
ATLANTIS: A BENCHMARK FOR SEMANTIC SEGMENTATION OF WATERBODY IMAGES

Mohammad H. Erfani

Department of Civil and Environmental Engineering
University of South Carolina
Columbia, SC
serfani@email.sc.edu

Zhenyao Wu, Xinyi Wu

Department of Computer Science and Engineering
University of South Carolina
Columbia, SC
{zhenyao, xinyiw}@email.sc.edu

Song Wang

Department of Computer Science and Engineering
University of South Carolina
Columbia, SC
songwang@cec.sc.edu

Erfan Goharian

Department of Civil and Environmental Engineering
University of South Carolina
Columbia, SC
goharian@cec.sc.edu

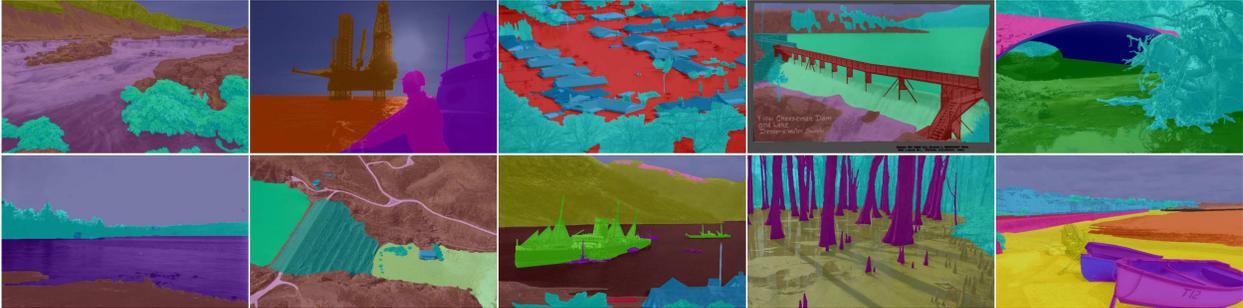


Figure 1: ATLANTIS: Artificial and Natural waterbodies dataset.

ABSTRACT

Vision-based semantic segmentation of waterbodies and nearby related objects provides important information for managing water resources and handling flooding emergency. However, the lack of large-scale labeled training and testing datasets for water-related categories prevents researchers from studying water-related issues in the computer vision field. To tackle this problem, we present ATLANTIS, a new benchmark for semantic segmentation of waterbodies and related objects. ATLANTIS consists of 5,195 images of waterbodies, as well as high quality pixel-level manual annotations of 56 classes of objects, including 17 classes of man-made objects, 18 classes of natural objects and 21 general classes. We analyze ATLANTIS in detail and evaluate several state-of-the-art semantic segmentation networks on our benchmark. In addition, a novel deep neural network, AQUANet, is developed for waterbody semantic segmentation by processing the aquatic and non-aquatic regions in two different paths. AQUANet also incorporates low-level feature modulation and cross-path modulation for enhancing feature representation. Experimental results show that the proposed AQUANet outperforms other state-of-the-art semantic segmentation networks on ATLANTIS. We claim that ATLANTIS is the largest waterbody image dataset for semantic segmentation providing a wide range of water and water-related classes and it will benefit researchers of both computer vision and water resources engineering.

1 Introduction

Every year, floods claim tens of billions of US dollars losses and thousands of lives globally [17]. Accurate detection, measurement, and track of the waterbodies can help both the public and decision-makers to take appropriate actions to minimize the risk and losses [19, 15]. With the popularity of smart phones and drones, various data at flooding sites can be collected rapidly and continuously to provide more useful and heterogeneous information source [12, 11, 5, 38], compared to the conventional gauge sensing and remote sensing [12, 27]. As a fundamental step to leverage such collected images for modeling and decision-making, we need first to conduct a refined semantic segmentation of included waterbodies and related objects in such scenes, which we focus on in this paper.

With the advancement of deep neural networks, semantic segmentation has achieved a great success in recent years on various kinds of images, such as natural images [26], street images [8, 46], and medical images [2, 22]. However, waterbody images pose many new unique challenges for semantic segmentation. In some forms, water preserves the intrinsic properties such as reflection, transparency, shapeless and colorless visual features; which in turn, brings difficulties to semantic segmentation of water and related objects. Moreover, in some other forms, these properties can be affected by illumination sources from surroundings, turbidity, and turbulence. Different-labeled waterbodies, such as river and canal, or lake and reservoir, often have similar visual characteristics that make task of semantic segmentation even harder. As shown in our later experiments, these unique challenges may significantly affect the performance of the existing semantic segmentation networks.

Meanwhile, deep-learning based approaches always require large-scale training data with necessary ground-truth annotations. Lack of such a public dataset for waterbody segmentation significantly impedes the research on this problem. The collection and annotation of such a dataset can be very laborious and time-consuming to cover a wide range of waterbodies and related objects. There is no specific repository providing relevant images. In addition, team members and annotators are required to have prior knowledge on water resources engineering to be capable of selecting and precisely annotating the images.

In this paper, we present a new benchmark, ATLANTIS (ArTificial And Natural waTer-bodIes dataSet). For the first time, this dataset has covered a wide range of natural and man-made (artificial) waterbodies such as sea, lake, river, canal, reservoir, and dam. ATLANTIS includes 5,195 pixel-wise annotated images split to 3,364 training, 535 validation and 1,296 testing images. As shown in the Table 1, in addition to 35 waterbody and water-related objects, ATLANTIS also covers 21 general labels. Moreover, we construct ATLANTIS Texture (ATeX) dataset, which consists of 12,503 patches for the water-bodies texture classification, sampled from 15 kinds of waterbodies in ATLANTIS.

Table 1: List of the ATLANTIS labels.

Artificial	breakwater; bridge; canal; culvert; dam; ditch; levee; lighthouse; pipeline; pier; offshore platform; reservoir; ship; spillway; swimming pool; water tower; water well.
Natural	cliff; cypress tree; fjord; flood; glaciers; hot spring; lake; mangrove; marsh; puddle; rapids; river; river delta; sea; shoreline; snow; waterfall; wetland.
General	road; sidewalk; building; wall; fence; pole; traffic sign; vegetation; terrain; sky; train; person; car; bus; truck; bicycle; parking meter; motorcycle; fire hydrant; boat; umbrella.

In order to tackle the inherent challenges in the segmentation of waterbodies, AQUANet is developed which takes an advantage of two different paths to process the aquatic and non-aquatic regions, separately. Each path includes low-level feature and cross-path modulation, to adjust features for better representation. The results show that the proposed AQUANet outperforms other ten state-of-the-art semantic-segmentation networks on ATLANTIS, and the ablation studies justify the effectiveness of the components of the proposed AQUANet.

2 Related Work

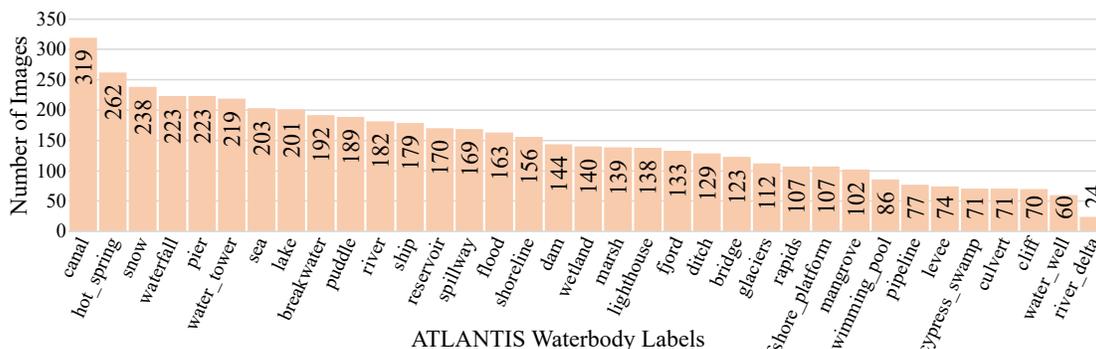
2.1 Semantic segmentation dataset

Large-scale annotated datasets, such as COCO [26], PASCAL Context [30], ADE20K [51], and more recently Mapillary Vistas Dataset [32] and BDD100K [46], make it possible for researchers to develop deep-learning based models for real-world applications. Considering the most related dataset to ATLANTIS, Gebrehiwot et al. [15] collected a small number of top-view waterbody dataset (100 images) using Unmanned Aerial Vehicles (UAVs) which contains only

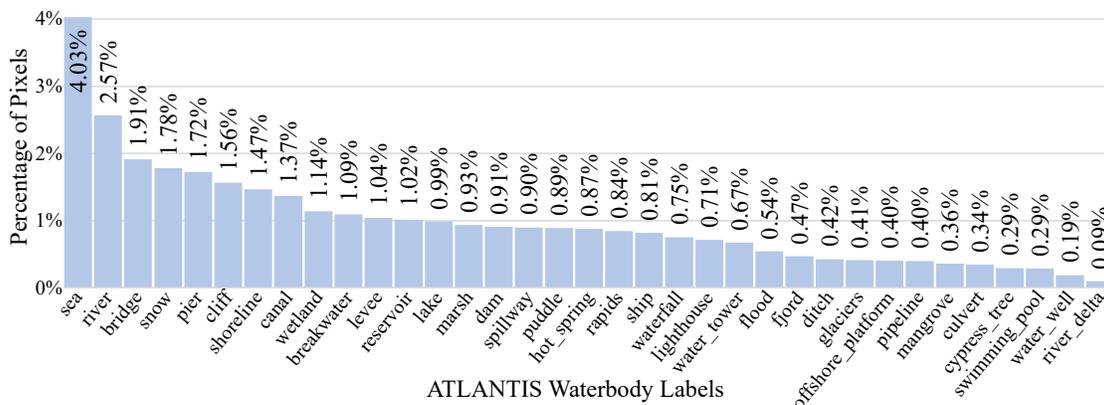
four categories (i.e., water, building, vegetation and road). In addition, Sazara et al. [37] introduced a larger dataset (253 images) which just focuses on flood region segmentation. More recently, Sarp et al. [36] provided a dataset that consists of 441 annotated roadway flood images. However, these datasets have limitations in either the number of annotated images, or the categories they covered, and none of those considers more complex classes of waterbodies such as sea, lake and waterfall. Therefore, we develop a new dataset, ATLANTIS, as the first large-scale annotated dataset to provide a wide-range of waterbodies and water-related objects.

2.2 Semantic segmentation network

All of existing semantic segmentation approaches share the same goal to classify each pixel of a given image but differ in the network design, including low-resolution representations learning [28, 6], high-resolution representations recovering [1, 33, 25], contextual aggregation schemes [47, 50, 48], feature fusion and refinement strategy [25, 20, 23, 52, 14]. Typically, method designs are dependent on their respective datasets and all the mentioned networks are developed by training on benchmark datasets such as Cityscapes [8], COCO [26] and VOC [13] where the inter-class boundary is clear even for the within-group categories (e.g., car and truck). As mentioned above, waterbody images pose new challenges to semantic segmentation. Previous works on waterbody segmentation mainly use satellite imagery [31, 24, 10]. In this work, we focus on natural waterbody images terrestrially captured by various cameras, and design AQUANet, a new two-path semantic segmentation network, by including an aquatic branch explicitly for waterbody classes.



(a)



(b)

Figure 2: (a) The frequency distribution of images for different waterbodies in ATLANTIS (b) The Percentage of pixels for waterbody labels.

3 ATLANTIS Dataset

The ATLANTIS dataset is designed and developed with the goal of capturing a wide-range of water-related objects, either those exist in natural environment or the infrastructure and man-made (artificial) water systems. In this dataset,

labels were first selected based on the most frequent objects, used in water-related studies or can be found in real-world scenes. Aside from the background objects, total of 56 labels, including 17 artificial, 18 natural waterbodies, and 21 general labels, are selected (Table 1). These general labels are considered for providing contextual information that most likely can be found in water-related scenes. After finalizing the selection of waterbody labels, a comprehensive investigation on each individual label was performed by annotators to make sure all the labels are vivid examples of those objects in real-world. Moreover, sometimes some of the water-related labels, e.g., levee, embankment, and floodbank, have been used interchangeably in water resources field; thus, those labels are either merged into a unique group or are removed from the dataset to prevent an individual object receives different labels.

In order to gather a corpus of images, we have used Flickr API to query and to collect “medium-sized” unique images for each label based on eight commonly used “Creative Commons”, “No Known Copyright Restrictions” and “United States Government Work” licenses. Downloaded images were then filtered by a two-stage hierarchical procedure. In the first stage, each annotator was assigned to review a specific list of labels and remove irrelevant images based on that specific list of labels. In the second stage, several meetings were held between the entire annotation team and the project coordinator to finalize the images which appropriately represent each of 56 labels.

This sieving procedure has been applied four times in order to meet the limit and to reach the current number of images. The percentage of image acceptance rate for the third and fourth phases are 14.41% and 5.06%, respectively. It means if we want to add 1000 more images to the dataset, we should process at least 20,000 images. Finally, images were annotated by annotators who have solid water resources engineering background as well as experience working with the CVAT [39], which is a free, open source, and web-based image/video annotation tool.

3.1 Dataset statistics

Figure 2 shows the frequency distribution of the number of images and the percentage of pixels for waterbody labels. Such a long-tailed distribution is common for semantic segmentation datasets [26, 51] even if the number of images that contain specific label are pre-controlled. Such frequency distribution for pixels would be inevitable for objects existing in real-world. Taking “water tower” as an example, despite having 219 images, the percentage of pixels are less than many other labels in the dataset. Figure 3 shows the positive but weak correlation between number of images for each label and the corresponding pixels. In total, only 4.89% of pixels are unlabeled, and 34.17% and 60.94% of pixels belong to waterbodies (natural and man-made) and general labels, respectively. As it is evident, the main proportion of pixels belongs to general labels. This clearly shows the importance of general labels for better scene understanding [3] and accurate object classification in semantic segmentation network.

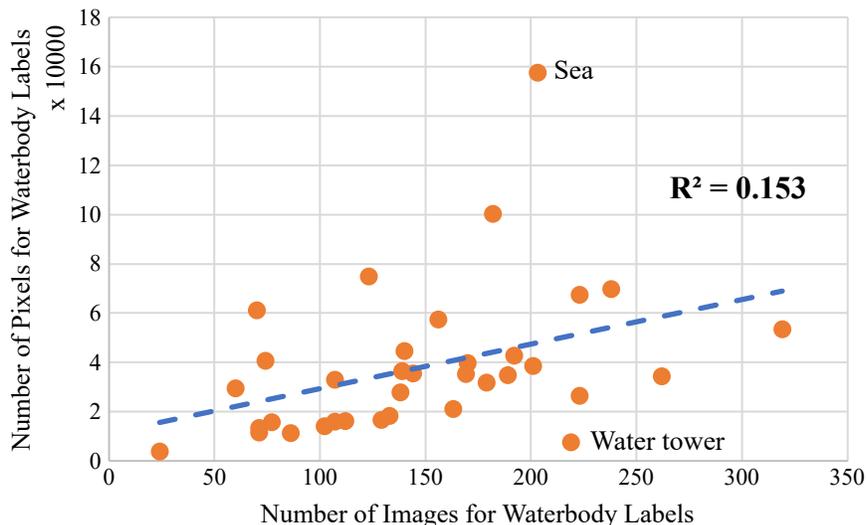


Figure 3: Correlation between number of images for each label and the corresponding pixels

3.1.1 Spatial analysis

Following ADE20K dataset [51], a spatial analysis, known as “mode of the object segmentation” has been done on the ground truth segmentation map for each label. Specifically, considering a waterbody label L with n images in ATLANTIS, we resize their corresponding n ground-truth segmentation maps to 512×512 pixels. We then count the

most frequent label at each pixel of the map and construct a “mode of segmentation” map for label L , as shown in Figure 4. This map demonstrates the spatial distributions of the most frequent co-occurred labels with respect to a given waterbody label. Based on this, we equipped our proposed network with cross-path feature modulation to cope with the difficulties associated with the recognition of waterbody labels having visual similarities.

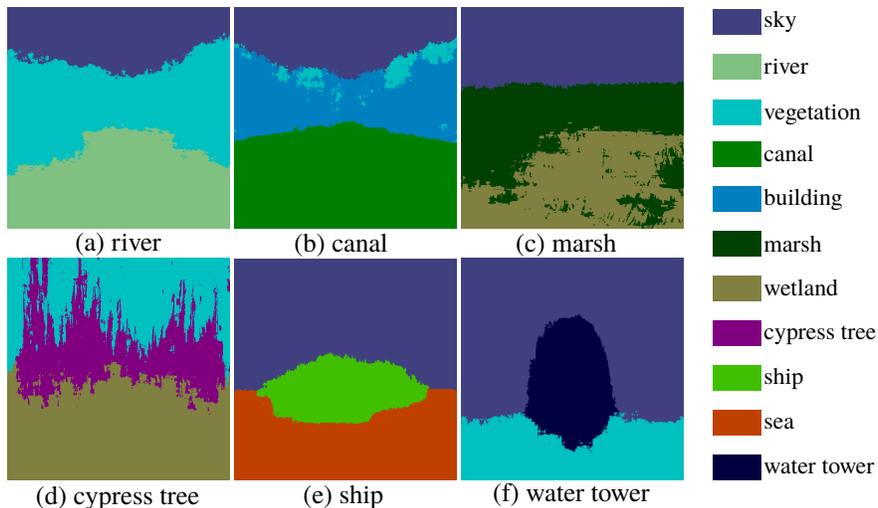


Figure 4: Spatial analysis of four different waterbody labels.

3.2 Image annotation

3.2.1 Annotation pipeline

ATLANTIS was annotated by six annotators having prior knowledge in the area of water resources. The goal in the annotation task is to balance speed and quality. Generally, time spent on a single image can range from 4 minutes to 25 minutes depending on the image complexity. In this project, each kind of waterbody was assigned to a specific annotator. Before annotation of a label, all the images of that label are scrutinized and discussed by a group of experts in water resources engineering. We can see that the annotation of complex flood scenes takes more time since such images are usually captured in urban areas and have more elaborated components as shown in Figure 5.

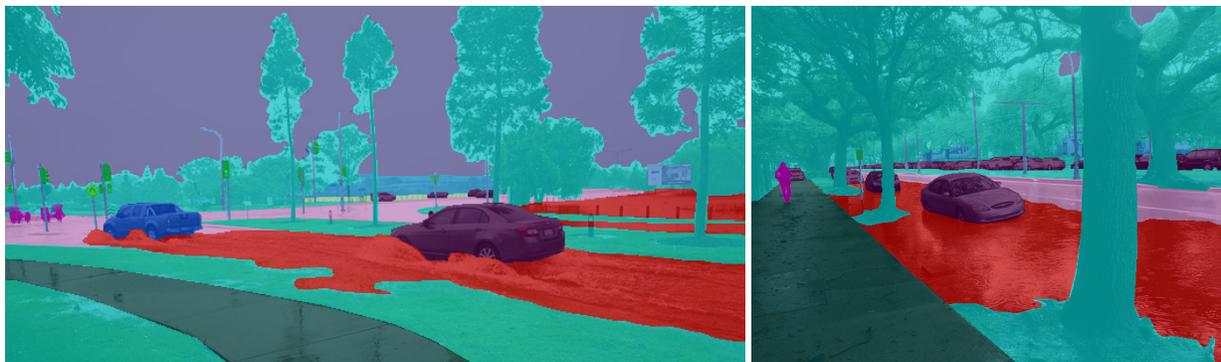


Figure 5: Two samples of complex flood scenes

3.2.2 Consistency analysis

While one image is annotated by one annotator for ATLANTIS, we perform additional consistency analysis across annotators and over time for an annotator. We choose 52 images from ATLANTIS, by including both images that are highly susceptible to wrong labelling and those contain objects prone to be either left unannotated or wrongly annotated. We ask three annotators to annotate them again and compare the results against the already approved ground truth in ATLANTIS. The accuracy and mIoU in terms of all 52 images (total) and the subsets of images that had been annotated

by themselves before (individual) are shown in Table 2. We can see that an annotator can process the images that he/she annotated before with much better consistency.

Table 2: Consistency analysis for three annotators.

		Annotator 1	Annotator 2	Annotator 3
Total	acc	84.00	79.93	79.00
	mIoU	62.29	59.33	57.34
Individual	acc	91.11	90.44	94.69
	mIoU	72.83	76.14	75.69

3.3 ATLANTIS Texture (ATeX)

Waterbodies usually bear texture appearance and it is an interesting problem to study whether different kinds of waterbodies may show subtle differences associated with the texture features. For this purpose, we construct a new waterbody texture dataset, ATeX, by cropping patches from ATLANTIS and take the corresponding annotated waterbody label as the label of the patch. We set patch size to 32×32 pixels and all pixels in a cropped patch must have the same waterbody label in the original image. We also ensure there is no partial overlap between any two patches. In total we collected 12,503 patches with 15 waterbody labels: Two waterbody labels “estuary” and “swamp” are added based on the nearby tree species – mangrove “estuary” and cypress for “swamp”, while four waterbody labels “canal”, “ditch”, “reservoir” and “fjord” are omitted because of high visual similarities with other labels. Sample images of ATeX dataset are shown in Figure 6. We split ATeX into 8,753 for training, 1,252 for validation and 2,498 for testing. In the later experiment, we train different models to evaluate their classification performance on ATeX images.

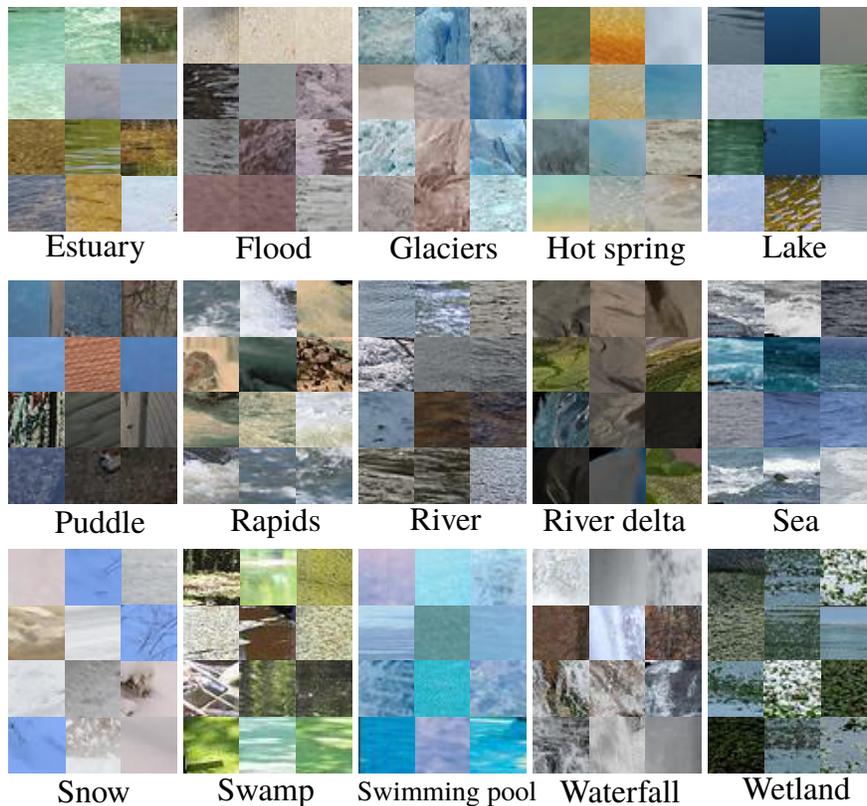


Figure 6: Samples of ATeX texture images.

4 AQUANet

Typically, existing semantic segmentation networks are designed based on a strong backbone (e.g. ResNet [16]) to extract features from images with additional feature aggregation schemes such as ASP-OC [47] and PPM [50]. Because

of difficulties associated with semantic segmentation of waterbodies, we design AQUANet to segment aquatic and non-aquatic categories, separately, as shown in Figure 7.

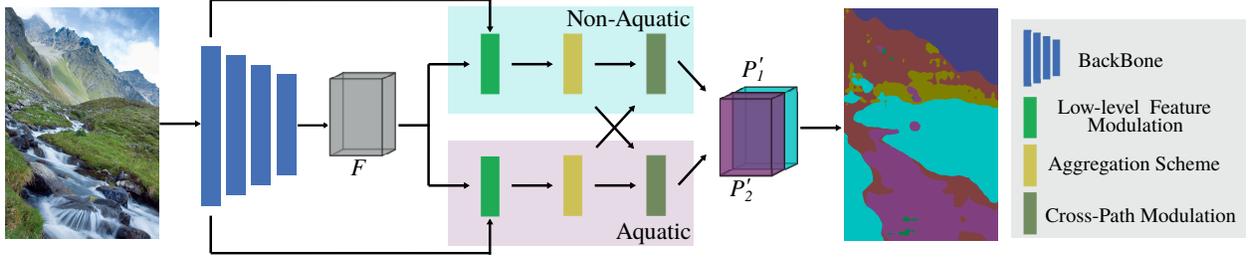


Figure 7: The network architecture of proposed AQUANet.

4.1 Network architecture

According to Figure 7, the input image is first fed into a ResNet-101 (pretrained on ImageNet [9]) to extract the feature F with a size of $C \times H \times W$. Then, the feature is sent into two separate paths for further processing. The aquatic path is to segment different types of waterbodies including sea, river, waterfall, wetland and etc., while the non-aquatic path is to segment other categories such as ship and bridge. In each path, the feature F is first modulated by the low-level feature F_l extracted from the third convolutional layer of ResNet-101, and then passed into ASP-OC [47] to produce the probability map. In the last step, two cross-path modulation blocks are applied to adjust the probability maps P_1 and P_2 in parallel. Finally, the resulted probability maps are concatenated and upsampled to the size of the original image.

4.2 Feature modulation

The goal of the feature modulation \mathcal{M} is to adjust a feature map F_1 given feature map F_2 to represent the adjusted feature F'_1 . It can be formulated as:

$$F'_1 = \mathcal{M}(F_1|F_2). \quad (1)$$

To generate the modulated feature F'_1 , the parameters α and β are learned from F_2 via the feature modulation that consists of three downsampling layers, six 1×1 convolutional layers and two leakyReLU layers as shown in Figure 8. The learned parameters α and β have the shape as the F_1 . Then, the resulting feature F'_1 is constructed as follows according to [43, 34]:

$$F'_1 = \alpha \cdot F_1 + \beta + F_1. \quad (2)$$

4.2.1 Low-level feature modulation

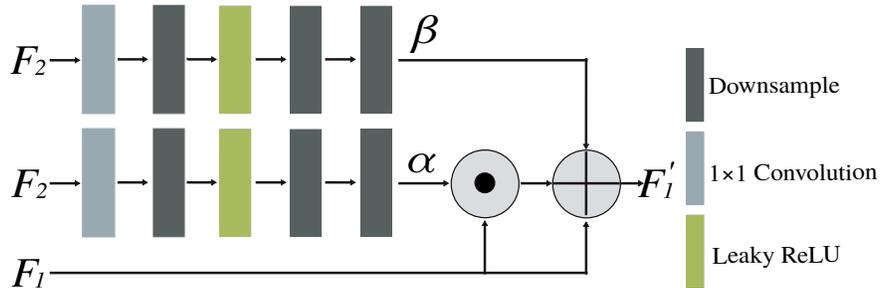


Figure 8: An illustration of the feature modulation.

To enhance the low-level texture representation of the water-bodies, we propose to use the low-level feature F_l to modulate the feature F and the resulted feature F' is defined as:

$$F' = \mathcal{M}(F|F_l). \quad (3)$$

Note that different channels of the receptive field in the third convolutional layer of ResNet-101 are used to construct the low-level feature F_l for the two paths.

4.2.2 Cross-path modulation

We also propose a cross-path modulation to aggregate the outputs of the probability maps. The adjusted probability maps resulted from this module are defined as:

$$P'_1 = \mathcal{M}(P_1|P_2), \quad (4)$$

$$P'_2 = \mathcal{M}(P_2|P_1), \quad (5)$$

where the P'_1 and P'_2 represent the final probability maps for the aquatic and non-aquatic objects, respectively. Note that the modulation layer is not sharing weights across the two paths.

4.3 Loss function

The proposed network is trained in a fully supervised fashion constrained by the widely-used cross-entropy loss function on both final prediction P (main loss) and the intermediate feature produced by the fourth block of the ResNet-101 (auxiliary loss). Following Zhao et al. [50], the weights of the main and the auxiliary loss are set to 1 and 0.4, respectively.

5 Experiments

To demonstrate the effectiveness of the proposed method for water-bodies semantic segmentation, we train AQUANet on the proposed ATLANTIS dataset. For performance evaluation, we take the mean of class-wise intersection over union (mIoU) and the per-pixel accuracy (acc) as the main evaluation metrics. To further evaluate the performance on the waterbodies, we calculate the mean IoU for aquatic categories (A-mIoU) and the accuracy in the aquatic region (A-acc). Aquatic categories include 17 labels showing just water content in different forms and bodies, e.g., sea, river, lake, etc.

Table 3: The per-category results on the ATLANTIS test set by current state-of-the-art methods and our AQUANet. The best and the second best results are highlighted with **bold** font and underline, respectively.

Method	canal	ditch	fjord	flood	glaciers	hot spring	lake	puddle	rapids	reservoir	river	river delta	sea	snow	swimming pool	waterfall	wetland	A-acc (%)	A-mIoU	acc (%)	mIoU
PSPNet	53.8	29.0	42.9	<u>46.5</u>	57.2	53.9	29.7	54.7	38.2	29.8	28.8	65.5	<u>63.5</u>	49.9	47.7	48.4	47.5	66.19	46.29	72.72	40.85
DeepLabv3	52.5	27.2	<u>52.3</u>	43.8	58.7	42.5	31.1	54.2	46.0	32.4	27.1	51.1	61.5	46.3	53.6	52.8	52.9	65.83	46.25	69.21	36.23
DANet	50.5	34.1	37.1	37.0	51.0	<u>61.6</u>	23.8	51.5	42.8	30.2	31.5	63.5	60.4	<u>50.8</u>	43.1	55.2	54.6	62.00	45.80	<u>74.12</u>	39.60
CCNet	41.1	17.4	35.2	26.9	43.7	47.9	18.6	43.8	29.9	16.6	23.7	48.3	53.3	47.6	38.4	51.1	34.1	51.98	36.33	70.84	36.11
EMANet	46.1	16.6	27.1	23.0	53.8	63.7	17.2	43.6	42.2	17.2	21.0	68.6	53.5	47.3	36.1	52.1	36.2	55.88	39.13	71.93	36.43
ANNet	50.9	22.8	31.6	32.0	53.1	58.1	25.6	52.9	48.4	20.8	28.6	56.8	60.4	51.1	43.9	57.9	51.4	61.51	43.90	74.06	39.79
GCNet	56.6	19.0	44.7	34.8	46.9	36.1	35.8	39.4	39.9	41.6	32.4	67.0	62.2	46.4	42.9	50.7	<u>59.7</u>	69.89	44.48	68.64	37.73
DNLNet	54.4	26.3	48.8	36.3	63.2	55.3	<u>35.5</u>	52.3	40.4	32.1	31.3	37.1	61.7	48.3	<u>52.4</u>	48.7	54.6	67.72	45.80	71.95	39.97
OCNet	<u>56.4</u>	<u>33.6</u>	48.0	37.3	57.7	55.2	29.2	50.6	43.8	35.1	35.6	<u>65.9</u>	62.7	47.2	47.9	53.1	54.9	67.97	<u>47.89</u>	73.54	<u>41.19</u>
OCRNet	52.4	19.4	46.9	34.9	48.3	58.8	30.4	39.7	42.5	29.8	31.9	55.5	55.4	47.3	43.6	<u>56.8</u>	51.5	65.90	43.83	71.66	36.17
Ours	<u>55.0</u>	27.7	53.4	47.0	<u>63.1</u>	60.5	33.2	<u>54.4</u>	<u>46.3</u>	<u>39.0</u>	<u>34.7</u>	63.2	64.2	50.3	44.9	53.0	66.1	<u>68.63</u>	50.34	75.18	42.22

5.1 Experimental settings

The AQUANet is implemented using PyTorch. During training, the base learning rate is set to 2.5×10^{-4} and it is decayed following the poly policy [50]. The network is optimized using SGD with a momentum of 0.9 and weight decay of 0.0001. In total, we train the network for 30 epochs, around 80K iterations with a batch size of 2. The training data are augmented with random horizontal flipping, random scaling ranging from 0.5 to 2.0 and random cropping with the size of 640×640 .

5.2 Comparisons

We use several state-of-the-art networks to perform training and testing on ATLANTIS, including PSPNet [50], DeepLabv3 [7], CCNet [20], EMANet [23], ANNet [52], DANet [14], DNLNet [45], GCNet [4], OCNet [47], OCRNet

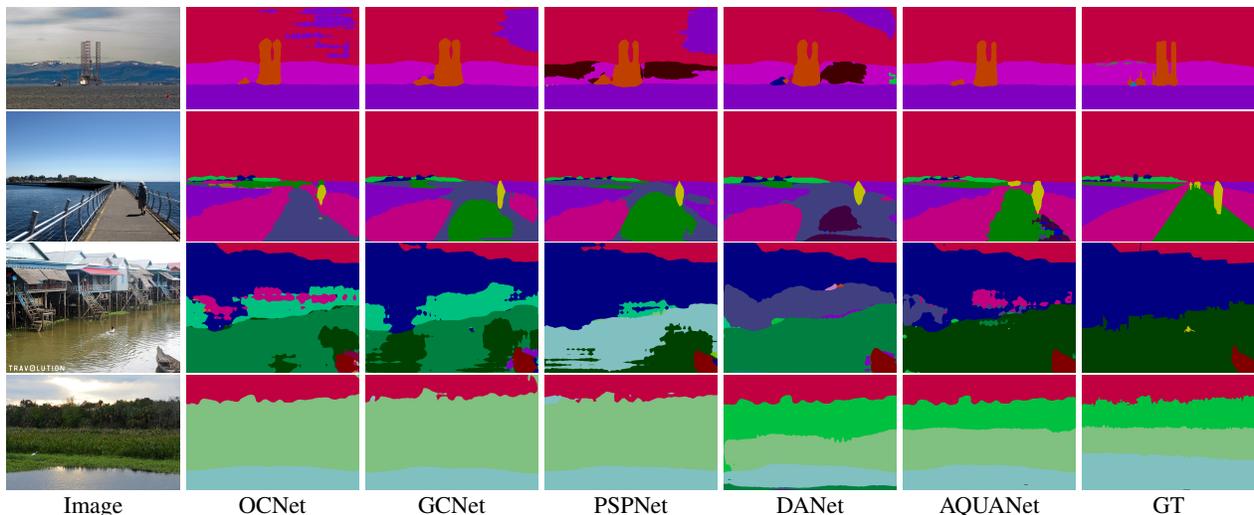


Figure 9: Visualisation comparison of AQUANet and four well-known methods on the ATLANTIS validation set.

[48]. For a fair comparison, we train all the networks with the same backbone (ResNet-101) for 30 epochs. As shown in Table 3, the proposed AQUANet outperforms all these networks on waterbody image semantic segmentation. Figure 9 shows the visualization results of some samples from ATLANTIS validation set. Considering the ground truth and comparing with other networks’ outputs, the boundaries between different classes are better preserved in AQUANet output. Compared with Chen et al. [7], Zhao et al. [50], Fu et al. [14], and Yuan and Wang [47], our method achieves better results on both the aquatic and non-aquatic regions.

5.3 Failure cases

Due to the challenges associated with segmentation of waterbody images, there are still many failure cases we found in the testing stage. Three failure examples are shown Figure 10, from which we observe that many aquatic classes are vulnerable to mis-classification – here sea is mis-classified to lake and river (row 1-2) and river is mis-classified to canal (row 3).

5.4 Ablation studies

We also conduct ablation studies to compare a number of different model variants of the proposed network, including the design of aquatic and non-aquatic paths, and the two feature modulations. The results are shown in Table 4. We can see that the design of two paths can improve the performance of waterbody image semantic segmentation. Moreover, both the proposed low-level feature modulation (LM) and the cross-path modulation (CM) can achieve certain performance gains in terms of acc and mIoU.

Table 4: Ablation study of each proposed component of AQUANet on the ATLANTIS dataset.

Two Paths	LM	CM	A-acc	A-mIoU	acc	mIoU
			67.27	44.53	73.28	38.81
✓			67.73	47.89	75.29	40.28
✓	✓		68.11	47.90	75.29	40.28
✓		✓	67.21	46.63	75.81	40.57
✓	✓	✓	68.85	48.42	76.18	40.83

Comparing to other state-of-the-art semantic segmentation models, AQUANet provides highest mIoU and accuracy. In addition, by considering just labels which includes water content, AQUANet still works better than others (mIoU=50.34%). In this regard, OCNet provides the second highest performance (mIoU=47.89%) and GCNet also provides the best results for canal, lake and reservoir. These results approved AQUANet is well customized for analysis of waterbodies and water-related scenes. However, considering mIoU as a more informative metric for semantic

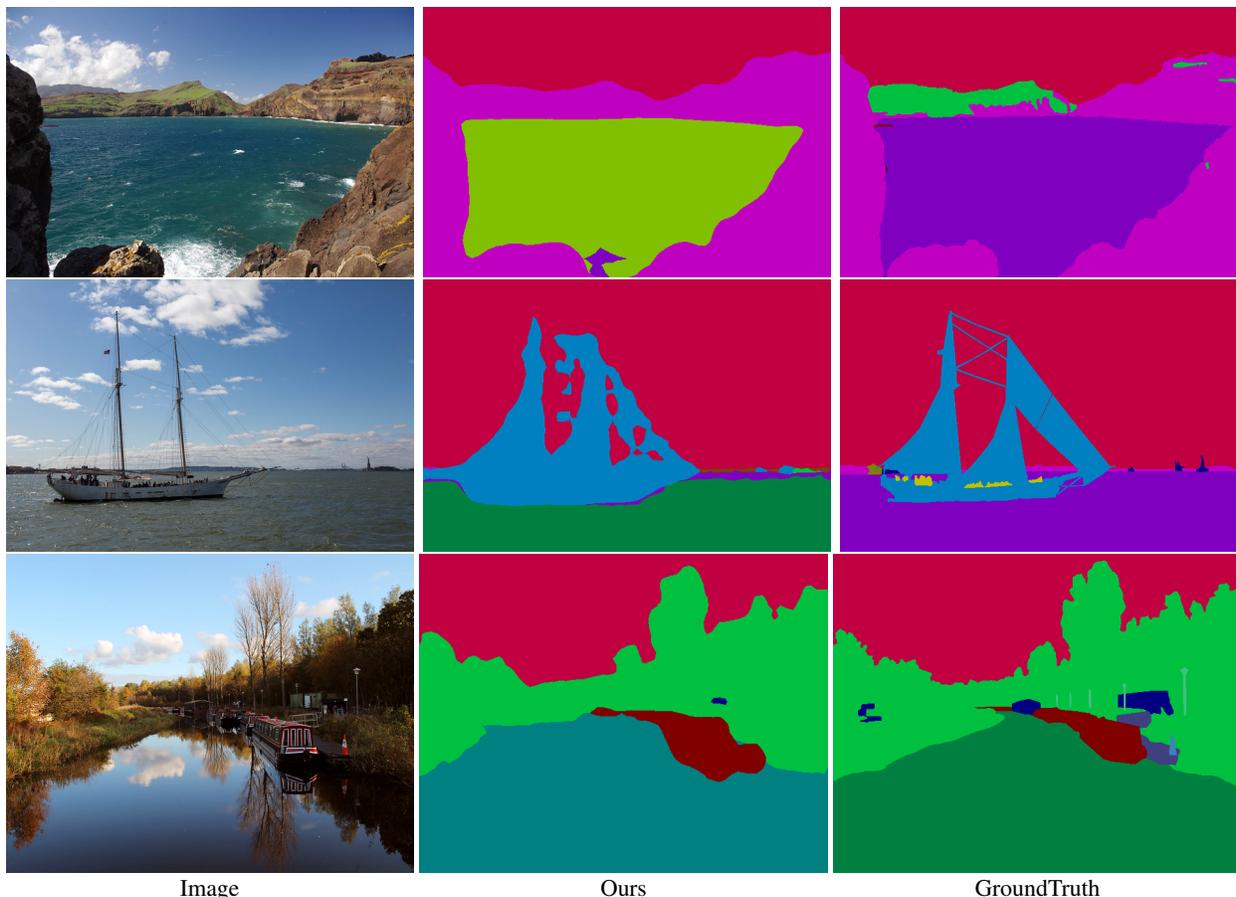


Figure 10: Failure cases from the ATLANTIS val set.

segmentation task, all approaches could only achieve 36.11%~42.22% mIoU on the proposed dataset. It shows that the semantic segmentation of waterbodies and related objects is a challenging task and needs more researches.

5.5 ATeX Experimental results

We further train ten well-known classification models including VGG [40], ResNet [16], SqueezeNet [21], DenseNet [18], GoogLeNet [41], ShuffleNet v2 [29], MobileNetV2 [35], ResNeXt [44], Wide ResNet [49] and EfficientNet [42] on the proposed ATeX dataset. All models are implemented using PyTorch. Cross-entropy loss function is applied for training networks. We train all the networks with the same 30 training epochs, SGD optimizer with a momentum of 0.9 and weight decay of 0.0001, and batch size is set to 64. For all networks, the learning rate is first set to 2.5×10^{-4} , then it is adjusted based on decaying rate of the resulted loss function during training. Table 5 shows the training time and learning rate for each models over certain 30 epochs.

Three common performance metrics including Precision, Recall and F1-score are reported to evaluate the performance of the models on ATeX. Table 5 shows weighted average (averaging the support-weighted mean per label) of these three metrics on the test set. Accordingly, EffNet-B7, EffNet-B0 and ShuffleNet V2 $\times 1.0$ provide the best results. Considering training time, ShuffleNet V2 $\times 1.0$ can be presented as the most efficient network.

6 Conclusion

In this paper, we introduced ATLANTIS, a large-scale dataset for semantic segmentation of waterbodies and water-related scenes, by carefully collecting images of diverse area from the internet and providing high-quality annotations with the help of annotators major in water resources engineering. We further provided comprehensive analysis of the characteristic of ATLANTIS and reported the performance of current state-of-the-arts by training and testing

Table 5: The performance result on ATeX test set by well-known classification models.

Networks	Time [mm:ss]	LR	Val Acc.	Prec.	Test Recall	F1
Wide-ResNet-50-2	06:56	2.5E-4	91	77	75	75
VGG-16	04:38	2.5E-4	90	75	72	72
SqueezeNet	00:47	7.5E-4	82	81	81	81
ShuffleNet V2 $\times 1.0$	01:46	1.0E-2	90	90	90	90
ResNeXt-50-32 $\times 4d$	03:15	2.5E-4	90	77	75	75
ResNet-18	01:28	2.5E-4	87	74	72	72
MobileNet V2	01:35	2.5E-4	88	74	72	72
GoogLeNet	01:51	5.0E-3	89	88	88	88
EffNet-B7	12:42	1.0E-2	90	91	91	91
EffNet-B0	02:38	7.5E-3	91	90	90	90
DenseNet-161	06:15	2.5E-4	91	81	79	79

the networks on our dataset. A novel baseline network AQUANet is also proposed for waterbody image semantic segmentation and achieves the best performance on ATLANTIS. Additionally, we constructed ATLANTIS Texture (ATeX) dataset that are sampled from ATLANTIS for texture-based classification and evaluated several baseline methods on it.

ATLANTIS posed significant challenges for semantic segmentation which we believe will boost new insights in both water resources engineering and computer vision communities.

References

- [1] Vijay Badrinarayanan, Ankur Handa, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for robust semantic pixel-wise labelling. *arXiv preprint arXiv:1505.07293*, 2015.
- [2] Jorge Bernal, Nima Tajkbaksh, Francisco Javier Sánchez, Bogdan J Matuszewski, Hao Chen, Lequan Yu, Quentin Angermann, Olivier Romain, Bjørn Rustad, Ilanko Balasingham, et al. Comparative validation of polyp detection methods in video colonoscopy: results from the miccai 2015 endoscopic vision challenge. *IEEE Transactions on Medical Imaging*, 36(6):1231–1249, 2017.
- [3] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1209–1218, 2018.
- [4] Yue Cao, Jiarui Xu, Stephen Lin, Fangyun Wei, and Han Hu. Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In *Int. Conf. Comput. Vis. Worksh.*, pages 0–0, 2019.
- [5] Guido Cervone, Elena Sava, Qunying Huang, Emily Schnebele, Jeff Harrison, and Nigel Waters. Using twitter for tasking remote-sensing data collection and damage assessment: 2013 boulder flood case study. *International Journal of Remote Sensing*, 37(1):100–124, 2016.
- [6] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(4):834–848, 2017.
- [7] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [8] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3213–3223, 2016.
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 248–255. IEEE, 2009.
- [10] Lunhao Duan and Xiangyun Hu. Multiscale refinement network for water-body segmentation in high-resolution satellite imagery. *IEEE Geoscience and Remote Sensing Letters*, 17(4):686–690, 2019.
- [11] Anette Eltner, Melanie Elias, Hannes Sardemann, and Diana Spieler. Automatic image-based water stage measurement for long-term observations in ungauged catchments. *Water Resources Research*, 54(12):10–362, 2018.
- [12] Anette Eltner, Patrik Olā Bressan, Thales Akiyama, Wesley Nunes Gonçalves, and José Marcato Junior. Using deep learning for automatic water stage measurements. *Water Resources Research*, 57(3):e2020WR027608, 2021.

- [13] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vis.*, 88(2):303–338, 2010.
- [14] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3146–3154, 2019.
- [15] Asmamaw Gebrehiwot, Leila Hashemi-Beni, Gary Thompson, Parisa Kordjamshidi, and Thomas E Langan. Deep convolutional neural network for flood extent mapping using unmanned aerial vehicles data. *Sensors*, 19(7):1486, 2019.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016.
- [17] Yukiko Hirabayashi, Roobavannan Mahendran, Sujan Koirala, Lisako Konoshima, Dai Yamazaki, Satoshi Watanabe, Hyungjun Kim, and Shinjiro Kanae. Global flood risk under climate change. *Nature Climate Change*, 3(9):816–821, 2013.
- [18] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4700–4708, 2017.
- [19] Xiao Huang, Cuizhen Wang, and Zhenlong Li. A near real-time flood-mapping approach by integrating social media and post-event satellite imagery. *Annals of GIS*, 24(2):113–123, 2018.
- [20] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Int. Conf. Comput. Vis.*, pages 603–612, 2019.
- [21] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [22] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *International Conference on Multimedia Modeling*, pages 451–462. Springer, 2020.
- [23] Xia Li, Zhisheng Zhong, Jianlong Wu, Yibo Yang, Zhouchen Lin, and Hong Liu. Expectation-maximization attention networks for semantic segmentation. In *Int. Conf. Comput. Vis.*, pages 9167–9176, 2019.
- [24] Ziyao Li, Rui Wang, Wen Zhang, Fengmin Hu, and Lingkui Meng. Multiscale features supported deeplabv3+ optimization scheme for accurate water semantic segmentation. *IEEE Access*, 7:155787–155804, 2019.
- [25] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1925–1934, 2017.
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Eur. Conf. Comput. Vis.*, pages 740–755. Springer, 2014.
- [27] Shi-Wei Lo, Jyh-Horng Wu, Fang-Pang Lin, and Ching-Han Hsu. Visual sensing for urban flood monitoring. *Sensors*, 15(8):20006–20029, 2015.
- [28] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3431–3440, 2015.
- [29] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Eur. Conf. Comput. Vis.*, pages 116–131, 2018.
- [30] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 891–898, 2014.
- [31] David F Muñoz, Paul Muñoz, Hamed Moftakhari, and Hamid Moradkhani. From local to regional compound flood mapping with deep learning and data fusion techniques. *Science of The Total Environment*, 782:146927, 2021.
- [32] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Buló, and Peter Kontschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *Int. Conf. Comput. Vis.*, pages 4990–4999, 2017.
- [33] Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han. Learning deconvolution network for semantic segmentation. In *Int. Conf. Comput. Vis.*, pages 1520–1528, 2015.
- [34] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

- [35] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4510–4520, 2018.
- [36] Salih Sarp, Murat Kuzlu, Mecit Cetin, Cem Sazara, and Ozgur Guler. Detecting floodwater on roadways from image data using mask-r-cnn. In *2020 International Conference on INnovations in Intelligent SysTems and Applications (INISTA)*, pages 1–6. IEEE, 2020.
- [37] Cem Sazara, Mecit Cetin, and Khan M Iftekharuddin. Detecting floodwater on roadways from image data with handcrafted features and deep transfer learning. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pages 804–809. IEEE, 2019.
- [38] Emily Schnebele and Guido Cervone. Improving remote sensing flood assessment using volunteered geographical data. *Natural Hazards and Earth System Sciences*, 13(3):669–677, 2013.
- [39] Boris Sekachev, Nikita Manovich, Maxim Zhiltsov, Andrey Zhavoronkov, Dmitry Kalinin, Ben Hoff, TOsmanov, Dmitry Kruchinin, Artyom Zankevich, DmitriySidnev, Maksim Markelov, Johannes222, Mathis Chenuet, a andre, telenachos, Aleksandr Melnikov, Jijoong Kim, Liron Ilouz, Nikita Glazov, Priya4607, Rush Tehrani, Seungwon Jeong, Vladimir Skubriev, Sebastian Yonekura, vugia truong, zliang7, lizhming, and Tritin Truong. opencv/cvat: v1.1.0, August 2020. URL <https://doi.org/10.5281/zenodo.4009388>.
- [40] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [41] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1–9, 2015.
- [42] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Int. Conf. on Mach. Learn.*, pages 6105–6114. PMLR, 2019.
- [43] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 606–615, 2018.
- [44] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1492–1500, 2017.
- [45] Minghao Yin, Zhulian Yao, Yue Cao, Xiu Li, Zheng Zhang, Stephen Lin, and Han Hu. Disentangled non-local neural networks. In *Eur. Conf. Comput. Vis.*, pages 191–207. Springer, 2020.
- [46] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2636–2645, 2020.
- [47] Yuhui Yuan and Jingdong Wang. Ocnet: Object context network for scene parsing. *arXiv preprint arXiv:1809.00916*, 2018.
- [48] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *Eur. Conf. Comput. Vis.*, 2020.
- [49] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [50] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2881–2890, 2017.
- [51] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *Int. J. Comput. Vis.*, 127(3):302–321, 2019.
- [52] Zhen Zhu, Mengde Xu, Song Bai, Tengting Huang, and Xiang Bai. Asymmetric non-local neural networks for semantic segmentation. In *Int. Conf. Comput. Vis.*, pages 593–602, 2019.