Hydrology and
Earth System
Sciences

# Eye of Horus: a vision-based framework for real-time water level measurement

**Seyed Mohammad Hassan Erfani**[1], **Corinne Smith**[2], **Zhenyao Wu**[3], **Elyas Asadi Shamsabadi**[4], **Farboud Khatami**[1], **Austin R. J. Downey**[1,2], **Jasim Imran**[1], and **Erfan Goharian**[1]

[1]Department of Civil & Environmental Engineering, University of South Carolina, Columbia, SC 29208, USA
[2]Department of Mechanical Engineering, University of South Carolina, Columbia, SC 29208, USA
[3]Department of Computer Science & Engineering, University of South Carolina, Columbia, SC 29201, USA
[4]School of Civil Engineering, Faculty of Engineering, The University of Sydney, Sydney, NSW 2006, Australia

**Correspondence:** Erfan Goharian (goharian@cec.sc.edu)

**Abstract.** Heavy rains and tropical storms often result in floods, which are expected to increase in frequency and intensity. Flood prediction models and inundation mapping tools provide decision-makers and emergency responders with crucial information to better prepare for these events. However, the performance of models relies on the accuracy and timeliness of data received from in situ gaging stations and remote sensing; each of these data sources has its limitations, especially when it comes to real-time monitoring of floods. This study presents a vision-based framework for measuring water levels and detecting floods using computer vision and deep learning (DL) techniques. The DL models use time-lapse images captured by surveillance cameras during storm events for the semantic segmentation of water extent in images. Three different DL-based approaches, namely PSPNet, TransUNet, and SegFormer, were applied and evaluated for semantic segmentation. The predicted masks are transformed into water level values by intersecting the extracted water edges, with the 2D representation of a point cloud generated by an Apple iPhone 13 Pro lidar sensor. The estimated water levels were compared to reference data collected by an ultrasonic sensor. The results showed that Seg-Former outperformed other DL-based approaches by achieving 99.55 % and 99.81 % for intersection over union (IoU) and accuracy, respectively. Moreover, the highest correlations between reference data and the vision-based approach reached above 0.98 for both the coefficient of determination ($R^2$) and Nash–Sutcliffe efficiency. This study demonstrates the potential of using surveillance cameras and artificial intelligence for hydrologic monitoring and their integration with existing surveillance infrastructure.

## 1 Introduction

Flood forecasts and flood inundation mapping (FIM) can play an important role in saving human lives and reducing damage by providing timely information for evacuation planning, emergency management, and relief efforts (Gebrehiwot et al., 2019). These models and tools are designed to identify and predict inundation areas and the severity of damage caused by storm events. Two primary sources of data for these models are in situ gaging networks and remote sensing. For example, in situ stream gages, such as those operated by the United States Geological Survey (USGS) provide useful streamflow information like water height and discharge at monitoring sites (Turnipseed and Sauer, 2010). However, they cannot provide an adequate spatial resolution of streamflow characteristics (Lo et al., 2015). The limitation of in situ stream gages is further exacerbated by the lack of systematic installation along the waterways and accessibility issues (Li et al., 2018; King et al., 2018). Satellite data and remote sensing can complement in situ gage data by providing information at a larger spatial scale (Alsdorf et al., 2007). However, continuous monitoring of data for a region of interest remains to be a problem due to the limited revisit intervals of satellites, cloud cover, and systematic departures or biases (Panteras and Cervone, 2018). Crowdsourcing methods

have gained attention as a potential solution but their reliability is questionable (Schnebele et al., 2014; Goodchild, 2007; Howe, 2008). To address these limitations and enhance real-time monitoring capabilities, surveillance cameras are investigated here as a new source of data for hydrologic monitoring and flood data collection. However, this requires a significant investment in computer vision (CV) and artificial intelligence (AI) techniques to develop reliable methods for detecting water in surveillance images and translating that information into numerical data.

Recent advances in CV offer new techniques for processing image data for the quantitative measurements of physical attributes from a site (Forsyth and Ponce, 2002). However, there is limited knowledge of how visual information can be used to estimate physical water parameters using CV techniques. Inspired by the principle of the float method, Tsubaki et al. (2011) used different image processing techniques to analyze images captured by closed-circuit television (CCTV) systems installed for surveillance purposes to measure the flow rate during flood events. In another example, Kim et al. (2011) proposed a method for measuring water level by detecting the borderline between a staff gauge and the surface of water based on image processing of the captured image of the staff gage installed in the middle of the river. As the use of images for environmental monitoring becomes more popular, several studies have investigated the source and magnitude of errors common in image-based measurement systems, such as the effect of image resolution, lighting effects, perspective, lens distortion, water meniscus, and temperature changes (Elias et al., 2020; Gilmore et al., 2013). Furthermore, proposed solutions to resolve difficulties originating from poor visibility have been developed to better identify readings on staff gages (Zhang et al., 2019). Recently, deep learning (DL) has become prevalent across a wide range of disciplines, particularly in applied sciences such as CV and engineering.

DL-based models have been utilized by the water resources community to determine the extent of water and waterbodies visible in images captured by surveillance camera systems. These models can estimate the water level (Pally and Samadi, 2022). In a similar vein, Moy de Vitry et al. (2019) and Vandaele et al. (2021) employed a DL-based approach to identify floodwater in surveillance footage and introduced a novel qualitative flood index, SOFI, to determine water level fluctuations. SOFI was calculated by taking the aspect ratio of the area of the water surface detected within an image to the total area of the image. However, these types of methods, which make prior assumptions and estimate water level fluctuation roughly, cannot serve as a vision-based alternative for measuring streamflow characteristics. More systematic studies adopted photogrammetry to reconstruct a high-quality 3D model of the environment with a high spatial resolution to have a precise estimation of real-world coordination while measuring streamflow rate and stage. For example, Eltner et al. (2018, 2021) introduced a method based on structure from motion (SfM) and photogrammetric techniques to automatically measure the water stage using low-cost camera setups.

Advances in photogrammetry techniques enable 3D surface reconstruction with a high temporal and spatial resolution. These techniques are adopted to build 3D surface models from RGB imagery (Westoby et al., 2012; Eltner and Schneider, 2015; Eltner et al., 2016). However, most of the photogrammetric methods are still expensive as they rely on differential global navigation satellite systems (DGNSS), ground control points (GCPs), commercial software, and data processing on an external computing device (Froideval et al., 2019). A lidar scanner, on the other hand, is now easily available since the introduction of the iPad Pro and iPhone 12 Pro in 2020 by Apple. This device is the first smartphone equipped with a native lidar scanner and offers a potential paradigm shift in digital field data acquisition, which puts these devices at the forefront of smartphone-assisted fieldwork (Tavani et al., 2022). So far, the iPhone lidar sensor has been used in different studies such as forest inventories (Gollob et al., 2021) and coastal cliff sites (Luetzenburg et al., 2021). The availability of lidar sensors to build 3D environments and advancements in DL-based models offer a great potential to produce numerical information from ground-based imageries.

This paper presents a vision-based framework for measuring water levels from time-lapse images. The proposed framework introduces a novel approach by utilizing the iPhone lidar sensor as a laser scanner, which is commonly available on consumer-grade devices, for scanning and constructing a 3D point cloud of the region of interest. During the data collection phase, time-lapse images and ground truth water level values were collected using an embedded camera and ultrasonic sensor. The water extent in the captured images was determined automatically using semantic segmentation DL-based models. For the first time, the performance of three different state-of-the-art DL-based approaches, including convolutional neural networks (CNNs), hybrid CNN transformer, and transformers–multilayer perceptron (MLP), was evaluated and compared. CV techniques were applied for camera calibration, pose estimation of the camera setup in each deployment, and 3D–2D reprojection of the point cloud onto the image plane. Finally, $K$-nearest neighbor (KNN) was used to find the nearest projected (2D) point cloud coordinates to the waterline on the riverbanks, for estimating the water level in each time-lapse image.

## 2   Deep learning architectures

Since this study tends to cover a wide range of DL approaches, this section solely focuses on reviewing different DL-based architectures. So far, different DL networks have been applied and evaluated for semantic segmentation of the waterbodies within the RGB images captured by cameras

(Erfani et al., 2022). All existing semantic segmentation approaches – CNN and transformer-based – share the same objective of classifying each pixel of a given image but differ in the network design.

CNN-based models were designed to imitate the recognition system of primates (Shamsabadi et al., 2022), while possessing different network designs such as low-resolution representation learning (Long et al., 2015; Chen et al., 2017), high-resolution representation recovery (Badrinarayanan et al., 2015; Noh et al., 2015; Lin et al., 2017), contextual aggregation schemes (Yuan and Wang, 2018; Zhao et al., 2017; Yuan et al., 2020), feature fusion and refinement strategy (Lin et al., 2017; Huang et al., 2019; Li et al., 2019; Zhu et al., 2019; Fu et al., 2019). CNN-based models follow local to global features in different layers of the forward pass, which used to be thought of as a general intuition of the human recognition system. In this system, objects are recognized through the analysis of texture and shape-based clues–local and global representations and their relationship in the entire field of view. Recent research, however, shows that significant differences exist between the visual behavioral system of humans and CNN-based models (Geirhos et al., 2018b; Dodge and Karam, 2017; De Cesarei et al., 2021; Geirhos et al., 2020, 2018a) and reveal higher sensitivity of the visual systems in humans to global features rather than local ones (Zheng et al., 2018). This fact drew attention to models that focus on the global context in their architectures.

Developed by Dosovitskiy et al. (2020), Vision Transformer (ViT) was the first model that showed promising results on a computer vision task (image classification) without using convolution operation in its architecture. In fact, ViT adopts "transformers", as a self-attention mechanism, to improve accuracy. Transformer was initially introduced for sequence-to-sequence tasks such as text translation (Vaswani et al., 2017). However, as applying the self-attention mechanism to all image pixels is computationally expensive, the transformer-based models could not compete with the CNN-based models until the introduction of ViT architecture which applies self-attention calculations to the low-dimension embedding of small patches originating from splitting the input image to extract global contextual information. Successful performance of ViT on image classification inspired several subsequent works on transformer-based models for different computer vision tasks (Liu et al., 2021).

In this study, three different DL-based approaches including CNN, hybrid CNN transformer, and transformers–multilayer perceptron (MLP) were trained and tested for semantic segmentation of water. For these approaches, the selected models were PSPNet (Zhao et al., 2017), TransUNet (Chen et al., 2021), and SegFormer (Xie et al., 2021), respectively. The performance of these models is evaluated and compared using conventional metrics, including class-wise intersection over union (IoU) and per-pixel accuracy (ACC).
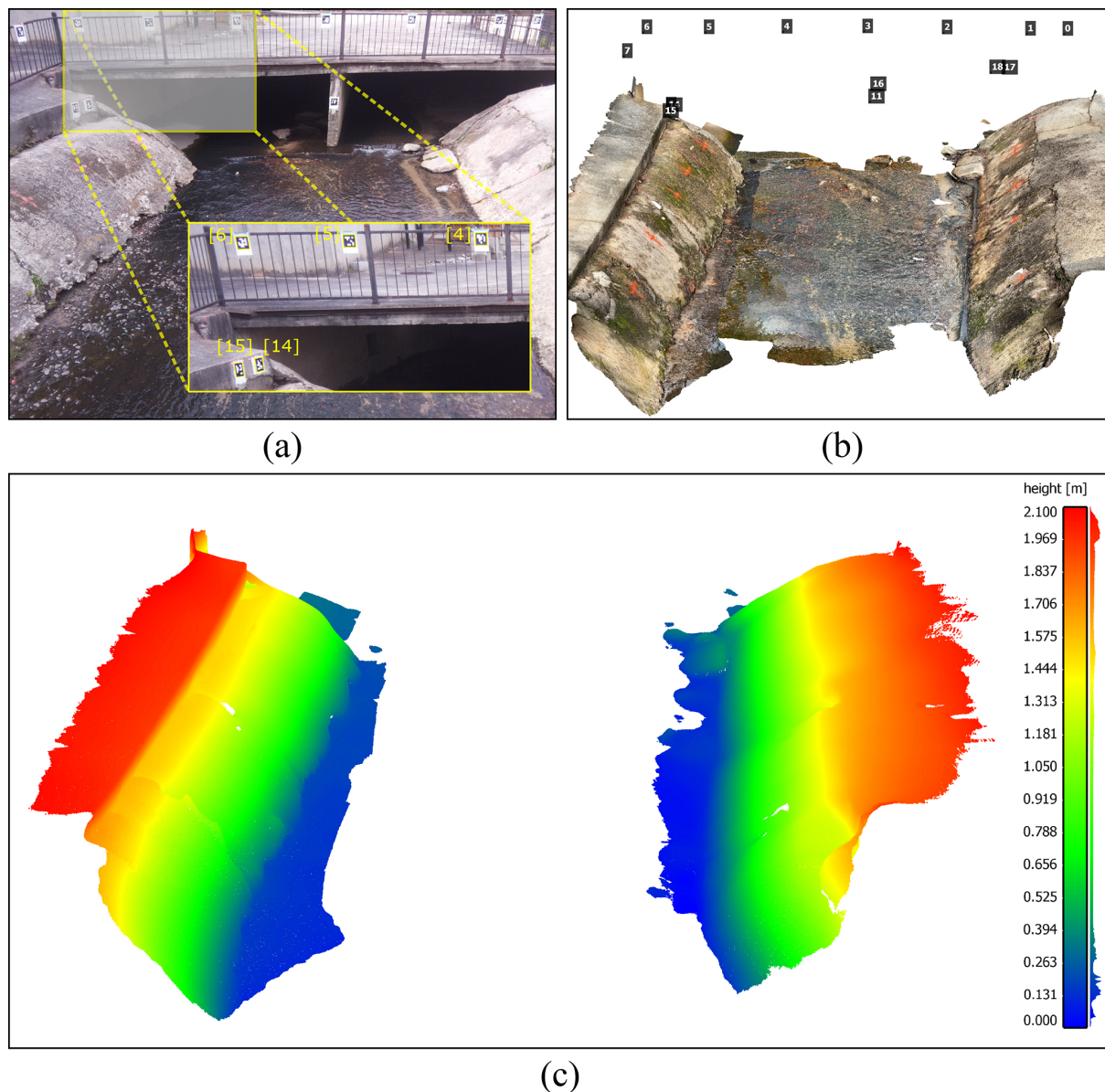
# 3 Study area

In order to evaluate the performance of the proposed framework for measuring the water levels in rivers and channels, a time-lapse camera system has been deployed at Rocky Branch, South Carolina. This creek is approximately 6.5 km long and collects stormwater from the University of South Carolina campus and the City of Columbia. Rocky Branch is subjected to rapid changes in water flow and discharges into the Congaree River (Morsy et al., 2016). The observation site is located within the University of South Carolina campus behind 300 Main Street (see Fig. 1a).

An Apple iPhone 13 Pro lidar sensor was used to scan the region of interest. Although there is no official information about the technology and hardware specifications, Gollob et al. (2021) reports that the lidar module operates at the 8XX nm wavelength and consists of an emitter (vertical cavity surface-emitting laser with diffraction optics element, VCSEL DOE) and a receptor (single photon avalanche diode array-based near-infrared complementary metal oxide semiconductor image sensor, SPAD NIR CMOS) based on direct-time-of-flight technology. Comparisons between the Apple lidar sensor and other types of laser scanners including handheld, industrial, and terrestrial have been conducted by several recent studies (Mokroš et al., 2021; Vogt et al., 2021). Gollob et al. (2021) tested and reported the performance of a set of eight different scanning apps and found three applications including 3D Scanner App, Polycam, and SiteScape suitable for actual practice tests. The objective of this study is not the evaluation of the iPhone lidar sensor and app performance. Therefore, the 3D Scanner App (LAAN LABS, 2022) was used with the following settings: confidence, high; range, 5.0 m; masking, none; and resolution, 5 mm for scanning and 3D reconstruction processing. The scanned 3D point cloud and its corresponding scalar field are shown in Fig. 1b and c, respectively.

As the lidar scanner settings were set at the highest level of accuracy and computational demand, scanning the whole region of interest at the same time was not possible. So, the experimental region was divided into several sub-regions and scanned in multiple steps. In order to assemble the sub-region lidar scans, several GCPs were considered in the study area. These GCPs were measured by a total station (Topcon GM Series) and used as landmarks to align distinct 3D point clouds with each other and create an integrated point cloud encompassing the entirety of the study area.

Moreover, several ArUco markers were installed for estimating camera (extrinsic) parameters. In each setup deployment, these parameters should be recalculated (additional information can be found in Sect. 4.3). Since it was not possible to accurately measure the real-world coordination of ArUco markers by the lidar scanner, the coordinates of the top-left corner of markers were also measured by the surveying total station. To establish a consistent coordinate system, the 3D point cloud scanned for each sub-region was transformed

(a)

(b)

(c)

**Figure 1.** Study area of the Rocky Branch creek. **(a)** View of the region of interest, **(b)** the scanned 3D point cloud of the region of interest including an indication of the ArUco markers' locations, and **(c)** the scalar field of left and right banks of Rocky Branch in the region of interest (the color bar and the frequency distribution of $z$ values for the captured points are shown on the right side).
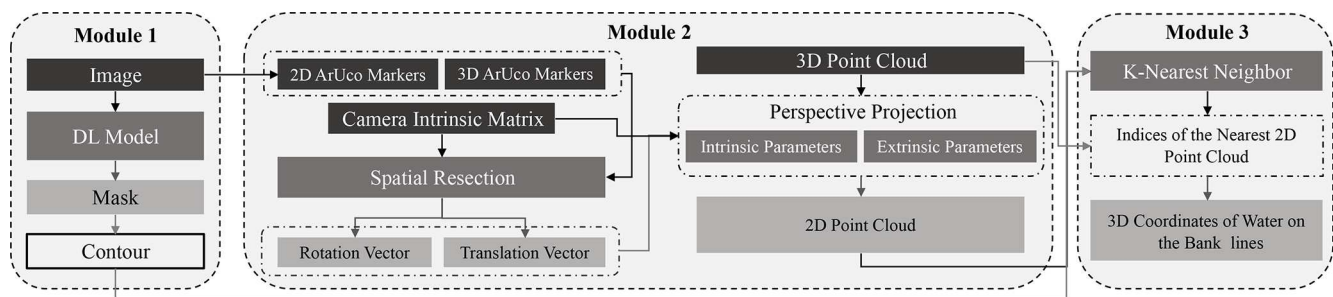
into the total station's coordinate system. The real-world coordinates of ArUco markers were then added to the 3D point cloud (see Fig. 1b).
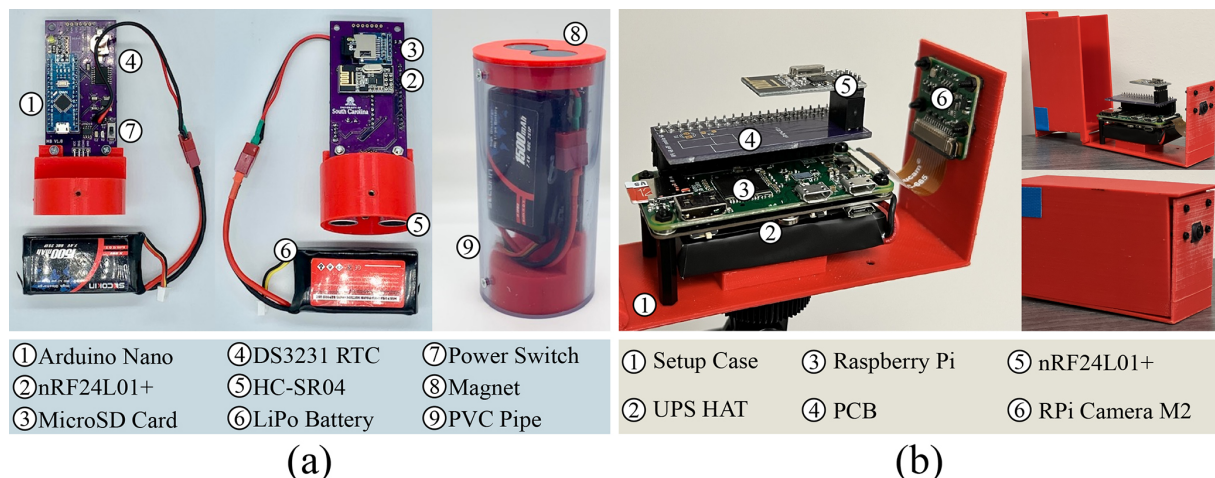
## 4   Methodology

This study introduces the Eye of Horus, a vision-based framework for hydrologic monitoring and real-time water level measurements in bodies of water. The proposed framework includes three main components. The first step is designing two deployable setups for data collection. These se-

tups consist of a programmable time-lapse camera run by Raspberry Pi and an ultrasonic sensor run by Arduino. After collecting data, the first phase (Module 1) involves configuring and training DL-based models for semantic segmentation of water in the captured images. In the second phase (Module 2), CV techniques for camera calibration, spatial resection, and calculating projection matrix are discussed. Finally, in the third phase (Module 3), an machine learning (ML)-based model uses the information achieved by CV models to find the relationships between real-world coordinates of water level in the captured images (see Fig. 2).

**Figure 2.** The Eye of Horus workflow includes three main modules starting from processing images captured by the time-lapse camera to estimating water level by projecting the waterline on riverbanks using CV techniques.



**Figure 3.** Data acquisition devices. **(a)** Beena, run by Raspberry Pi (Zero W) for capturing time-lapse images of the river scene, and **(b)** Aava, run by Arduino Nano for measuring water level correspondence.

## 4.1 Data acquisition

Two different single-board computers (SBCs) were used in this study: Raspberry Pi (Zero W) for capturing time-lapse images of a river scene and Arduino (Nano 3.x) for measuring water level as the ground truth data. These devices were designed to communicate with each other, i.e., to trigger the other to start or stop recording. While capturing time-lapse images, the Pi camera device triggers the ultrasonic sensor to measure the corresponding water level. The camera device is equipped with the Raspberry Pi Camera Module 2 which has a Sony IMX219 8 MP sensor. This sensor is able to capture an image size of $4256 \times 2832$ pixels. However, in this study, the image resolution was set to $1920 \times 1440$ pixels to balance image quality and computational cost in subsequent image processing steps. This setup is also equipped with a 1200 mAh UPS lithium battery power module to provide uninterrupted power to the Pi SBC (see Fig. 3a).

The Arduino-based device records the water level. The design is based on a drone-deployable sensor created by Smith et al. (2022). The nRF24L01+ single-chip 2.4 GHz transceiver allows the Arduino and Raspberry Pi to communicate via radio frequency (RF). The chip is housed in both packages and the channel, pipe addresses, data rate, and transceiver–receiver configuration are all set in the software. The HC-SR04 ultrasonic sensor is mounted on the base of the Arduino device and provides a contactless water level measurement. Two permanent magnets at the top of the housing attach to a ferrous structure and allow the ultrasonic sensor to be suspended up to 14 ft over the surface of the water. The device also includes a microSD card module and DS3231 real-time clock, which enable data logging and storage on-device as well as transmission. The device is powered by a rechargeable 7.4 V 1500 mAh lithium polymer battery (see Fig. 3b).

The Arduino device waits to receive a ping from the Raspberry Pi device to initiate data collection. The ultrasonic sensor measures the distance from the sensor transducer to the surface of the water. The nRF24L01+ transmits this distance to the Raspberry Pi device and saves the measurement and a time stamp from the real-time clock to an onboard microSD card. This acts as backup data storage, in case transmission to the Raspberry Pi fails. The nRF24L01+ RF transceivers have an experimentally determined range of up to 30 ft, which al-

**Table 1.** The configuration of models trained and tested in this study.

| Model names | Params (M) | Total size (MB) | Batch size $(B, H, W, C)$ | Loss function | Optimizer | LR |
|---|---|---|---|---|---|---|
| PSPNet | 66.2 | 7178 | $2 \times 500 \times 500 \times 3$ | Binary cross entropy | SGD | $2.50 \times 10^{-4}$ |
| TransUNet | 20.1 | 6017 | $2 \times 448 \times 448 \times 3$ | Cross entropy + dice | SGD | $2.50 \times 10^{-4}$ |
| SegFormer-B0 | 3.7 | 2217 | $2 \times 512 \times 512 \times 3$ | Cross entropy | AdamW | $6.00 \times 10^{-5}$ |
| SegFormer-B5 | 82.0 | 27 666 | $2 \times 1024 \times 1024 \times 3$ | Cross entropy | AdamW | $6.00 \times 10^{-5}$ |

lows flexibility in the relative placement of the camera to the measuring site.

A dataset for semantic segmentation was created by collecting images from a specific region of interest at different times of the day and under various flow regimes. This dataset includes 1172 images, with manual annotations of the streamflow in the creek for all of them. The dataset is further divided into 812 training images, 124 validation images, and 236 testing images.

## 4.2 Deep learning model for water segmentation

The water extent can be automatically determined on the 2D image plane with the help of DL-based models. The task of semantic segmentation was applied within the framework of this study to delineate the waterline on the left and right banks of the channel. Three different DL-based models were trained and tested in this study. PSPNet, the first model, is a CNN-based semantic segmentation multi-scale network that can better learn the global context representation of a scene (Zhao et al., 2017). ResNet-101 (He et al., 2016) was used as the backbone of this model to encode input images into the features. ResNet architecture takes the advantage of "residual blocks" that assist the flow of gradients during the training stage allowing effective training of deep models even up to hundreds of layers. These extracted features are then fed into a pyramid pooling module in which feature maps produced by small to large kernels are concatenated to distinguish patterns of different scales (Minaee et al., 2022).

TransUNet, the second model, is a U-shaped architecture that employs a hybrid of CNN and transformers as the encoder to leverage both the local and global contexts for precise localization and pixel-wise classification (Chen et al., 2021). In the encoder part of the network, CNN is first used as a feature extractor to generate a feature map for the input image, which is then fed into transformers to extract long-range dependencies. The resulting features are upsampled in the decoding path and combined with detailed high-resolution spatial information skipped from the CNN to make estimations on each pixel of the input image.

SegFormer, the third model, unifies a novel hierarchical transformer, which does not require the positional encodings used in standard transformers, and multilayer perceptron (MLP) performs efficient segmentation (Xie et al.,

2021). The hierarchical transformer introduced in the encoder of this architecture gives the model the attention ability to multi-scale features (high-resolution fine- and low-resolution coarse information) in the spatial input without the need for positional encodings that may adversely affect a model's performance when testing on a different resolution from training. Moreover, unlike other segmentation models that typically use deconvolutions in the decoder path, a lightweight MLP is employed as the decoder of this network that inputs the features extracted at different stages of the encoder to generate a prediction map faster and more efficiently. Two different variants, i.e., SegFormer-B0 and SegFormer-B5, were applied in this study. The configuration of the models implemented in this study is elaborated in Table 1. The total number of parameters (Params), occupied memory size on GPU (total size), and input image size (batch size) are reported in million (M), megabyte (MB), and batch size × height × width × channel $(B, H, W, C)$, respectively.

The models were implemented using PyTorch. During the training procedure, the loss function, optimizer, and learning rate were set individually for each model based on the results of preliminary runs used to find the optimal hyperparameters. In the case of PSPNet and TransUNet, the base learning rate was set to $2.5 \times 10^{-4}$ and decayed using the poly-policy (Zhao et al., 2017). These networks were optimized using stochastic gradient descent (SGD) with a momentum of 0.9 and weight decay of 0.0001. For SegFormer (B0 and B5), a constant learning rate of $6.0 \times 10^{-5}$ was used, and the networks were trained with the AdamW optimizer (Loshchilov and Hutter, 2017). All networks were trained for 30 epochs with a batch size of two. The training data for PSPNet and TransUNet were augmented with horizontal flipping, random scaling, and random cropping.

## 4.3 Projective geometry

In this study, CV techniques are used for different purposes. First, CV models were used for camera calibration. They include focal length, optical center, radial distortion, camera rotation, and translation. These parameters provide the information (parameters or coefficients) about the camera that is required to determine the relationship between 3D object points in the real-world coordinate system and its corresponding 2D projection (pixel) in the image captured by

that calibrated camera. Generally, camera calibration models estimate two kinds of parameters. First, the intrinsic parameters of the camera (e.g., focal length, optical center, and radial distortion coefficients of the lens). Second, extrinsic parameters (referring to the orientation – rotation and translation – of the camera) with respect to the real-world coordinate system.

To estimate the camera intrinsic parameters, built-in OpenCV was applied for camera calibration using a 2D checkerboard (Bradski, 2000). The focal length ($f_x$, $f_y$), optical centers ($c_x$, $c_y$), and the skew coefficient ($s$) can be used to create a camera intrinsic matrix **K**:

$$\mathbf{K} = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}. \tag{1}$$

The camera extrinsic parameters were determined using the pose computation problem, Perspective-n-Point (PnP), which consists of solving for the rotation, and translation that minimizes the reprojection error from 2D–3D point correspondences (Marchand et al., 2015). The PnP estimates the extrinsic parameters given a set of "object points", their corresponding "image projections", and the camera intrinsic matrix and the distortion coefficients. The camera extrinsic parameters can be represented as a combination of a $3 \times 3$ rotation matrix **R** and a $3 \times 1$ translation vector **t**:

$$[\mathbf{R}|\boldsymbol{t}] = \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \end{bmatrix}. \tag{2}$$

Equation (3) represents the "projection matrix" in a homogeneous coordinate system. The projection matrix consists of two parts: the intrinsic matrix (**K**), containing intrinsic parameters, and the extrinsic matrix ([**R**|**t**]), which can be represented as follows:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \overbrace{\begin{bmatrix} f_x & s & c_x & 0 \\ 0 & f_y & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}^{\mathbf{K}} \overbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix}}^{[\mathbf{R}|\boldsymbol{t}]} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}. \tag{3}$$

Direct linear transformation (DLT) is a mathematical technique commonly used to estimate the parameters of the projection matrix. The DLT method requires a minimum of six pairs of known 3D–2D correspondences to establish 12 equations and estimate all parameters of the projection matrix. Generally, the intrinsic parameters remain constant for a specific camera model, such as the Raspberry Pi Camera Module 2, and can be reused for all images captured by that camera. However, the extrinsic parameters change whenever the camera's location is altered. Consequently, for each setup deployment, recalculation of the extrinsic parameters is necessary to reconstruct the projection matrix. To simplify this process, the PnP method was replaced with DLT. It can reduce the required number of 3D–2D correspondence pairs to three by reusing the intrinsic parameters.

Additionally, ArUco markers were incorporated to represent pairs of known 3D–2D correspondences. For this purpose, the pixel coordinates of ArUco markers were determined using the OpenCV ArUco marker detection module on the 2D image plane, and the corresponding 3D real-world coordinates were measured by the total station. With these 3D–2D point correspondences, the spatial position and orientation of the camera can be estimated for each setup deployment. After retrieving all the necessary parameters, a full-perspective camera model can be generated. Using this model, the 3D point cloud is projected onto the 2D image plane. The projected (2D) point cloud represents the 3D real-world coordinates of the nearest 2D pixel correspondence on the image plane.

### 4.4 Machine learning for image measurements

Using the projection matrix, the 3D point cloud is projected on the 2D image plane (see Fig. 4). The projected (2D) point cloud is intersected with the waterline pixels, the output of the DL-based model (Module 1), to find the nearest point cloud coordinate. To achieve this objective, we utilize the $K$-nearest-neighbors (KNN) algorithm. Notably, the indices of the selected points remain consistent for both the 3D point cloud and the projected (2D) correspondences. As a result, by utilizing the indices of the chosen projected (2D) points, the corresponding real-world 3D coordinates can be retrieved.
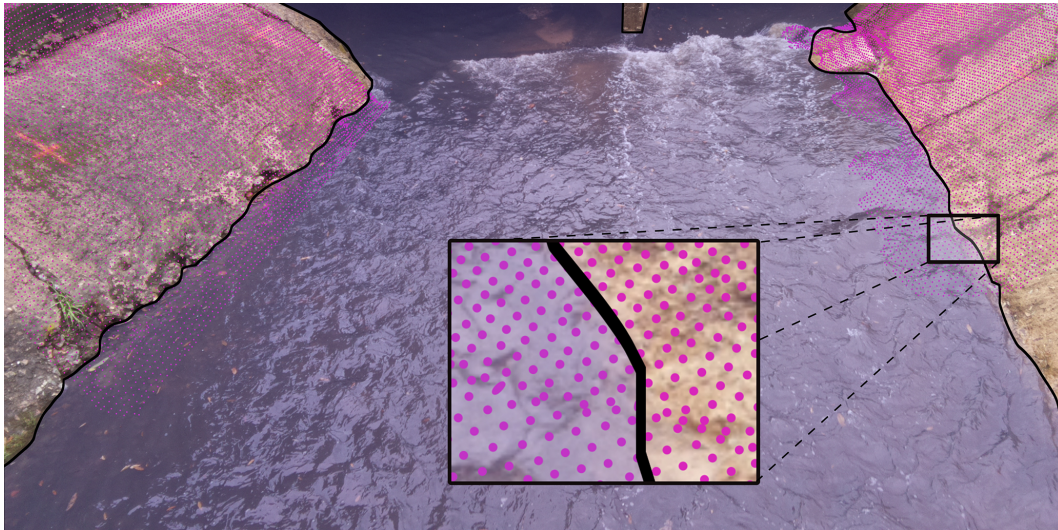
### 4.5 Performance metrics

The performance of the proposed framework is evaluated based on four different metrics including coefficient of determination ($R^2$), Nash–Sutcliffe efficiency (NSE), root mean square error (RMSE), and percent bias (PBIAS). $R^2$ is a widely used metric that quantifies how much of the observed dispersion can be explained in a linear relationship by the prediction.

$$r^2 = \left( \frac{\sum\limits_{i=1}^{n} \left( O_i - \overline{O} \right) \left( P_i - \overline{P} \right)}{\sqrt{\sum\limits_{i=1}^{n} \left( O_i - \overline{O} \right)^2 \cdot \sum\limits_{i=1}^{n} \left( P_i - \overline{P} \right)^2}} \right)^2 \tag{4}$$

However, if the model systematically over- or underestimates the results, $R^2$ will still be close to 1.0 as it only takes dispersion into account (Krause et al., 2005). NSE, another commonly used metric in hydrology, presents the model performance with an interpretable scale and is used to differentiate between "good" and "bad" models (Knoben et al., 2019).

$$\text{NSE} = 1 - \frac{\sum\limits_{i=1}^{n} \left( O_i - P_i \right)^2}{\sum\limits_{i=1}^{n} \left( O_i - \overline{O} \right)^2} \tag{5}$$

**Figure 4.** KNN is used to find the nearest projected (2D) point cloud (magenta dots) to the waterline (black line) on the image plane.

RMSE represents the square root of the average of squares of the errors, the differences between predicted values and observed values.

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(O_i - P_i)^2} \qquad (6)$$

The PBIAS of estimated water level, compared against the ultrasonic sensor data, was also used to show where the two estimates are close to each other and where they significantly diverge (Lin et al., 2020).

$$\text{PBIAS} = \frac{100}{n}\sum_{i=1}^{n}\frac{(O_i - P_i)}{\sum_{i=1}^{n}O_i}, \qquad (7)$$

where $n$ is the number of data points and $O$ and $P$ are observed and predicted values, respectively.

## 5    Results and discussion

The results of this study are presented in two sections. First, the performance of DL-based models is discussed. Then, in the second section, the performance of the proposed framework is evaluated for five different deployments.
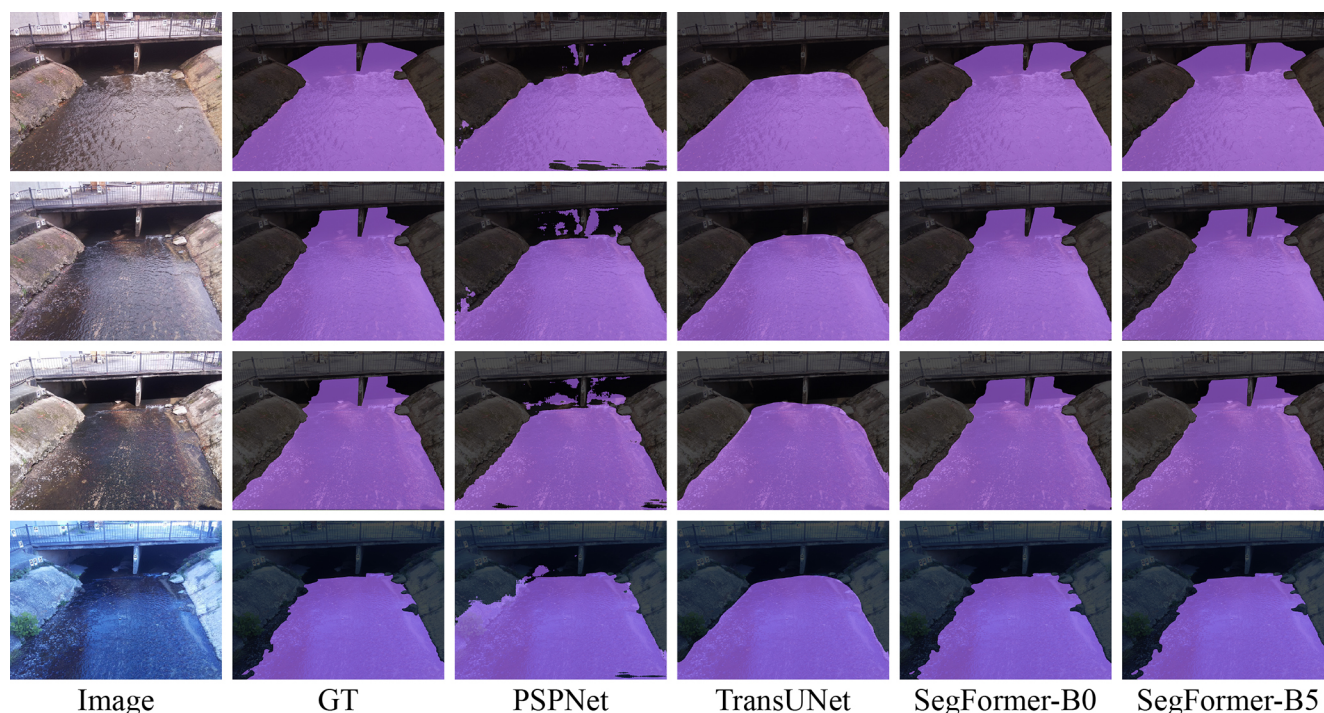
### 5.1    DL-based models results

The performance of DL-based models for the task of semantic segmentation is evaluated and compared in this section. Since the proposed dataset includes just two classes, "river" and "non-river", non-river was omitted from the evaluation process, and the performance of models is only reported for the river class of the test set. The class-wise intersection over

union (IoU) and the per-pixel accuracy (ACC) were considered the main evaluation metrics in this study. According to Table 2, both variants of SegFormer – SegFormer-B0 and SegFormer-B5 – outperform other semantic segmentation networks on the test set. Considering the models' configurations detailed in Table 1, SegFormer-B0 can be considered the most efficient DL-based network, as it is comprised of only 3.7 M trainable parameters and occupies just 2217 MB of GPU ram during training. In Fig. 5, four different visual representations of the models' performance on the validation set of the proposed dataset are presented. Since the water level is estimated by intersecting the waterline on riverbanks with the projected (2D) point cloud, precise delineation of the waterline is of utmost importance to achieve better results in the following steps. This means that estimating the correct location of the waterline on creek banks in each time-lapse image plays a more significant role than performance metrics in this study. Taking the quality of waterline detection into account and based on the visual representations shown in Fig. 5, SegFormers' variants still outperform DL-based approaches. In this regard, a comparison of PSP-Net and TransUNet showed that PSPNet can delineate the waterline more clearly, while the segmented area is more integrated for TransUNet outputs.

CNNs are typically limited by the nature of their convolution operations, leading to architecture-specific issues such as locality (Geirhos et al., 2018a). Consequently, CNN-based models may achieve high accuracy on training data, but their performance can decrease considerably on unseen data. Additionally, compared to transformer-based networks, they perform poorly at detecting semantics that require combining long- and short-range dependencies. Transformers can relax the biases of DL-based models induced by convolutional operations, achieving higher accuracy in localiza-

**Figure 5.** Visual representations of different DL-based image segmentation approaches on the validation dataset.

**Table 2.** The performance metrics of different DL-based approaches.

| Model names | IoU (River) | ACC (River) |
|---|---|---|
| PSPNet | 94.88 % | 95.84 % |
| TransUNet | 93.54 % | 96.89 % |
| SegFormer-B0 | 99.38 % | 99.77 % |
| SegFormer-B5 | 99.55 % | 99.81 % |

tion of target semantics and pixel-level classification with lower fluctuations in varied situations through the leverage of both local and global cues (Naseer et al., 2021). Yet, various transformer-based networks may perform differently depending on the targeted task and the network's architecture. TransUNet adopts transformers as part of its backbone; however, transformers generate single-scale low-resolution features as output (Xie et al., 2021), which may limit the accuracy when multi-scale objects or single objects with multi-scale features are segmented. The problem of producing single-scale features in standard transformers is addressed in SegFormer variants through the use of a novel hierarchical transformer encoder (Xie et al., 2021). This approach has resulted in human-level accuracy being achieved by SegFormer-B0 and SegFormer-B5 in the delineation of the waterline, as shown in Fig. 5. The predicted masks are in satisfactory agreement with the manually annotated images.

## 5.2 Water level estimation

This section reports the framework performance based on several deployments in the field. The performance results are separately shown for the left and right banks and compared with ultrasonic sensor data as the ground truth. The ultrasonic sensor was evaluated previously and documented an average distance error of 6.9 mm (Smith et al., 2022). The setup was deployed on several rainy days. The results of each deployment are reported in Table 3.

In addition to Table 3, the results of each deployment are visually demonstrated in Fig. 6. The scatterplots show the relationships between the ground truth data (measured by the ultrasonic sensor) and the banks of the river. The scatterplots visually present whether the camera readings overestimate or underestimate the ground truth data. Moreover, the time-series plot of water level is shown for each deployment separately. A hydrograph, showing changes in the water level of a stream over time, can be a useful tool for demonstrating whether camera readings can satisfactorily capture the response of a catchment area to rainfall. The proposed framework can be evaluated in terms of its ability to accurately track and identify important characteristics of a flood wave, such as the rising limb, peak, and recession limb.

The first deployment was done on 17 August 2022 (see Fig. 6a). The initial water level of the base flow and parts of the rising limb were not captured in this deployment. Table 3 shows that the performance results of the right-bank camera readings are better than those of the left bank. $R^2$

**Figure 6.** Scatterplot and time series plot for estimated water level by the proposed framework and measured by the ultrasonic sensor for setup deployment on **(a)** 17 August, **(b)** 19 August, and **(c)** 25 August 2022.
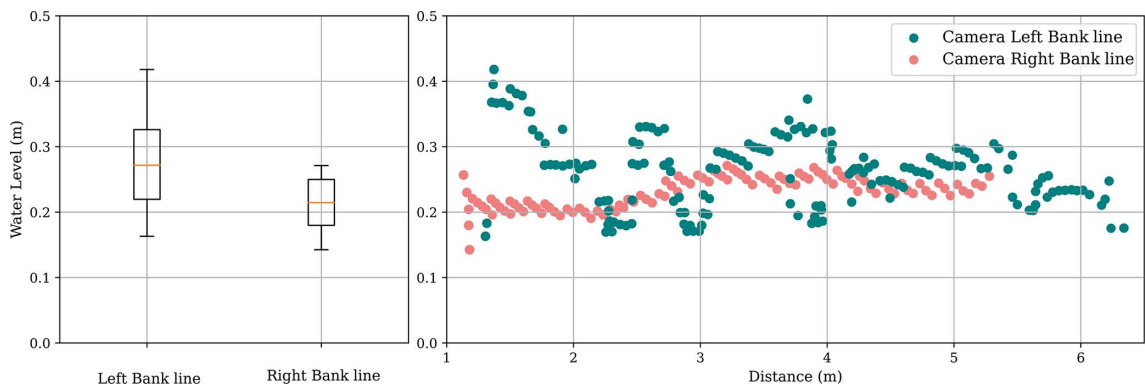
for both banks was about 0.80 showing a strongly related correlation between the water level estimated by the framework and ground truth data. Figure 6a shows how the left and right-bank camera readings perform during the rising limb; the right-bank camera readings still underestimated the water level during this time frame, and during the recession limb, the left-bank camera readings overestimated the water level. However, the hydrograph plot shows that both left and right-

bank camera readings were able to capture the peak water level.

The second deployment was done on 19 August 2022. In this deployment, all segments of the hydrograph were captured. According to Table 3, the performance of the right-bank camera readings was better than the left-bank one; more than 0.95 was reported for $R^2$ and the NSE of the right bank line. During the rising limb and crest segment, Fig. 6b shows

**Table 3.** The performance metrics of the framework for 5 different days of setup deployment.

| Deployment date | Position | Metrics | | | |
|---|---|---|---|---|---|
| | | $R^2$ | NSE | RMSE | PBIAS |
| 17 Aug 2022 | Left bank | 0.8019 | 0.5258 | 0.0409 | 10.6401 |
| | Right bank | 0.7932 | 0.7541 | 0.0294 | −0.4848 |
| 19 Aug 2022 | Left bank | 0.7701 | 0.5713 | 0.0647 | 16.1015 |
| | Right bank | 0.9678 | 0.9588 | 0.0201 | −3.4752 |
| 25 Aug 2022 | Left bank | 0.7690 | 0.5700 | 0.0435 | −7.7091 |
| | Right bank | 0.8922 | 0.8711 | 0.0238 | −1.7738 |
| 10 Nov 2022 | Left bank | 0.9461 | 0.8129 | 0.0511 | −13.1183 |
| | Right bank | 0.9857 | 0.9790 | 0.0171 | −1.5210 |
| 11 Nov 2022 | Left bank | 0.9588 | 0.8881 | 0.0397 | −10.3656 |
| | Right bank | 0.9855 | 0.9829 | 0.0155 | −1.7987 |



**Figure 7.** Water level fluctuation along both left and right banks for the flow regime for an image captured at 13:29 LT on 19 August 2022.

that both banks estimated a water level similar to ground truth. During the recession limb, the right-bank water level estimation remained coincident with ground truth, while the left bank overestimated the water level. The third deployment was on 25 August 2022. This time, the water level of the recession limb and the following base flow were captured (see Fig. 6c). The right-bank camera readings with $R^2$ of 0.89 performed better than the left bank. This time, left-bank camera readings underestimated the water level over the recession limb, but during the following base flow, the water level was estimated correctly by cameras on both banks.

The results indicate that the right-bank camera readings performed better than the left bank. Further investigation of the field conditions revealed that stream erosion had a more significant impact on the concrete surface of the left bank, resulting in patches and holes that were not scanned by the iPhone lidar. As a result, the KNN algorithm used to find the nearest (2D) point cloud coordinates to the waterline could not accurately represent the corresponding real-world coordinates of these locations. Figure 7 shows a box plot and scatterplot of the estimated water level for a time-lapse im-

age captured at 13:29 LT on 19 August 2022. The patches and holes on the left-bank surface caused instability in water level estimation for the region of interest. The box plot of the left bank (Cam-L-BL) was taller than that of the right bank (Cam-R-BL), indicating that the estimated water level was spread over larger values on the left bank due to the presence of these irregularities.

After analyzing the initial results, the deployable setups were modified to enhance the quality of data collection. The programming code of the Arduino device, Aava, was modified to measure five different records for water level each time it is triggered by the camera device, Beena, and to transmit the average distance to the Raspberry Pi device. This modification decreased the number of noise spikes in the measured data and allowed a better comparison between camera readings and ground truth data. The case of the camera device, Beena, was redesigned to protect the single board against rain without requiring an umbrella, which makes the camera setup unstable in stormy weather and causes a decrease in the precision of measurements. Moreover, an opening is incorporated into the redesigned case to connect an ex-

**Figure 8.** Scatterplot and time series plot for estimated water level by the proposed framework and measured by the ultrasonic sensor for setup deployment on **(a)** 10 November and **(b)** 11 November 2022.

ternal power bank to enhance the run time. Finally, the viewpoint of the camera was subtly shifted to the right to adjust the share of the riverbanks on the camera's field of view.

The results of the deployments on 10 and 11 November 2022 demonstrate that modifications to the setup have significantly improved the results of the left bank (as shown in Table 3). NSE improved from approximately 0.55 for the first three setup deployments to over 0.80 for the modified deployments. Figure 8 shows the setup performances during all segments of the flood wave. The peaks were captured by the right bank line on both deployment dates, and there was no effect of noisy spikes on either camera readings or ground truth data. However, the right-bank images still underestimated the water level during the rainstorms.

## 6   Conclusion

This study introduced Eye of Horus, a vision-based framework for hydrologic monitoring and measuring of real-time water-related parameters, e.g., water level, from surveillance images captured during flood events. Time-lapse images and real water level correspondences were collected by a Rasp-

berry Pi camera and an Arduino HC-SR05 ultrasonic sensor, respectively. Moreover, computer vision and deep learning techniques were used for semantic segmentation of the water surface within the captured images and for reprojecting the 3D point cloud constructed with an iPhone lidar scanner, on the (2D) image plane. Eventually, the $K$-nearest neighbor algorithm was used to intersect the projected (2D) point cloud with the waterline pixels extracted from the output of the deep learning model to find the real-world 3D coordinates.

A vision-based framework offers a new alternative to current hydrologic data collection and real-time monitoring systems. Hydrological models require geometric information for estimating discharge routing parameters, stage, and flood inundation maps. However, determining bankfull characteristics is a challenge due to natural or anthropogenic downcutting of streams. Using visual sensing, stream depth, water velocity, and instantaneous streamflow at bankfull stage can be reliably measured.

# References

Alsdorf, D. E., Rodríguez, E., and Lettenmaier, D. P.: Measuring surface water from space, Rev. Geophys., 45, RG2002, https://doi.org/10.1029/2006RG000197, 2007.

Badrinarayanan, V., Handa, A., and Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for robust semantic pixel-wise labelling, arXiv [preprint], https://doi.org/10.48550/arXiv.1505.07293, 2015.

Bradski, G.: The OpenCV Library, Dr. Dobb's Journal of Software Tools, https://opencv.org/ (last access: 4 November 2023), 2000.

Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A. L., and Zhou, Y.: Transunet: Transformers make strong encoders for medical image segmentation, arXiv [preprint], https://doi.org/10.48550/arXiv.2102.04306, 2021.

Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., and Yuille, A. L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, IEEE T. Pattern Anal. Mach. Intel., 40, 834–848, 2017.

De Cesarei, A., Cavicchi, S., Cristadoro, G., and Lippi, M.: Do humans and deep convolutional neural networks use visual information similarly for the categorization of natural scenes?, Cognit. Sci., 45, e13009, https://doi.org/10.1111/cogs.13009, 2021.

Dodge, S. and Karam, L.: A study and comparison of human and deep learning recognition performance under visual distortions, in: IEEE Int. Conf. Comput. Communication and Networks, Vancouver, BC, Canada, 1–7, https://doi.org/10.1109/ICCCN.2017.8038465, 2017.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N.: An image is worth $16 \times 16$ words: Transformers for image recognition at scale, arXiv [preprint], https://doi.org/10.48550/arXiv.2010.11929, 2020.

Elias, M., Eltner, A., Liebold, F., and Maas, H.-G.: Assessing the influence of temperature changes on the geometric stability of smartphone-and raspberry pi cameras, Sensors, 20, 643, https://doi.org/10.3390/s20030643, 2020.

Eltner, A. and Schneider, C.: Analysis of different methods for 3d reconstruction of natural surfaces from parallel-axes uav images, Photogram. Rec., 30, 279–299, 2015.

Eltner, A., Kaiser, A., Castillo, C., Rock, G., Neugirg, F., and Abellán, A.: Image-based surface reconstruction in geomorphometry – merits, limits and developments, Earth Surf. Dynam., 4, 359–389, https://doi.org/10.5194/esurf-4-359-2016, 2016.

Eltner, A., Elias, M., Sardemann, H., and Spieler, D.: Automatic image-based water stage measurement for long-term observations in ungauged catchments, Water Resour. Res., 54, 10362–10371, https://doi.org/10.1029/2018WR023913, 2018.

Eltner, A., Bressan, P. O., Akiyama, T., Gonçalves, W. N., and Marcato Junior, J.: Using deep learning for automatic water stage measurements, Water Resour. Res., 57, e2020WR027608, https://doi.org/10.1029/2020WR027608, 2021.

Erfani, S. M. H.: smhassanerfani/horus: Pre-release version (v1.0.0-alpha), Zenodo [data set], https://doi.org/10.5281/zenodo.10071662, 2023.

Erfani, S. M. H., Wu, Z., Wu, X., Wang, S., and Goharian, E.: Atlantis: A benchmark for semantic segmentation of waterbody images, Environ. Model. Softw., 149, 105333, https://doi.org/10.1016/j.envsoft.2022.105333, 2022.

Forsyth, A. A. and Ponce, J.: Computer vision: a modern approach, Prentice hall professional technical reference, ISBN 0130851981, 2002.

Froideval, L., Pedoja, K., Garestier, F., Moulon, P., Conessa, C., Pellerin Le Bas, X., Traoré, K., and Benoit, L.: A low-cost open-source workflow to generate georeferenced 3d sfm photogrammetric models of rocky outcrops, Photogram. Rec., 34, 365–384, 2019.

Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., and Lu, H.: Dual attention network for scene segmentation, in: IEEE Conf. Comput. Vis. Pattern Recog., Long Beach, CA, USA, 3141–3149, https://doi.org/10.1109/CVPR.2019.00326, 2019.

Gebrehiwot, A., Hashemi-Beni, L., Thompson, G., Kordjamshidi, P., and Langan, T. E.: Deep convolutional neural network for flood extent mapping using unmanned aerial vehicles data, Sensors, 19, 1486, https://doi.org/10.3390/s19071486, 2019.

Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W.: Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness,

arXiv [preprint], https://doi.org/10.48550/arXiv.1811.12231, 2018a.

Geirhos, R., Temme, C. R. M., Rauber, J., H Schütt, H., Bethge, M., and Wichmann, F. A.: Generalisation in humans and deep neural networks, Adv. Neural Inform. Process. Syst., 31, 7538–7550, ISBN 9781510884472, 2018b.

Geirhos, R., Meding, K., and Wichmann, F. A.: Beyond accuracy: quantifying trial-by-trial behaviour of cnns and humans by measuring error consistency, Adv. Neural Inform. Process. Syst., 33, 13890–13902, 2020.

Gilmore, T. E., Birgand, F., and Chapman, K. W.: Source and magnitude of error in an inexpensive image-based water level measurement system, J. Hydrol., 496, 178–186, 2013.

Gollob, C., Ritter, T., Kraßnitzer, R., Tockner, A., and Nothdurft, A.: Measurement of forest inventory parameters with Apple iPad pro and integrated LiDAR technology, Remote Sens., 13, 3129, https://doi.org/10.3390/rs13163129, 2021.

Goodchild, M. F.: Citizens as sensors: the world of volunteered geography, Geo J., 69, 211–221, 2007.

He, K., Zhang, X., Ren, S., and Sun, J.: Deep residual learning for image recognition, in: IEEE Conf. Comput. Vis. Pattern Recog., Las Vegas, NV, USA, 770–778, https://doi.org/10.1109/CVPR.2016.90, 2016.

Howe, J.: Crowdsourcing: How the power of the crowd is driving the future of business, Random House, https://doi.org/10.2146/ajhp100029, 2008.

Huang, Z., Wang, X., Huang, L., Huang, C., Wei, Y., and Liu, W.: Ccnet: Criss-cross attention for semantic segmentation, in: Int. Conf. Comput. Vis., Seoul, South Korea, 603–612, https://doi.org/10.1109/ICCV.2019.00069, 2019.

Kim, J., Han, Y., and Hahn, H.: Embedded implementation of image-based water-level measurement system, IET Comput. Vis., 5, 125–133, 2011.

King, T. V., Neilson, B. T., and Rasmussen, M. T.: Estimating discharge in low-order rivers with high-resolution aerial imagery, Water Resour. Res., 54, 863–878, 2018.

Knoben, W. J. M., Freer, J. E., and Woods, R. A.: Technical note: Inherent benchmark or not? Comparing Nash–Sutcliffe and Kling–Gupta efficiency scores, Hydrol. Earth Syst. Sci., 23, 4323–4331, https://doi.org/10.5194/hess-23-4323-2019, 2019.

Krause, P., Boyle, D. P., and Bäse, F.: Comparison of different efficiency criteria for hydrological model assessment, Adv. Geosci., 5, 89–97, https://doi.org/10.5194/adgeo-5-89-2005, 2005.

LAAN LABS: 3D Scanner App – LiDAR Scanner for iPad Pro & iPhone Pro, https://3dscannerapp.com/ (last access: 16 September 2022), 2022.

Li, X., Zhong, Z., Wu, J., Yang, Y., Lin, Z., and Liu, H.: Expectation-maximization attention networks for semantic segmentation, in: Int. Conf. Comput. Vis., Seoul, South Korea, 9166–9175, https://doi.org/10.1109/ICCV.2019.00926, 2019.

Li, Z., Wang, C., Emrich, C. T., and Guo, D.: A novel approach to leveraging social media for rapid flood mapping: a case study of the 2015 south carolina floods, Cartogr. Geogr. Inform. Sci., 45, 97–110, 2018.

Lin, G., Milan, A., Shen, C., and Reid, I.: Refinenet: Multi-path refinement networks for high-resolution semantic segmentation, in: IEEE Conf. Comput. Vis. Pattern Recog., Honolulu, HI, USA, 5168–5177, https://doi.org/10.1109/CVPR.2017.549, 2017.

Lin, P., Pan, M., Allen, G. H., de Frasson, R. P., Zeng, Z., Yamazaki, D., and Wood, E. F.: Global estimates of reach-level bankfull river width leveraging big data geospatial analysis, Geophys. Res. Lett., 47, e2019GL086405, https://doi.org/10.1029/2019GL086405, 2020.

Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows, in: Int. Conf. Comput. Vis., Montreal, QC, Canada, 9992–10002, https://doi.org/10.1109/ICCV48922.2021.00986, 2021.

Lo, S.-W., Wu, J.-H., Lin, F.-P., and Hsu, C.-H.: Visual sensing for urban flood monitoring, Sensors, 15, 20006–20029, 2015.

Long, J., Shelhamer, E., and Darrell, T.: Fully convolutional networks for semantic segmentation, in: IEEE Conf. Comput. Vis. Pattern Recog., Boston, MA, USA, 3431–3440, https://doi.org/10.1109/CVPR.2015.7298965, 2015.

Loshchilov, I. and Hutter, F.: Decoupled weight decay regularization, arXiv [preprint], https://doi.org/10.48550/arXiv.1711.05101, 2017.

Luetzenburg, G., Kroon, A., and Bjørk, A. A.: Evaluation of the apple iphone 12 pro lidar for an application in geosciences, Sci. Rep., 11, 1–9, 2021.

Marchand, E., Uchiyama, H., and Spindler, F.: Pose estimation for augmented reality: a hands-on survey, IEEE T. Pattern Anal. Mach. Intel., 22, 2633–2651, 2015.

Minaee, S., Boykov, Y. Y., Porikli, F., Plaza, A. J., Kehtarnavaz, N., and Terzopoulos, D.: Image segmentation using deep learning: A survey, IEEE T. Pattern Anal. Mach. Intel., 44, 3523–3542, https://doi.org/10.1109/TPAMI.2021.3059968, 2022.

Mokroš, M., Mikita, T., Singh, A., Tomaštík, J., Chudá, J., Wężyk, P., Kuželka, K., Surovỳ, P., Klimánek, M., Zięba-Kulawik, K., Bobrowski, R., and Liang, X.: Novel low-cost mobile mapping systems for forest inventories as terrestrial laser scanning alternatives, Int. J. Appl. Earth Obs. Geoinf., 104, 102512, https://doi.org/10.1016/j.jag.2021.102512, 2021.

Morsy, M. M., Goodall, J. L., Shatnawi, F. M., and Meadows, M. E.: Distributed stormwater controls for flood mitigation within urbanized watersheds: case study of rocky branch watershed in columbia, south carolina, J. Hydrol. Eng., 21, 05016025, https://doi.org/10.1061/(ASCE)HE.1943-5584.0001430, 2016.

Moy de Vitry, M., Kramer, S., Wegner, J. D., and Leitão, J. P.: Scalable flood level trend monitoring with surveillance cameras using a deep convolutional neural network, Hydrol. Earth Syst. Sci., 23, 4621–4634, https://doi.org/10.5194/hess-23-4621-2019, 2019.

Naseer, M. M., Ranasinghe, K., Khan, S. H., Hayat, M., Shahbaz Khan, F., and Yang, M.-H.: Intriguing properties of vision transformers, Adv. Neural Inform. Process. Syst., 34, 23296–23308, 2021.

Noh, H., Hong, S., and Han, B.: Learning deconvolution network for semantic segmentation, in: Int. Conf. Comput. Vis., Santiago, Chile, 1520–1528, https://doi.org/10.1109/ICCV.2015.178, 2015.

Pally, R. J. and Samadi, S.: Application of image processing and convolutional neural networks for flood image classification and semantic segmentation, Environ. Model. Softw., 148, 105285, https://doi.org/10.1016/j.envsoft.2021.105285, 2022.

Panteras, G. and Cervone, G.: Enhancing the temporal resolution of satellite-based flood extent generation using crowdsourced data

for disaster monitoring, Int. J. Remote Sens., 39, 1459–1474, 2018.

Schnebele, E., Cervone, G., and Waters, N.: Road assessment after flood events using non-authoritative data, Nat. Hazards Earth Syst. Sci., 14, 1007–1015, https://doi.org/10.5194/nhess-14-1007-2014, 2014.

Shamsabadi, E. A., Xu, C., and Dias-da Costa, D.: Robust crack detection in masonry structures with transformers, Measurement, 200, 111590, https://doi.org/10.1016/j.measurement.2022.111590, 2022.

Smith, C., Satme, J., Martin, J., Downey, A. R. J., Vitzilaios, N., and Imran, J.: UAV rapidly-deployable stage sensor with electro-permanent magnet docking mechanism for flood monitoring in undersampled watersheds, HardwareX, 12, e00325, https://doi.org/10.1016/j.ohx.2022.e00325, 2022.

Tavani, S., Billi, A., Corradetti, A., Mercuri, M., Bosman, A., Cuffaro, M., Seers, T., and Carminati, E.: Smartphone assisted fieldwork: Towards the digital transition of geoscience fieldwork using lidar-equipped iphones, Earth-Sci. Rev., 227, 103969, https://doi.org/10.1016/j.earscirev.2022.103969, 2022.

Tsubaki, R., Fujita, I., and Tsutsumi, S.: Measurement of the flood discharge of a small-sized river using an existing digital video recording system, J. Hydro-Environ. Res., 5, 313–321, 2011.

Turnipseed, D. P. and Sauer, V. B.: Discharge measurements at gaging stations, Technical report, US Geological Survey, https://doi.org/10.3133/tm3A8, 2010.

Vandaele, R., Dance, S. L., and Ojha, V.: Deep learning for automated river-level monitoring through river-camera images: an approach based on water segmentation and transfer learning, Hydrol. Earth Syst. Sci., 25, 4435–4453, https://doi.org/10.5194/hess-25-4435-2021, 2021.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I.: Attention is all you need, Adv. Neural Inform. Process. Syst., 30, 5998–6008, ISBN 9781510860964, 2017.

Vogt, M., Rips, A., and Emmelmann, C.: Comparison of ipad pro's lidar and truedepth capabilities with an industrial 3d scanning solution, Technologies, 9, 25, https://doi.org/10.3390/technologies9020025, 2021.

Westoby, M. J., Brasington, J., Glasser, N. F., Hambrey, M. J., and Reynolds, J. M.: 'structure-from-motion' photogrammetry: A low-cost, effective tool for geoscience applications, Geomorphology, 179, 300–314, 2012.

Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., and Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers, Adv. Neural Inform. Process. Syst., 34, 12077–12090, 2021.

Yuan, Y. and Wang, J.: Ocnet: Object context network for scene parsing, arXiv [preprint], https://doi.org/10.48550/arXiv.1809.00916, 2018.

Yuan, Y., Chen, X., and Wang, J.: Object-contextual representations for semantic segmentation, in: Eur. Conf. Comput. Vis., Springer, 173–190, https://doi.org/10.1007/978-3-030-58539-6_11, 2020.

Zhang, Z., Zhou, Y., Liu, H., and Gao, H.: In-situ water level measurement using nir-imaging video camera, Flow Meas. Instrum., 67, 95–106, 2019.

Zhao, H., Shi, J., Qi, X., Wang, X., and Jia, J.: Pyramid scene parsing network, in: Proceedings of the IEEE conference on computer vision and pattern recognition, Honolulu, HI, USA, 6230–6239, https://doi.org/10.1109/CVPR.2017.660, 2017.

Zheng, Y., Huang, J., Chen, T., Ou, Y., and Zhou, W.: Processing global and local features in convolutional neural network (cnn) and primate visual systems, Mobile Multimed./Image Process. Secur. Appl., 10668, 44–51, 2018.

Zhu, Z., Xu, M., Bai, S., Huang, T., and Bai, X.: Asymmetric non-local neural networks for semantic segmentation. in: Int. Conf. Comput. Vis., Seoul, South Korea, 593–602, https://doi.org/10.1109/ICCV.2019.00068, 2019.