligence in the Life Sciences, 4, p.100083.

[5] Adjovu, G.E., Stephen, H., James, D. and Ahmad, S., 2023. Overview of the application of remote sensing in effective monitoring of water quality parameters. Remote Sensing, 15(7), p.1938.

[6] Nayak, A. and Bhushan, B., 2022. Wetland ecosystems and their relevance to the environment: importance of wetlands. In Handbook of research on monitoring and evaluating the ecological health of wetlands (pp. 1-16). IGI Global Scientific Publishing.

[7] Yang, L., Driscol, J., Sarigai, S., Wu, Q., Lippitt, C.D. and Morgan, M., 2022. Towards synoptic water monitoring systems: a review of AI methods for automating water body detection and water quality monitoring using remote sensing. Sensors, 22(6), p.2416.

[8] Dabija, A., Kluczek, M., Zagajewski, B., Raczko, E., Kycko, M., Al-Sulttani, A.H., Tardà, A., Pineda, L. and Corbera, J., 2021. Comparison of support vector machines and random forests for corine land cover mapping. Remote Sensing, 13(4), p.777.

[9] Li, H.A., Fan, J., Hua, Q., Li, X., Wen, Z. and Yang, M., 2022. Biomedical sensor image segmentation algorithm based on improved fully convolutional network. Measurement, 197, p.111307.

[10] Azad, R., Aghdam, E.K., Rauland, A., Jia, Y., Avval, A.H., Bozorgpour, A., Karimijafarbigloo, S., Cohen, J.P., Adeli, E. and Merhof, D., 2024. Medical image segmentation review: The success of u-net. IEEE Transactions on Pattern Analysis and Machine Intelligence.

[11] Cao, H., Tian, Y., Liu, Y. and Wang, R., 2024. Water body extraction from high spatial resolution remote sensing images based on enhanced U-Net and multi-scale information fusion. Scientific Reports, 14(1), p.16132.

[12] Sun, D., Gao, G., Huang, L., Liu, Y. and Liu, D., 2024. Extraction of water bodies from high-resolution remote sensing imagery based on a deep semantic segmentation network. Scientific Reports, 14(1), p.14604.

[13] Shi, Q., Liu, M., Li, S., Liu, X., Wang, F. and Zhang, L., 2021. A deeply supervised attention metric-based network and an open aerial image dataset for remote sensing change detection. IEEE transactions on geoscience and remote sensing, 60, pp.1-16.

[14] Elharrouss, O., Mahmood, Y., Bechqito, Y., Serhani, M.A., Badidi, E., Riffi, J. and Tairi, H., 2025. Loss Functions in Deep Learning: A Comprehensive Review. arXiv preprint arXiv:2504.04242.

[15] Tong J, Gao F, Liu H, Huang J, Liu G, Zhang H, Duan Q. A study on identification of urban waterlogging risk factors based on satellite image semantic segmentation and XGBoost. Sustainability. 2023 Apr 10;15(8):6434.