

# Deep Semantic Segmentation of Trees Using Multispectral Images

Irem Ulku , Erdem Akagündüz , *Member, IEEE*, and Pedram Ghamisi , *Senior Member, IEEE*

**Abstract**—Forests can be efficiently monitored by automatic semantic segmentation of trees using satellite and/or aerial images. Still, several challenges can make the problem difficult, including the varying spectral signature of different trees, lack of sufficient labelled data, and geometrical occlusions. In this article, we address the tree segmentation problem using multispectral imagery. While we carry out large-scale experiments on several deep learning architectures using various spectral input combinations, we also attempt to explore whether hand-crafted spectral vegetation indices can improve the performance of deep learning models in the segmentation of trees. Our experiments include benchmarking a variety of multispectral remote sensing image sets, deep semantic segmentation architectures, and various spectral bands as inputs, including a number of hand-crafted spectral vegetation indices. From our large-scale experiments, we draw several useful conclusions. One particularly important conclusion is that, with no additional computation burden, combining different categories of multispectral vegetation indices, such as NVDI, atmospherically resistant vegetation index, and soil-adjusted vegetation index, within a single three-channel input, and using the state-of-the-art semantic segmentation architectures, tree segmentation accuracy can be improved under certain conditions, compared to using high-resolution visible and/or near-infrared input.

**Index Terms**—Satellite imagery, semantic segmentation, tree segmentation, vegetation indices (VIs).

## I. INTRODUCTION

**F**ORESTS are one of the principal factors that maintain Earth's climate stability [1]. Accurate monitoring of their state and sustainability is a global concern. Compared to other more traditional methods, such as aerial surveys or plot-based analyses, remote sensing has proven to be the most efficient way to monitor forest cover change processes. Conventionally, remote monitoring of processes, such as deforestation or forest degradation has been interpreted manually by expert analysts. However, with recent developments in the field of artificial intelligence, deep learning (DL) algorithms using satellite and/or

aerial imagery are becoming the de facto tool in forest monitoring [2].

Remote sensing imagery can be categorized by the altitude of the aircraft/spacecraft or the type of the optical system (or, simply put, the sensor) [3]. These two factors largely determine the tradeoff between the area covered and the amount of detail (i.e., the resolution) of the constructed image. Depending on the type of the sensor, the optical system can sense a part of the electromagnetic spectrum, sample different spectral bands separately and simultaneously, and hence, create multi or hyperspectral imagery. In addition, various spectral bands are used to construct the so-called “spectral indices,” which are mathematical combinations of spectral reflectance of different wavelengths. These hand-crafted indices are used to detect the presence of objects or situations, such as vegetation [4], water [5], fire [6], and landslide [7]. Although these hand-crafted indices are valuable features for pattern recognition, today's trend is moving toward building nontransparent, end-to-end DL models to detect any feature of interest, usually using supervised techniques [8].

In this article, we address the problem of segmenting forest areas (or simply trees) in both satellite and aerial images using deep semantic segmentation techniques and multispectral imagery. We attack the problem in multiple dimensions. To begin with, we examine the effect of altitude using satellite and aerial images with various metric and spectral resolutions. Second, we benchmark different deep semantic segmentation architectures, including models with pretrained convolutional encoders. We utilize different spectral bands, such as visible [red-green-blue (RGB)], near-infrared (NIR), and short-wave infrared (SWIR) (and their multispectral combinations), as raw inputs to these models. Moreover, we experiment with hand-crafted spectral vegetation indices (VIs) as input and analyze their performance compared to raw spectral input images. Our goal is to determine how to efficiently utilize DL-based architectures for tree segmentation, what practical issues arise with limited remote sensing data, and which methods should be utilized. In order to accomplish this, we conduct large-scale experiments.

The rest of this article is organized as follows. In the following section, we refer to the literature on the subject and underline our contribution. The third section provides the details of our benchmarking environment with respect to the image sets, the segmentation models, and the spectral indices we utilize in our experiment. Section IV presents the experimental results and covers related analyses and discussions. Finally, Section V concludes this article.

Manuscript received 21 December 2021; revised 5 March 2022, 4 April 2022, 10 May 2022, and 14 June 2022; accepted 25 August 2022. Date of publication 31 August 2022; date of current version 14 September 2022.

Irem Ulku is with the Department of Computer Engineering, Ankara University, 06560 Ankara, Turkey (e-mail: iremulku@gmail.com).

Erdem Akagündüz is with the Graduate School of Informatics, Middle East Technical University, 06800 Ankara, Turkey (e-mail: akaerdem@metu.edu.tr).

Pedram Ghamisi is with the Institute of Advanced Research in Artificial Intelligence, 1030 Vienna, Austria, and also with the Helmholtz-Zentrum Dresden-Rossendorf, Helmholtz Institute Freiberg for Resource Technology, 09599 Freiberg, Germany (e-mail: p.ghamisi@gmail.com).

Digital Object Identifier 10.1109/JSTARS.2022.3203145

## II. RELATED LITERATURE

Survey studies on object detection in optical remote sensing images [9], [10], [11] show that efforts to detect trees in satellite images date back to the early 90s [12], [13], [14]. Previously, traditional machine learning methods, such as multinomial logistic regression [15], random forest (RF) [16], [17], [18], and support vector machine [19], [20], [21] dominated the area of tree detection and segmentation within remote imaging. These studies relied not only on hand-crafted feature extractions but also on feature encoding [22] and feature pooling [23]. Because of the low-level characteristics of the extracted features, these methods were not able to capture high-level features that reside in complex high-dimensional (spatial, temporal, and spectral) satellite or aerial imagery. Recently, DL approaches have emerged to tackle such limitations by making use of hierarchical learning processes while extracting high level, complex abstractions from the data [24], [25], [26], [27], [28], [29]. Convolutional neural networks (CNNs) [30], with their proven ability to learn these high-level abstract feature representations with the help of convolutional and pooling layers, are the most commonly used DL models for semantic segmentation.

Following the recent interest in deep semantic segmentation architectures [31], image context exploitation and fusing features from different stages for image segmentation have given rise to various methods. Pioneered by the fully convolutional network (FCN) [32] and followed by many improvements [33], [34], [35], [36], [37], [38], the encoder–decoder (ED) architectures show promising success in solving semantic segmentation tasks. In fact, the ED-based CNN architectures are utilized for semantic segmentation of Earth observation data [37], [39], [40], [41], [42], [43], [44], [45], allowing the fusion of high- and low-level feature maps while improving spatial accuracy in semantic segmentation.

It is also worth mentioning that classical tree segmentation methods commonly utilized 3-D airborne light detection and/or range (LIDAR) data [46], [47], [48], [49], [50], [51], [52], [53], [54], [55]. For instance, when using 3-D LIDAR data to perform tree segmentation, data-specific methods that relied on crafted features were employed, such as canopy height model [56], based watershed algorithm [57], or normalized point cloud-based layer stacking algorithm [58]. Besides, in some studies both LIDAR and multispectral data were utilized by applying different algorithms for each data type. While implementing machine/DL algorithms that are well-suited to multispectral data, the studies concurrently applied methods native to LIDAR data [59], [60], [61]. Nevertheless, due to its flexibility in both spatial and spectral resolution, and its compatibility with different remote sensing systems (UAV, satellite, etc.), multispectral and/or hyperspectral imagery are today's main standards for not only tree segmentation, but for several other related problems, such as crop, water, and field segmentation [3].

The specific problem of detecting trees presents some challenges, such as geometrical complexity, color variations, and/or self-occlusion of the tree branches [62], just to name a few. As a result of these challenges, the problem proves to be considerably difficult. For example, regarding self-occlusion, even a slight

change in the viewpoint could cause the same structures to appear significantly different. The tree segmentation problem is also challenging due to the environmental background, which may contain other types of green vegetation, such as small bushes. A closely related problem is differentiating between the view of a large canopy, hence, the overlap between adjacent individual trees. A satellite-retrieved spectral signature is the combination of different reflectance sources, such as tree crown, tree crown shadow, soil, and herbaceous vegetation [17]. Another key factor that limits individual tree mapping is that the spatial resolutions of most satellite/aerial sensors are usually not capable of capturing objects with sizes less than one square meter. However, thanks to commercial satellites and UAVs with multispectral cameras, the improved spatial resolution makes it possible to map every tree on large scale [63].

Instead of using individual spectral bands, VIs have been proposed in the literature [64]. VIs are derived from varying mathematical combinations of two or more spectral reflectance values to examine the presence of the vegetation within a pixel. VIs have several drawbacks, notably their inadequacy at interpreting the information content for forest canopy, thereby insufficiently quantifying the effects among different variables, and hence, making them less favorable for rapidly changing biophysical factors, such as chlorophyll content per unit leaf area, leaf angles, and fractional cover [65]. Performances of VIs are compared to that of the spectral mixture analysis in terms of tree canopy cover and are found to be poor [66]. Since there is a high correlation between many VIs, they can be used together and complement each other in different aspects by addressing specific needs that characterize tree canopy.

Most of the DL architectures that process satellite/aerial imagery for the segmentation of objects, including trees, are end-to-end models that take multi- or hyperspectral bands as input channels [67]. However, due to both their multimodal, geolocated, and multitemporal nature and the insufficient amount of labelled remote-sensing data, applying DL models tends to be a challenging task [25]. Improvements in the high-performance computing architectures, such as GPUs lag far behind the growth of high-dimensional big satellite/aerial data; ergo, recent efforts toward building large-scale DL models fall short of training big remote-sensing data, unless a breakthrough is achieved in the acceleration of computing power. One solution to this problem, which we propose in this article, is to utilize VIs as a priori knowledge to the DL models. VIs by definition, provide insight into spectral characteristics of a surface by taking into account the interaction between electromagnetic radiation and vegetation, and hence, can be considered physics-based modelled features.

Our study explores the effect of fusing VIs as a physical model and various ED architectures as DL models in order to perform tree semantic segmentation tasks effectively and accurately. For this purpose, we implement various ED-based semantic segmentation architectures, such as U-Net [34], SegNet [35], DeepLabv3+ [36], DLinkNet [37], and DFANet [38], as well as a pre-DL era ensemble learning algorithm, namely RF [68]. We utilize specific VIs so as to integrate physics-based models as additional input information to the benchmarked models. A framework is constructed for satellite and aerial images that

not only provides a comparison between different ED-based semantic segmentation architectures but also leads to insights about fusing various VIs or spectral bands for the purpose of pixel-wise tree segmentation.

### III. METHODOLOGY

In this section, we provide a methodological background for the experiments conducted within our study. We commence by introducing the multispectral image sets utilized in our experiments. We then provide details on the benchmarked semantic segmentation architectures of our experiments and the types of input we provide to these models.

#### A. Multispectral Image Sets

The previous studies in the literature have proven that high spectral resolution is a major factor in solving the tree segmentation problem, as is the requirement of high spatial resolution [16]. The earlier airborne cameras and satellite-borne multispectral sensors provided high spatial resolution, while their spectral bands were limited to blue (0.45–0.51  $\mu\text{m}$ ), green (0.51–0.58  $\mu\text{m}$ ), red (0.6–0.69  $\mu\text{m}$ ), and NIR 1 (0.77–0.90  $\mu\text{m}$ ). Because several satellite-borne sensors, such as WorldView-2 and WorldView-3 were launched after 2009, numerous other spectral bands that are strongly related to vegetation properties have also been included. One is the Coastal band (0.40–0.45  $\mu\text{m}$ ), where the associated reflectance is highly bound up with the chlorophyll content of the vegetation. Another band, the red edge (0.71–0.75  $\mu\text{m}$ ), leads to intuitive discrimination of healthy trees, which is the primary task in precision agriculture. Aside from the detection of healthy trees, infected trees can also be exposed if the Yellow band (0.59–0.63  $\mu\text{m}$ ) is utilized to extract the yellowness of their crowns. Apart from the NIR 1 band, the NIR 2 band (0.86–1.04  $\mu\text{m}$ ) can provide additional information on the vegetation analysis since it is less influenced by atmospheric conditions [69]. The SWIR (SWIR) band (1.19–2.36  $\mu\text{m}$ ) is useful to characterize vegetation water content since it is sensitive to water absorption.

In order to conduct supervised tree segmentation experiments, multispectral image sets with pixel-label ground truth are required. Table I shows detailed information about the two image sets that we use in our conducted experiments. In Table I, ground sampling distance (GSD) refers to the distance between the centres of two adjacent pixels corresponding to their representation in the real world [70]. The improvement of segmentation performance is significantly related to the spatial resolution, with GSD being the most important descriptor [71]. Additionally, ground field-of-view (FOV) refers to the metric area covered by the sensor at the chosen operational altitude of the sensor during image collection. The reader may refer to [72], [73], [74], [75] for a list of available multispectral image sets with pixel labels. The two image sets used in our experiments are discussed as follows.

1) *DSTL Satellite Imagery Feature Detection Image Set*: The DSTL Satellite Imagery Feature Detection image set is provided by the Defense Science and Technology Laboratory [76], and includes 25 very high-resolution training images. Each image

TABLE I  
FEATURES OF BENCHMARK MULTISPECTRAL SEMANTIC SEGMENTATION  
IMAGE SETS

Features	Datasets	
	DSTL (satellite)	RIT18 (aerial)
<b>GSD</b>	0.31 m (RGB) 1.24 m (VIS+NIR) 7.5 m (SWIR)	0.047 m (VNIR)
<b>Ground FOV</b>	25 km <sup>2</sup>	2.88 km <sup>2</sup>
<b>Sensor</b>	World-View 3	TetraCam MicroMCA6
<b>Year</b>	2017	2017
<b>No. of spectral bands</b>	16	6
<b>Visible subbands</b>	400–452 nm 448–510 nm 518–586 nm 590–630 nm 632–692 nm	485–495 nm 545–555 nm 675–685 nm
<b>NIR subbands</b>	706–746 nm 772–890 nm 866–954 nm	715–725 nm 795–805 nm 890–910 nm
<b>SWIR subbands</b>	1195–1225 nm 1550–1590 nm 1640–1680 nm 1710–1750 nm 2145–2185 nm 2185–2225 nm 2235–2285 nm 2295–2365 nm	none
<b>Bit resolution</b>	11-bits/pixel (Visible and NIR) 14-bits/pixel (SWIR)	10-bits/pixel
<b>No. of 224 × 224 patches</b>	5985	1778

covers an area of 1 km × 1 km of the Earth's surface. In this image set, both 3-band and 16-band formats are provided with different spatial resolutions. The images are captured by WorldView-3, which is a commercial satellite. The WorldView-3 optical system provides very high-resolution satellite imagery with 31 cm panchromatic resolution, 1.24 m multispectral resolution, and 7.5 m SWIR resolution [77]. WorldView-3 satellite captures the same area with different types of satellite images, such as a panchromatic band (high-resolution), an RGB band (high-resolution), and a multispectral band (*M*-band) (lower resolution), and an SWIR band (*A*-band) (lowest resolution). The advantageous feature of having 11- and 14-bit as image color depths enables more information from each pixel to be used in a DL model. There are a total of 10 labeled object types, including trees.

2) *RIT-18 (The Hamlin State Beach Park) Aerial Image Set*: The RIT-18 aerial image set [73] was acquired in 2017 by mounting the Tetracam MicroMCA6 multispectral sensor on board the DJI-S1000 octocopter. This unmanned aircraft system is much cheaper than those based on manned aircraft and satellite systems. RIT-18 has been collected with a ground sample distance of 0.047 m and includes six very high-resolution multispectral bands. Three of them are visible RGB bands, while the other three are NIR bands, with all having very high spatial resolution. Due to their high spatial resolution,<sup>1</sup> the NIR bands

<sup>1</sup>10 and 20 nm band-pass filters were used for this sensor since wider bandwidths tend to create oversaturation [73].

introduced by the RIT-18 are strengthening the discriminative power of DL architectures, especially for cases where vegetation exists [73].

### B. Semantic Segmentation Architectures

In the past decade, studies on different semantic segmentation architectures have shown tremendous development [31]. Mostly driven by industry, the majority of these studies were designed for medical applications or unmanned systems, such as driverless cars. As mentioned in Section II, we benchmark different DL-based semantic segmentation architecture and a pre-DL ensemble learning method called the RF algorithm. In this part, we provide details on the semantic segmentation techniques that we utilize in our experiments.

The pioneering DL-based semantic segmentation architecture is the FCN [32]. This network is considered to be the ancestor of today's widely used ED architectures in semantic segmentation [31]. Although there are many different architectures and approaches in DL-based semantic segmentation literature, it is safe to say that the ED architecture is an off-the-shelf solution to many semantic segmentation problems.

Due to its role in recovering spatial resolution in extracted feature maps, the decoder is especially important to implementing ED architectures for multispectral imagery so that the high-resolution details can be recovered effectively. Given this effect, we considered various ED architectures, while paying attention to the different aspects of decoders. Upsampling layers are one of the most crucial parts of decoders, as they increase the spatial resolution [78]. Therefore, upsampling layers are the basis for selecting the proper architectures in this article.

Below we explain the five fundamental DL-based semantic segmentation architectures that we utilize in our experiments: U-Net[34], SegNet [35], DeepLabv3+ [36], DLinkNet [37], and DFANet [38]. The first two are considered to be the basic ED examples, while the latter three are advanced state-of-the-art models that have been proven to be more successful in today's semantic segmentation challenges [79]. Finally, we provide information about the RF algorithm, which is one of the strongest pixel-wise tree segmentation methods of the pre-DL era.

1) *U-Net* [34]: U-Net is the one of the pioneering and most-used ED approaches in the semantic segmentation literature [see Fig. 1(a)]. It is notable for its unique architecture, in which the whole feature map is transported from the encoder to the corresponding building block in the decoder via skip connections. Compared to state-of-the-art architectures, the decoder part of U-Net is computationally inefficient since it reconstructs the segmentation map, the size of which is equal to that of the original image through many up-scaling operations between high-level semantic information and precise local information. However, U-Net is still considered to be a de facto solution to any semantic segmentation problem.

The implemented U-Net model is a four-level-depth symmetric ED model. At each depth, two repeated blocks of convolutional layers are followed by a batch normalization operation during training and a rectified linear unit (ReLU) activation function in the experiment itself. Attached to these two blocks,

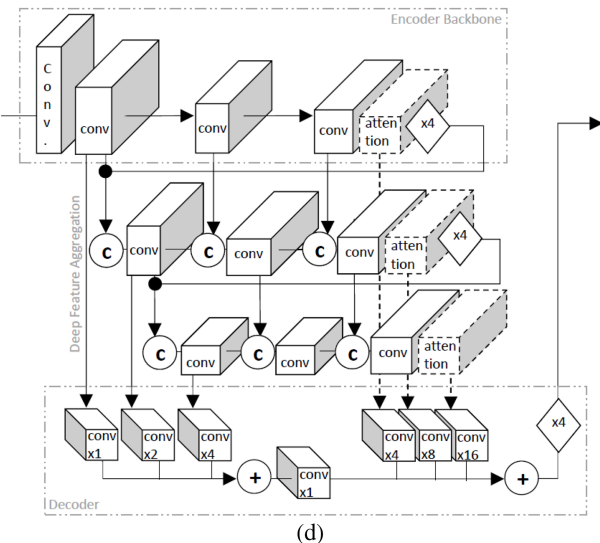
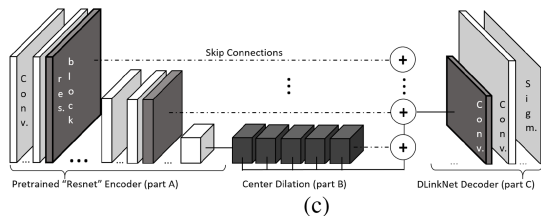
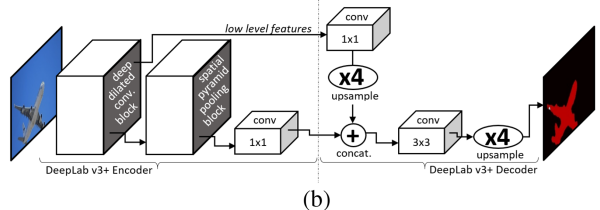
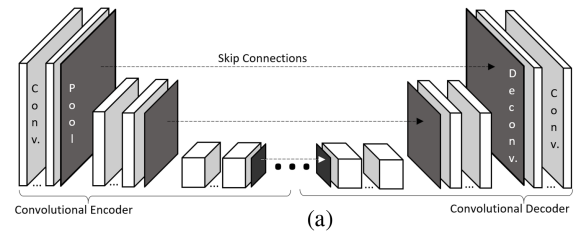


Fig. 1. Different semantic segmentation architectures were utilized in the experiments. (a) Basic ED architecture: U-Net (features as skip connections) or SegNet (pooling indices as skip connections). (b) DeepLabv3+ architecture. (c) DlinkNet architecture. (d) DFANet architecture.

each level is concluded by a  $2 \times 2$  max-pooling layer (stride 2) and a dropout layer with rate = 0.5, hence halving the spatial dimension at each depth. For each depth, the feature channels vary among 64, 128, 256, and 512, which is also symmetrical in the decoder. In the decoder part of the process, instead of the max-pooling layer, there are up-convolution layers that help double the spatial dimensions. For all convolutional layers, kernel, stride, and padding sizes are set as  $3 \times 3$ , 1  $\times$  1, and

$1 \times 1$ , respectively, thus keeping the spatial dimensions fixed within a block. In order to implement the binary classification of trees, the final decoder block is followed by a single depth  $1 \times 1$  convolution layer and pixel-wise sigmoid layers, which provides an input-sized single channel map of the detection measure. Skip connections between each depth carry feature activation values, a trademark of the original U-Net architecture.

Among the many different ED architectures, we consider U-Net a fitting choice for benchmarking, given that it is a basic example of using transposed convolution as the up-sampling layer [78].

2) *SegNet* [35]: SegNet is very similar to U-Net architecture. The only difference is that, rather than passing the entire feature map of the encoder to the decoder, SegNet uses a very simple yet computationally efficient transfer of information. Instead of unpooled features, maximum pooling indices are passed [see Fig. 1(a)].

In the implemented SegNet model, the architecture and the related parameters are the same as in the implemented U-Net model described earlier. The only difference is, as the trademark of the SegNet architecture, skip connections carry only the pooling indices, which is why SegNet is considered to be an efficient implementation of the U-Net.

Rather than using the transposed convolution as the upsampling layer, such as U-Net, SegNet employs interpolation by implementing the unpooling operation [80], [81]. Therefore, the criterion for selecting SegNet is to compare the performance of this upsampling operation at the decoder.

3) *DeepLabv3+* [36]: One of the most advanced, successful, and efficient semantic segmentation frameworks in the literature is the DeepLabv3+ [82] [see Fig. 1(b)], which is the final architecture of the DeepLab family [82], [83], [84]. It is originally trained with general-purpose image sets, such as PASCAL VOC 2012 [79] and COCO [85]. The main idea of DeepLabv3+ is to consolidate multiscale contextual information by using atrous spatial pyramid pooling layers. Compared to a default ED architecture, the decoder part is extremely lightweight, although with better segmentation accuracy.

As semantic segmentation architectures developed rapidly, recent studies have focused on utilizing state-of-the-art architectures, such as ResNet [86], Inception [87], and Xception [88] in the encoder, referred to as encoder backbone networks. The goals of this usage are to improve the generalization ability and increase the prediction accuracy of the model. Like many other semantic segmentation architectures, DeepLabv3+ can also function with pretrained encoders. In scenarios, where the labeled data is limited, different encoder backbone networks may help the model perform with increased efficacy. In the implemented DeepLabv3+ model, the first two levels of the ResNet34 or ResNet101 encoder are transferred for this purpose. However, because we utilized pretrained encoders that originally accepted three-channel (RGB) input, the experiments with this architecture are limited to spectral input with 3 or fewer channels.

We utilize DeepLabv3+ in our experiments because its up-sampling layer is completely different from that of the other ED architectures. DeepLabv3+ employs bilinear upsampling +

convolution in its decoder design [78], mainly for computational efficiency.

*D-LinkNet* [37]: Another state-of-the-art model in the literature is DLinkNet. This model is originally trained with RGB satellite images for pixel-wise road segmentation, which makes it special in the sense that it is one of the rare models specifically designed and trained using image sets with almost-zero zenith angles. The model is based on an approach focused on additional dilated convolution layers in the centre part of an ED architecture [see Fig. 1(c)]. DLinkNet is an improvement of the Linknet architecture [89], which is also an ED-based semantic segmentation model.

Similarly to the implemented DeepLabv3+ model, DLinkNet's utilized encoder is transferred and fine-tuned. The first two levels of the ResNet34 encoder is transferred to the implemented DLinkNet model in the conducted experiments. Hence, the experiments with this architecture are also limited to spectral input with less than or equal to 3 channels (please refer to Section III-C for model input types).

DLinkNet is also another example of the ED architecture that relies on transposed convolution in the decoder design [37]. However, it is not considered to be as computationally efficient as the DeepLab family.

4) *DFANet* [38]: The majority of recent semantic segmentation studies are targeted at real-time applications like autonomous driving that require fast understanding of the surrounding scenes while achieving high performance [90], [91], [92], [93]. In our study, we chose the state-of-the-art DFANet as a benchmark model because it has a sophisticated ED structure similar to DLinkNet, but is still real-time in performance.

The main contribution of DFANet is the so-called "Cross-level feature aggregation based ED architecture" [see Fig. 1(d)] that aims to reduce the number of parameters, while preserving a certain semantic segmentation performance. DFANet encoder is constructed by first starting from the Xception model as a lightweight backbone, and then providing the high-level feature maps obtained from this backbone as an output to the next Xception backbone. This process, which is called subnetwork aggregation, continues for three parallel Xception backbones. The goal of adding a fully connected attention module to the tail of each backbone is to acquire the maximum receptive field. In addition to this subnetwork aggregation, another module, namely the substage aggregation, helps us to ensure the combination of multiscale information. By recovering the lost spatial information caused by deep architectures, the aim of the substage aggregation module is quite similar to that of the skip connections except that skip connections are not able to keep the large-scale object and edge information in very deep architectures. The substage aggregation module is based on the fusion of different stages of the same depth in subnetworks. In other words, the output of the layer at a certain resolution from the previous backbone is contributing to the input of the corresponding layer of the next backbone.

The encoder, which is composed of the aggregation of three lightweight Xception backbones by means of subnetwork and substage modules, is followed by a slight decoder designed particularly to tackle real-time inference concerns. This decoder

is conceived as an efficient feature upsampling module that is characterized by the convolution and bilinear upsampling layers. This simple decoder structure fuses the high-level features and the low-level features of the three backbones within itself. After the high-level features are bilinearly upsampled by a factor of 4, these are added to the low-level features. Finally, this total sum is upsampled by a factor of 4 and the final prediction is obtained.

5) *Random Forest [68]*: In order to compare the data-driven learning power of deep neural networks relative to traditional machine learning algorithms, the RF method, which relies on hand-crafted features, is also introduced for segmentation performance evaluation. The RF method includes a set of decision trees in which each acts as a base classifier and depends on the independently sampled values of a random vector from the same distribution. As an ensemble learning algorithm, RF builds many decision trees at the time of training and relies on the maximum voting of the decision trees in the forest. Its robustness to overfitting and good generalization ability are among the advantages of the RF algorithm in semantic segmentation [94], [95].

The hand-crafted features that we fed to the RF model are basic image features, specifically derivatives (up to second order) in three different scales. The number of decision trees in the RF ensemble is selected as 20, whereas the minimum number of observations per tree leaf is set to 60.

### C. Model Input

In this study, we are also interested in exploring the correlation of spectral bands and VIs with tree segmentation performance. For that purpose, the bench-marked models are trained and tested with a set of input data with varying channel depth. We categorize the utilized model input into two main categories, three-channel and multispectral input. The first, the three-channel input, includes RGB, single VIs (three channels being the same), the NIR input (including separate three NIR sub-bands) and a three-channel vegetation index combination, including NDVI, atmospherically resistant vegetation index (ARVI), and soil-adjusted vegetation index (SAVI) indices (shown below). Since this input group has a depth of three channels, DL models that comprise pretrained encoders, such as ResNet or Xception can be used in experiments.

The second group, multispectral input, includes image sets with a higher number of channel depths. The characteristics of the multispectral band input depend on the image set (i.e., the sensor) and are explained in detail as follows. As expected, experiments using this input group were only conducted for the U-Net and the SegNet DL models, which can accept any number of input channel depths. On the other hand, RF experiments are applied to both model input groups.

1) *Red-Green-Blue*: As the name implies, these images consist of only red, green, and blue channels. The difference between the reflected radiations in visible (i.e., RGB) and NIR wavelengths give more information about vegetation in particular. A large difference suggests that the pixels are more likely to belong to dense vegetation, such as forest area, while a small one is likely to indicate sparse vegetation, such as grassland [96]. On the other

hand, for most commercial satellites (such as world-view-3), the spatial resolution of RGB channels is exceedingly higher than the original multispectral output of the sensor due to a resolution enhancement postprocess applied to visible band output.

2) *Vegetation Indices*: Vegetation covers exhibit a unique spectral behaviour, by which they can be differentiated from other ground elements [64]. Commonly, chlorophyll concentration is responsible for absorbing the radiation in the red band, while the leaf cellular structure is responsible for reflecting the radiation in the shorter NIR bands (700 to 750 nm). A deviation between the red and NIR is observed in the spectral reflectance curve, providing deeper insight into the existence of vegetation. The key idea is to make use of this deviation to differentiate vegetation from bare soil based on the contrast between the spectral reflectance behaviours. Another spectral radiance difference measured between the red and the blue bands acts as a self-corrector by reducing the atmospheric scattering effects in the red band, making it possible to define atmospherically resilient VIs.

Current VIs are commonly categorized in the following three fundamental families: mean VIs, atmospherically resilient VIs, and soil-adjusted VIs [97]. For this reason, we utilized three widely-used indices from each family in our experiments. The key characteristics of these indices are delineated in the following:

3) *Normalized Difference Vegetation Index (NDVI)*: is selected to represent the mean vegetation index category, for that it is one of the most widely used and well-known indices of this category [98]. NDVI for the DSTL image set is defined by

$$\text{NDVI}^D = \frac{\text{NIR1} - \text{Red}}{\text{NIR1} + \text{Red}} \quad (1)$$

where the WorldView-3 NIR1 band covers the range 0.77–0.90  $\mu\text{m}$  and the red band covers the range 0.63–0.69  $\mu\text{m}$ .

In the case of the RIT-18 image set, NDVI is calculated as follows:

$$\text{NDVI}^R = \frac{\text{NIR} - \text{Red}}{\text{NIR} + \text{Red}} \quad (2)$$

where the RIT-18 NIR band covers the range 0.71–0.91  $\mu\text{m}$  and the red band covers the range 0.67–0.68  $\mu\text{m}$ .

4) *Atmospherically Resistant Vegetation Index*: ARVI is particularly designed to correct atmospheric effects. In order to obtain such an atmospherically resilient VI, blue or green spectral bands are taken into account along with the red and NIR bands. There are also some other related VIs that reduce atmospheric turbidity [99]. ARVI for the DSTL image set is given by

$$\text{ARVI}^D = \frac{\text{SWIR 1} - \text{Red Edge} - \gamma(\text{Red Edge} - \text{Coastal})}{\text{SWIR 1} + \text{Red Edge} - \gamma(\text{Red Edge} - \text{Coastal})} \quad (3)$$

where the WorldView-3 SWIR 1 band covers the range 1.19–1.22  $\mu\text{m}$ , red edge band covers 0.71–0.75  $\mu\text{m}$  and the coastal band covers the range 0.40–0.45  $\mu\text{m}$ , and  $\gamma$  is taken as 1.0.

For the RIT-18 image set, ARVI calculation should be arranged according to the existing NIR bands due to the absence of the SWIR wavelengths. The following shows the calculation

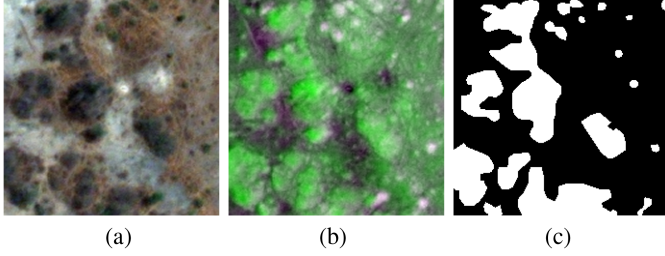


Fig. 2. Example of tree cover in satellite imagery. (a) RGB image. (b) NDVI, ARVI, and SAVI VIs depicted in three-channel false color. (c) Tree cover ground truth for semantic segmentation.

of ARVI for the RIT-18 image set

$$\text{ARVI}^R = \frac{\text{NIR} - \text{Red} - \gamma(\text{Red} - \text{Blue})}{\text{NIR} + \text{Red} - \gamma(\text{Red} - \text{Blue})} \quad (4)$$

where the RIT-18 covers the NIR range 0.71–0.91  $\mu\text{m}$ , the red range 0.67–0.68  $\mu\text{m}$ , and the blue range 0.48–0.49  $\mu\text{m}$ ;  $\gamma$  is taken as 1.0.

5) *Soil-Adjusted Vegetation Index*: Besides the two categories described earlier, there are many other VIs that have been proposed to reduce the soil background noise on NDVI by making use of a parameter called  $L$ , which incorporates the area density factor. One of these, SAVI enables the soil brightness effect to be assessed where the vegetation density is low in the area of interest [100]. For the DSTL image set, SAVI is given by

$$\text{SAVI}^D = \frac{\text{NIR1} - \text{Red}}{\text{NIR1} + \text{Red} + L} (1 + L) \quad (5)$$

where  $L$  is taken as 0.5.

Using the NIR range of 0.71–0.91  $\mu\text{m}$  and a red range of 0.67–0.68  $\mu\text{m}$ , SAVI can be formulated for the RIT-18 image set as follows:

$$\text{SAVI}^R = \frac{\text{NIR} - \text{Red}}{\text{NIR} + \text{Red} + L} (1 + L) \quad (6)$$

where  $L$  is taken as 0.5.

6) *Mixed-Vegetation input (NDVIA + ARVI + SAVI)*: In addition to utilizing the aforementioned VIs separately as input to the benchmarked models, we propose a new three-channel input, called mixed-vegetation input, which uses three different VIs to test the combined effect of minimized soil brightness influences, corrected atmospheric scattering effects, and minimized topographic effects. The VIs selected from each of the categories are as follows: NDVI from mean VIs, ARVI from atmospherically resilient VIs, and SAVI from soil-adjusted VIs. The combined three-channel VI image is fed into the DL architectures just like an RGB input. A false-color image of this fused input is depicted in Fig. 2.

7) *Near Infrared input (NIRi)*: Both the DSTL and RIT-18 datasets utilized in our experiments are shot with sensors that output three NIR channels, details of which are provided in Table I. It can be seen from this table that, the spectral resolution of the RIT-18 dataset (TetraCam MicroMCA6) is finer than DSTL’s sensor (World-View-3). On the other hand, World-View-3 NIR bands cover a larger spectrum. Hence, we believe that

testing these two 3-channel NIR input types in our experiments will provide us with an insight into the effect of NIR spectral resolution on tree segmentation.

8) *Visible + Near Infrared input (VNIR)*: VNIR is one of the multispectral input groups utilized in our experiments. DSTL VNIR input includes eight subbands (five visible + three NIR), whereas RIT-18 includes six subbands (three visible + three NIR) as a result of their individual sensory output (see Table I).

9) *Visible + NIR + SWIR (VNIR + SWIR)*: VNIR + SWIR is the second multispectral input group utilized only for the DSTL image set in our experiments since SWIR bands are out of the RIT-18 sensor’s spectral range. For this purpose, two separate groups are constructed, the first being the VNIR + SWIR<sup>2</sup> input that includes 8 DSTL + VNIR subbands and the two shortest SWIR bands (i.e., SWIR-1 and SWIR-2). The second group is referred to as the VNIR + SWIR<sup>8</sup> and includes all of the original multispectral bands of the World-View-3 sensor in 16 channels.

#### IV. EXPERIMENTAL RESULTS AND DISCUSSIONS

Before presenting the experimental results and the related discussions, we provide the details of the experimental setup and our evaluation methodology as follows.

##### A. Experimental Setup

A CUDA-enabled NVIDIA Quadro RTX 5000 GPU with 16 GB memory is employed in the experiments. Other hardware configurations are: Intel(R) Xeon(R) Gold 6240R CPU @ 2.40 GHz and 128 GB RAM. The PyTorch [101] DL framework is used for training, validation, and testing on a Windows 10 operating system. In both of the image sets, a simple data preprocessing scheme is employed by dividing images into 224-by-224 image patches. The total number of image patches belonging to each image set can be seen in Table I.

The DSTL image set is split into training, validation, and test sets in which 20% of the data is reserved as a test set to perform cross-validation. The remaining 80% of the data is further partitioned into 72% training and 8% validation sets so that hyperparameter tuning can also be implemented. Five-fold cross-validation is applied to obtain a less biased estimate of the benchmarked methods. In order to improve segmentation performance, manual hyperparameter tuning is applied. All ED architectures are trained using the adaptive moment estimation (Adam) [102] algorithm. For the DSTL image set, the initial learning rate is chosen as  $10^{-4}$  for U-Net and DFANet, while the initial value of  $5 \times 10^{-5}$  is used for the rest of the architectures, SegNet, DLinkNet, and DeepLabv3+. For DLinkNet architecture, the learning rate is reduced by 9% in every ten iteration steps, while for the others it is reduced by the same amount in every five iteration steps. Except for the SegNet, the mini-batch size is set to 8 and Xavier uniform is chosen for initialization. A Mini-batch size of 4 is chosen for SegNet. By observing the learning curves, 70 epochs were found to be sufficient for convergence.

Due to the fact that test image labelling is not provided for the RIT-18 aerial imagery, it is possible to use the already partitioned training and validation sets [73], where the validation

TABLE II  
RESULTS FOR THE DSTL IMAGE SET USING THE “THREE-CHANNEL INPUT” GROUP

	Jaccard Index					
	RGB	NIRi	NDVI	ARVI	SAVI	NDVI+ARVI+SAVI
U-Net	<b>0.570 ± 0.206</b>	0.493 ± 0.228	0.462 ± 0.253	0.493 ± 0.235	0.453 ± 0.259	0.524 ± 0.231
SegNet	0.492 ± 0.247	0.429 ± 0.259	0.415 ± 0.271	0.458 ± 0.246	0.416 ± 0.270	0.462 ± 0.246
DLinkNet (ResNet-34)	<b>0.571 ± 0.211</b>	0.526 ± 0.223	0.442 ± 0.262	0.485 ± 0.253	0.456 ± 0.257	<b>0.561 ± 0.226</b>
DeepLabv3+ (ResNet-34)	0.517 ± 0.212	0.504 ± 0.207	0.427 ± 0.244	0.435 ± 0.239	0.425 ± 0.246	0.505 ± 0.212
DeepLabv3+ (ResNet-101)	0.457 ± 0.217	0.466 ± 0.222	0.431 ± 0.252	0.440 ± 0.240	0.424 ± 0.244	0.475 ± 0.222
DeepLabv3+ (Xception)	0.481 ± 0.218	0.428 ± 0.232	0.313 ± 0.295	0.379 ± 0.248	0.309 ± 0.294	0.443 ± 0.244
DFANet	0.355 ± 0.239	0.350 ± 0.236	0.299 ± 0.264	0.309 ± 0.247	0.244 ± 0.263	0.313 ± 0.244
RF	0.371 ± 0.183	0.300 ± 0.177	0.219 ± 0.176	0.232 ± 0.126	0.206 ± 0.172	0.316 ± 0.180

Bold entities highlight the three best values among all results.

set is utilized as the test set in our experiments. For comparison, the same DSTL hyperparameters are also tuned for RIT-18. In order to handle the “black” (i.e., no data) regions in the image set (which are formed after image rectification), we did not utilize patches with more than 50% black region for training or testing. Thus, the training set is left with 814 image patches, while the test set contains 964. In the case of the RIT-18 image set, the same initial learning rate of  $10^{-4}$  is used for SegNet, DLinkNet, and DFANet architectures, while  $5 \times 10^{-5}$  is utilized for U-Net and DeepLabv3+. SegNet is the only model with its learning rate reduced by 9% in every 5 iteration steps, while other architectures have the same drop rate in every 10 iteration steps. The mini-batch size takes the value of 4 for SegNet and the DSTL image sets, while it is 8 for the rest. Xavier uniform is used as the initialization method and 100 epochs are evaluated with respect to convergence.

1) *Evaluation Metric*: Intersection over Union (IoU or Jaccard Index) considers false alarms and missed values simultaneously by counting the total number of mislabeled pixels. The Jaccard Index is defined as

$$\text{Jaccard Index} = \frac{\text{TP}}{\text{TP} + \text{FP} + \text{FN}} \quad (7)$$

where TP denotes the true positive pixels, FP denotes the false positive pixels, and FN denotes the false-negative pixels. Jaccard Index is introduced as the evaluation metric since it is becoming the de-facto standard metric for semantic segmentation evaluation tasks in the literature. Since 2008, The Jaccard Index has been used in the PASCAL VOC challenge as the fundamental evaluation standard [103]. The average Jaccard Index (mJI) for all test images is provided in our results.

## B. Results

Table II shows average Jaccard index values that are obtained from five-fold cross-validation results for the DSTL image set. As seen from the table, the best performance achieved is 0.571 (with a standard deviation of  $\pm 0.211$  among all test images) by applying DLinkNet with the ResNet34 backbone to RGB images. The U-Net architecture applied to RGB images yields very close performance to that of DLinkNet, with 0.570 mJI. The closest performance of a non-RGB input type is again by the DLinkNet+ResNet34 architecture using the mixed VIs, which yields an mJI of 0.561. We believe that the strength of the RGB input type on the DSTL image set is mainly because of the enhanced spatial resolution of the World-View-3’s RGB

channels. Although obtained from the original (i.e., unenhanced, lower) spatial resolution of the World-View-3 sensor, the mixed VIs input type shows a competitive performance. These results show that for this sensor (i.e., with this spatial resolution and altitude) when fused within a three-channel signal, the hand-crafted VIs provide valuable low-level features to the DL architectures and, hence, perform as well as the enhanced RGB input type. We should note that the resolution enhancement postprocess of RGB images from the World-View-3 sensor is an additional computation, which is carried out offline and may not always be supported for edge systems like satellites or UAVs.

On the other hand, even though the NIR reflectance information is quite important for assessing and distinguishing trees, the NIRi input type shows a low performance for all segmentation models for the DSTL image set. We believe that this is mainly due to two reasons: first, the aforementioned lower spatial resolution of the original NIR bands in the DSTL image set; and second, the coarse spectral resolution of the NIR subbands of the World-View-3 sensor, compared to the RIT-18’s performance with finer NIR bands as discussed in the following.

The segmentation performance results belonging to the RIT-18 image set are presented in Table III. There is a clear trend in NIRi results, with the highest values for most of the architectures. Judging by these superior mJI values, NIR reflectance shows a promising performance for tree segmentation. The peak value of 0.921 mJI is observed in the NIR reflectance band when DLinkNet with Resnet-34 encoder is applied. NIR wavelengths also work well for the U-Net architecture, with 0.885 mJI, giving evidence of the contribution from both high spatial (GSD 0.047 m) and high spectral resolution (3 channels between 485–685 nm).

To this extent, the segmentation performance of DLinkNet (Resnet-34) in the RIT-18 image set is consistent with those obtained in the DSTL image set. DLinkNet pretrained with ResNet-34 reliably outperforms other ED architectures and is better able to preserve detailed spatial information. The combination of the VIs (NDVI + SAVI + ARVI) does not have a significant impact on the performance if an image set already has a very high spatial and spectral resolution for NIR reflectance; thus, the NIR band alone can be a better choice for aerial remote sensing.

Another significant difference between the results of the DSTL and the RIT-18 image set experiments is the performance of the ARVI input type. Although performing comparatively well on the DSTL dataset, ARVI shows the lowest mJI values



TABLE III  
RESULTS FOR THE RIT-18 IMAGE SET USING THE “THREE-CHANNEL INPUT” GROUP

	Jaccard Index					
	RGB	NIR <sub>i</sub>	NDVI	ARVI	SAVI	NDVI+ARVI+SAVI
U-Net	0.860 ± 0.285	<b>0.885 ± 0.263</b>	0.841 ± 0.306	0.758 ± 0.371	0.827 ± 0.315	0.715 ± 0.293
SegNet	0.852 ± 0.293	0.829 ± 0.330	0.835 ± 0.310	0.753 ± 0.340	0.828 ± 0.330	0.513 ± 0.327
DLinkNet (ResNet-34)	<b>0.893 ± 0.233</b>	<b>0.921 ± 0.208</b>	0.861 ± 0.289	0.748 ± 0.368	0.842 ± 0.317	0.803 ± 0.360
DeepLabv3+ (ResNet-34)	0.841 ± 0.311	0.872 ± 0.284	0.810 ± 0.356	0.684 ± 0.465	0.820 ± 0.336	0.682 ± 0.464
DeepLabv3+ (ResNet-101)	0.852 ± 0.294	0.869 ± 0.296	0.787 ± 0.373	0.685 ± 0.464	0.811 ± 0.348	0.684 ± 0.465
DeepLabv3+ (Xception)	0.858 ± 0.277	0.836 ± 0.308	0.815 ± 0.345	0.684 ± 0.465	0.811 ± 0.346	0.681 ± 0.465
DFANet	0.807 ± 0.342	0.837 ± 0.306	0.796 ± 0.322	0.710 ± 0.461	0.821 ± 0.314	0.707 ± 0.408
RF	0.852 ± 0.237	0.866 ± 0.236	0.775 ± 0.288	0.732 ± 0.313	0.773 ± 0.287	0.821 ± 0.258

Bold entities highlight the three best values among all results.

TABLE IV  
RESULTS FOR BOTH SETS USING THE “MULTISPECTRAL INPUT” GROUP

	Jaccard Index			
	DSTL			RIT-18
	VNIR (5+3)	VNIR+SWIR <sup>2</sup> (5+3+2)	VNIR+SWIR <sup>8</sup> (5+3+8)	VNIR (3+3)
U-Net	0.546 ± 0.231	0.508 ± 0.237	0.538 ± 0.220	0.797 ± 0.351
SegNet	0.397 ± 0.268	0.427 ± 0.262	0.423 ± 0.263	0.841 ± 0.306
RF	0.360 ± 0.184	0.358 ± 0.185	0.352 ± 0.185	0.862 ± 0.236

for the RIT-18 dataset experiments. As explained in the previous sections, atmospheric correction requires information from the SWIR band. Hence, the lower performance of ARVI on the RIT-18 image set is due to the lack of SWIR spectral bands in the image set’s sensory output. When calculated properly using the required spectral bands, we observe a clearly improved performance in the hand-crafted VIs for tree segmentation.

The results of the experiments that utilize multispectral input types are provided in Table IV. Even though there are a number of band fusion modules applied by slightly modifying the networks [104], [105], the motivation in this study is to implement only the original architectures so only U-Net, SegNet, and RF are used. Although tested for a limited number of architectures, the results are consistent with three-channel input type experiments. The lowest performance is obtained with the VNIR+SWIR<sup>8</sup> input type for the DSTL image set. We see that the effect of increasing the number of input channels by itself is not significant to improve the segmentation performance; however, the optimum combination of these input channels is significant [106], [107], [108]. The Literature already shows that because of an insufficient amount of labeled remote-sensing data, DL models are usually difficult to be trained in an end-to-end fashion using multispectral input [25]. It is one of our fundamental research questions in this article to address the need for using VIs as hand-crafted features, instead of end-to-end training. In [109], it has already been shown that adding the input of extra-bands causes an inference to the network parameter learning and this, in turn, causes a degradation of the network performance. We believe that even if the input of all band data can increase the global information, those extra bands may be insufficient to identify detailed small tree objects in the scene [110].

Another point that cannot be ignored is most of these multiple spectral bands, such as NIR and SWIR in the DSTL image set have low spatial resolutions compared to RGB, leading to a degradation in the overall segmentation performance [106],

[111]. This is another indicator that supports the idea of using VIs as hand-crafted features when the spatial or spectral resolution of the visible and/or NIR bands is not sufficiently fine and the scale of the labelled image set is limited.

As obvious from Table III, the Jaccard index values of the RF algorithm are almost comparable to those of the ED architectures for the RIT-18 image set while this is not the case in Table II, i.e., for the DSTL image set. As already shown in [112], RF can considerably improve the performance by increasing the number of spectral bands in high spatial resolution images. Moreover, the RF algorithm not only offers the significant performance for dealing with multidimensional complex data [113], [114], but also requires only slight parameter tuning [115]. Therefore, RF is more likely to be robust to the performance degradation when multiple bands are used. This is why the best result of the multispectral input in Table IV belongs to the RF for the RIT-18 image set, where the spatial resolution is high.

The tree segmentation problem has many similarities to that of road segmentation in cases where multispectral remote sensing imagery is utilized. Occlusion is one of them, in which the object is partially or completely covered [116]. In the case of remote sensing imagery, the need for large receptive fields is crucial, since the input images are high resolution. The benefit of preserving detailed spatial information is additionally important, thus making the segmentation capable of acquiring complex geometries with highly discriminative feature information. We believe that these similarities, related to occlusions, high-resolution, and geometrical complexity, are the main reasons why a state-of-the-art road segmentation network, such as DLinkNet is able to provide excellent results for tree segmentation as well.

The similarities of our results can further be generalized to other road segmentation studies [80], [81], [117] as well. We observe that the performance of the DeepLabv3+ architecture employed with any of the state-of-the-art encoder backbones

	RGB	NIRi	NDVI	ARVI	SAVI	Mixed-Veg.	VNIR+SWIR <sup>2</sup>
Input							
U-Net							
SegNet							
DlinkNet (ResNet34)							
DeepLab v3+ (ResNet-34)							
DeepLab v3+ (ResNet-101)							
DeepLab v3+ (Xception)							
DFANet							
Random Forest							

Fig. 3. Semantic segmentation results for all implemented models and utilized input types, depicted on a selected  $224 \times 224$  pixels resolution region (approx.  $70 \text{ m} \times 70 \text{ m}$ ) from the DSTL Image Set. The first row consists of different input images; the rest are the semantic segmentation results. On semantic segmentation results, light green regions are the hits, dark green regions are the misses, and pink regions are the false alarms.

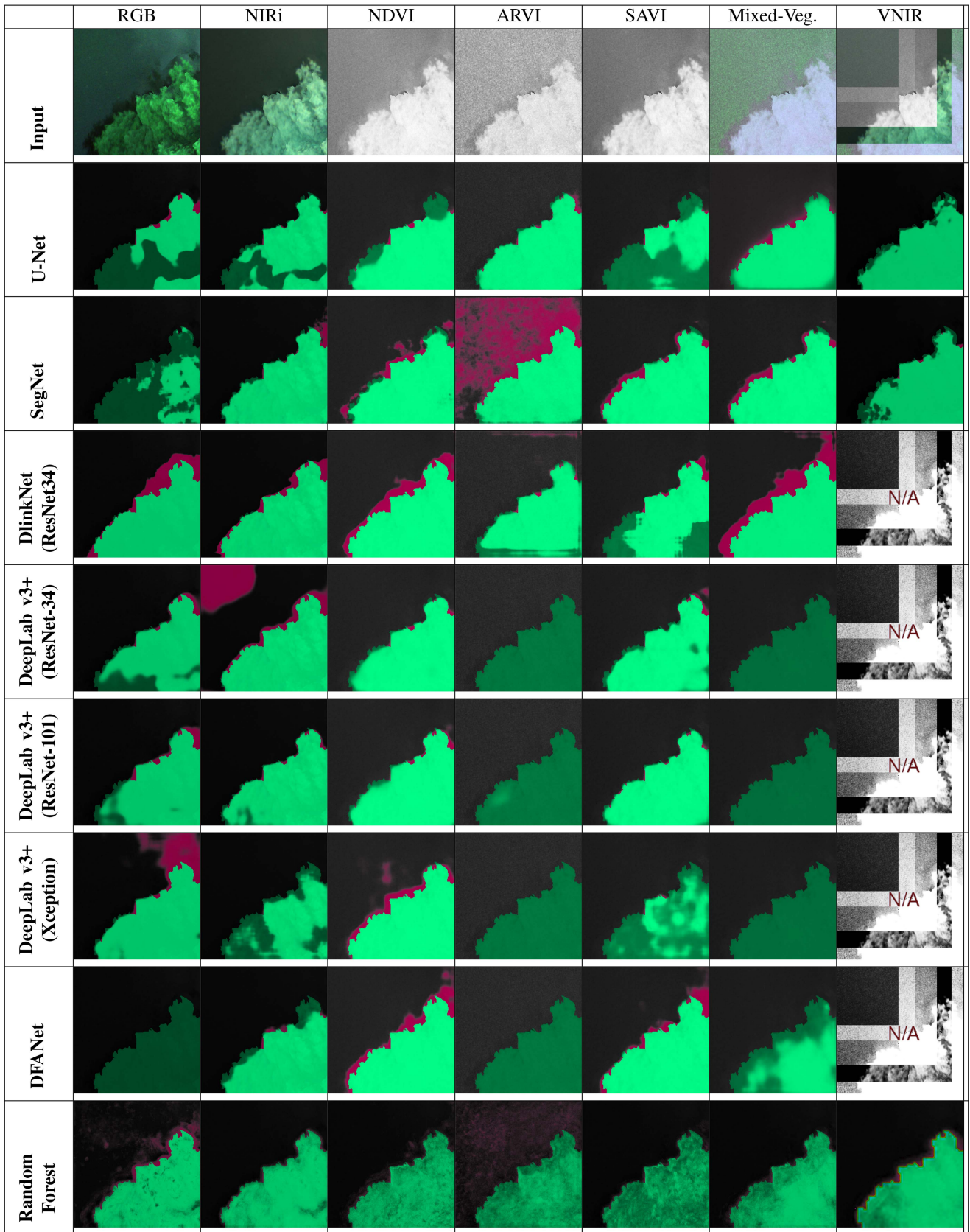


Fig. 4. Semantic segmentation results for all implemented models and utilized input types, depicted on a selected  $224 \times 224$  pixels resolution region (appr.  $11 \text{ m} \times 11 \text{ m}$ ) from the RIT-18 image set. The first row consists of different input images; the rest are the semantic segmentation results, light green regions are the hits, dark green regions are the misses, and pink regions are the false alarms.

is falling short of effective occlusion handling. DeepLabv3+ models are highly optimized to the particular problem definition and so it is unsurprising that they show performance degradation in diverse remote sensing imagery. Hence, DeepLabv3+ cannot handle the occlusion problem and may perform below its potential due to the high demand for training samples [81].

Since our selection criterion in the ED architectures is mainly related to the decoder design, segmentation performance results should also be appraised with respect to the types of upsampling layers. U-Net and DLinkNet are expected to provide broadly similar results, as they both use the same upsampling layer of transposed convolution. Segmentation performance results confirm that the architectures designed with transposed convolution, i.e., U-Net and DLinkNet, achieve superior performance over other ED architectures. This behaviour is persistent regardless of the image sets and no matter what kind of input is fed to these networks. DLinkNet achieves a comparable and even better performance than U-Net by exploiting dilated convolution to expand the receptive field and by using pretrained ResNet-34 as its encoder. The reason why SegNet fails to achieve the superior performance is that the recovery of high-resolution spatial details using “index-based” skip connections is not sufficiently effective [81]. In DeepLabv3+, the implementation of bilinear interpolation upsampling may lead to a smooth segmentation so that the decoder is unable to recover pixel-wise tree predictions accurately [118], [119]. This is mainly due to the fact that the decoder in DeepLabv3+ is not suitable for processing high-resolution remote sensing imagery, where the tree pixels are small compared to the overall image [120].

The performance of the DFANet architecture should be further examined against that of the other architectures, since it is specifically tailored for real-time semantic segmentation. Just like the DeepLabv3+, DFANet also employs bilinear upsampling operation in the decoder part such that both have similar segmentation performances. However, as shown in Table II, DFANet Jaccard index values are even lower than those of the DeepLabv3+ due to the low spatial resolution of the DSTL image set. As a result of its lack of network depth and reduced number of parameters, DFANet cannot capture high-dimensional features like U-Net [121]. The DFANet experimental results are notable for their relatively high Jaccard index values for NDVI+SAMI+ARVI combination, which is consistent with our results obtained from other DL architectures. Moreover, the experimental results in Table III show that the segmentation performance of DFANet architecture is almost close to that of U-Net architecture, since the RIT-18 image set has high spatial resolution. We believe that, for the DFANet to be able to exhibit better semantic segmentation performance in high-resolution remote sensing imagery input, some improvements on the DFANet architecture should be done which allow the model to better fit to this kind of data [122]. Obviously, such improvements will also affect its real-time nature.

Although it is highly preferable to draw conclusions from averaged results, such as those from Tables II–IV, it is worth noting that the single image results introduced in Figs. 3 and 4 depict explicit behavior in which ARVI introduces significant false alarms. The pixels marked with pink are the false alarms,

which can be seen to be quite dense in ARVI compared to other VIs. ARVI becomes sensitive to all kinds of green vegetation other than trees while reducing atmospheric scattering effects and, hence, presents false alarms by brightening the pixels belonging to other types of green vegetation, including small bushes and grasslands.

Tables II–IV show that the RF algorithm achieves comparable segmentation performances, even though the results are not as good as those from DL-based ED architectures. Even the best of the RF algorithm results (RGB, with 0.371 mJI), shown in Table II for the DSTL image set, is still lower than all the DL-based results. On the other hand, the achieved results shown in Table III for the RIT-18 are quite high, with the highest being 0.866 mJI for NIR reflectance, indicating how well the RF algorithm can perform when spatial resolution is sufficiently high.

For all input types and utilized methods, results exhibit high standard deviations of Jaccard index values. This leads us to a conclusion that individual maximum or minimum performances of all input types and/or methods are statistically close to each other. We believe that the reason behind this fact is due to the tested images being significantly diverse in nature. In RIT-18 image set, some images cover only a car park, whereas some images only cover vegetation. In DSTL, some images cover only buildings, whereas some others cover only traffic and roads. Due to this extreme diversity, we believe that the high deviation in the individual results are justifiable. However, the deviation in Jaccard values is calculated based on the results that are obtained from thousands of images. The average Jaccard values (mJIs) obtained from this dataset are considerably different (reaching 0.3 mJI of difference between some input types and methods). We believe that, although there is a deviation of success related to the diverse nature of individual images, the mJI values obtained from thousands of samples are a true indicator of the capabilities of different input types and/or segmentation methods.

## V. CONCLUSION

In this article, we study different semantic segmentation models and different multispectral input combinations for the pixel-wise tree segmentation problem on remote sensing imagery. For this purpose, we utilize a set of comparative experiments where we benchmark the selected models and the given input types. The essence of the comparison for the segmentation models lies in the selection of ED architectures for tree segmentation. We specifically analyze the descriptive power of these models to overcome issues, such as tree occlusion and geometric complexity. Different decoder designs, which make use of high-resolution information of remote sensing imagery and enable ED architectures to handle these problems, are explored. The comparative results procure significant conclusions, the most important of which are summarized as follows.

- 1) The spatial resolution of the remote sensing imagery is the most important factor and has a greater influence on segmentation performance than spectral resolution and applied architectures. It is necessary to utilize a very high-resolution remote sensing image set: otherwise, it

is not possible to reach considerable segmentation accuracy, even when a powerful, DL-based ED architecture is implemented. Therefore, an aerial image set is a better choice than satellite imagery, which is not only insufficient in terms of spatial resolution but also expensive.

- 2) The second most important factor is the spectral resolution of the NIR reflectance, for increasing the spectral discrimination between trees and other green vegetation. Although it is known that NIR reflectance information is significantly valuable for discriminating trees, the experiments demonstrate that if the spectral resolution of the NIR band is not sufficiently fine, sufficiently accurate tree segmentation cannot be obtained.
- 3) DLinkNet architecture consistently outperforms the compared semantic segmentation models, mainly for the cases including tree occlusion and grassland. Thanks to the decoder design with transposed convolution layers, further improved with dilated convolution and pretrained ResNet-34, the large receptive field of the DLinkNet is best suited to high-resolution remote sensing images.
- 4) It can be advantageous to use a combination of VIs in cases where the spatial resolution of the remote sensing image set is insufficient, such as the case we observe in our experiments with the DSTL image set. In Table II, we see that RGB and mixed-Veg input types perform similarly because the spatial resolution is not sufficiently high. This is not the case in Table III. Thus, we conclude that the following three factors should be taken into account when VIs are to be used for semantic segmentation of trees:
  - a) the input signal should include a combination of VIs of diverse characters;
  - b) the spatial resolution is relatively lower as in the case of satellite imagery;
  - c) a DL-based model architecture with sufficiently descriptive strength should be utilized.
- 5) To assess trees in the presence of the most common types of aerosols, such as smoke or sulfates, it is especially essential to utilize the SWIR band while calculating VIs that reduce atmospheric effects, such as ARVI. Due to the sensitivity of the SWIR band to the liquid water content of the tree leaves, ARVI calculated with the SWIR band shows a better segmentation performance than the one that does not utilize SWIR reflectance.

Regarding the future directions for this study, our first plan is to focus on improving the backbone architecture. The pretrained models, such as the ResNet are trained on RGB data, and they can potentially downgrade the performance when fine-tuning on small remote sensing datasets due to the large domain gap. Because, as we saw in our preliminary experiments, it is not possible to obtain convergence for training, when we train some deep architectures (such as DeepLabv3+ or DLinkNet) from scratch with a limited scale of data. Hence, the only ideal option can be to use encoders pretrained with satellite or aerial imagery with the same sensors, or in other words, the same multispectral bands trained with the same architectures. Such pretrained encoders, to the best of our knowledge, do not publicly exist. However, there are efforts to implement domain adaptation to

multispectral images obtained from different sensors, so that larger-scale experiments can be carried out [123]. Following this direction, a domain (i.e., sensor) independent multispectral backbone can be obtained.

As an alternative, converting the RGB backbones to a compatible model that can process multispectral data is also a viable direction. Additive group normalization method discussed in [104] is a recent example of this approach. By analyzing which bands contribute more to segmentation, they attempt to obtain DL-based VI designs that will be compatible with three-channel backbones, such as the ResNet. Following this path, we may be able to utilize multispectral band input to advanced architectures, such as DLinkNet or DeepLabv3+, and achieve better performance.

Another promising direction, not only for this particular problem at hand but in DL is utilizing vision transformers (ViT) instead of convolution-based architectures. Following the very recent ViT design [124], computer vision problems are being attacked with transformer-based architectures. As of the time this article is prepared, many transformed-based multispectral applications including segmentation [125], and many others [126], [127], are proposed, and will likely increase in near future.

## REFERENCES

- [1] T. M. Lenton et al., "Tipping elements in the earth's climate system," *Proc. Nat. Acad. Sci. USA*, vol. 105, no. 6, pp. 1786–1793, 2008. [Online]. Available: <https://www.pnas.org/content/105/6/1786>
- [2] L. Zhang, Z. Shao, J. Liu, and Q. Cheng, "Deep learning based retrieval of forest aboveground biomass from combined LiDAR and landsat 8 data," *Remote Sens.*, vol. 11, no. 12, 2019, Art. no. 1459.
- [3] D. J. Mulla, "Twenty five years of remote sensing in precision agriculture: Key advances and remaining knowledge gaps," *Biosyst. Eng.*, vol. 114, no. 4, pp. 358–371, 2013.
- [4] S. Barati, B. Rayegani, M. Saati, A. Sharifi, and M. Nasri, "Comparison the accuracies of different spectral indices for estimation of vegetation cover fraction in sparse vegetated areas," *Egyptian J. Remote Sens. Space Sci.*, vol. 14, no. 1, pp. 49–56, 2011. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1110982311000147>
- [5] K. Herndon, R. Muench, E. Cherrington, and R. Griffin, "An assessment of surface water detection methods for water resource management in the Nigerien Sahel," *Sensors*, vol. 20, no. 2, 2020, Art. no. 431. [Online]. Available: <https://www.mdpi.com/1424-8220/20/2/431>
- [6] A. A. Alkhatib, "A review on forest fire detection techniques," *Int. J. Distrib. Sensor Netw.*, vol. 10, no. 3, 2014, Art. no. 597368. [Online]. Available: <https://doi.org/10.1155/2014/597368>
- [7] O. Ghorbanzadeh, K. Gholamnia, and P. Ghamisi, "The application of ResU-net and OBIA for landslide detection from multi-temporal sentinel-2 images," *Big Earth Data*, vol. 6, pp. 1–26, 2022.
- [8] X. Yuan, J. Shi, and L. Gu, "A review of deep learning methods for semantic segmentation of remote sensing imagery," *Expert Syst. Appl.*, vol. 169, 2021, Art. no. 114417. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0957417420310836>
- [9] G. Cheng and J. Han, "A survey on object detection in optical remote sensing images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 117, pp. 11–28, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0924271616300144>
- [10] H. Mayer, "Automatic object extraction from aerial imagery—A survey focusing on buildings," *Comput. Vis. Image Understanding*, vol. 74, no. 2, pp. 138–149, 1999. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1077314299907506>
- [11] J. Mena, "State of the art on automatic road extraction for GIS update: A novel classification," *Pattern Recognit. Lett.*, vol. 24, no. 16, pp. 3037–3058, 2003. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167865503001648>
- [12] F. A. Gougeon, "A crown-following approach to the automatic delineation of individual tree crowns in high spatial resolution aerial images," *Can. J. Remote Sens.*, vol. 21, no. 3, pp. 274–284, 1995.

- [13] T. Brandtberg and F. Walter, "Automated delineation of individual tree crowns in high spatial resolution aerial images by multiple-scale analysis," *Mach. Vis. Appl.*, vol. 11, no. 2, pp. 64–73, 1998.
- [14] P. Meyera, K. Staenzb, and K. Ittena, "Semi-automated procedures for tree species identification in high spatial resolution data from digitized colour infrared-aerial photography," *ISPRS J. Photogrammetry Remote Sens.*, vol. 51, no. 1, pp. 5–16, 1996.
- [15] J. Wu, W. Yao, and P. Polewski, "Mapping individual tree species and vitality along urban road corridors with LiDAR and imaging sensors: Point density versus view perspective," *Remote Sens.*, vol. 10, no. 9, 2018, Art. no. 1403.
- [16] M. Immitzer, C. Atzberger, and T. Koukal, "Tree species classification with random forest using very high spatial resolution 8-band worldview-2 satellite data," *Remote Sens.*, vol. 4, no. 9, pp. 2661–2693, 2012.
- [17] V. Plakman, T. Janssen, N. Brouwer, and S. Veraverbeke, "Mapping species at an individual-tree scale in a temperate forest, using sentinel-2 images, airborne laser scanning data, and random forest classification," *Remote Sens.*, vol. 12, no. 22, 2020, Art. no. 3710.
- [18] L. Soleimannejad, S. Ullah, R. Abedi, M. Dees, and B. Koch, "Evaluating the potential of sentinel-2, landsat-8, and IRS satellite images in tree species classification of hyrcanian forest of Iran using random forest," *J. Sustain. Forestry*, vol. 38, no. 7, pp. 615–628, 2019.
- [19] C. Sothe et al., "Comparative performance of convolutional neural network, weighted and conventional support vector machine and random forest for classifying tree species using hyperspectral and photogrammetric data," *GIScience Remote Sens.*, vol. 57, no. 3, pp. 369–394, 2020.
- [20] R. Al-Ruzouq et al., "Image segmentation parameter selection and ant colony optimization for date palm tree detection and mapping from very-high-spatial-resolution aerial imagery," *Remote Sens.*, vol. 10, no. 9, 2018, Art. no. 1413.
- [21] M. S. Colgan, C. A. Baldeck, J.-B. Féret, and G. P. Asner, "Mapping savanna tree species at ecosystem scales using support vector machine classification and BRDF correction on airborne hyperspectral and LiDAR data," *Remote Sens.*, vol. 4, no. 11, pp. 3462–3480, 2012.
- [22] K. Liyanage and B. M. Whitaker, "Satellite image classification using LC-KSVD sparse coding," in *Proc. Intermountain Eng. Technol. Comput.*, 2020, pp. 1–6.
- [23] L. Jin, S. Gao, Z. Li, and J. Tang, "Hand-crafted features or machine learnt features? Together they improve RGB-D object recognition," in *Proc. IEEE Int. Symp. Multimedia*, 2014, pp. 311–319.
- [24] L. Zhang, L. Zhang, and B. Du, "Deep learning for remote sensing data: A technical tutorial on the state of the art," *IEEE Geosci. Remote Sens. Mag.*, vol. 4, no. 2, pp. 22–40, Jun. 2016.
- [25] X. X. Zhu et al., "Deep learning in remote sensing: A comprehensive review and list of resources," *IEEE Geosci. Remote Sens. Mag.*, vol. 5, no. 4, pp. 8–36, Dec. 2017.
- [26] M. Wu, C. Zhang, J. Liu, L. Zhou, and X. Li, "Towards accurate high resolution satellite image semantic segmentation," *IEEE Access*, vol. 7, pp. 55609–55619, 2019.
- [27] T. Hoese and C. Kuenzer, "Object detection and image segmentation with deep learning on earth observation data: A review-Part I: Evolution and recent trends," *Remote Sens.*, vol. 12, no. 10, 2020, Art. no. 1667.
- [28] T. Hoese, F. Bachofer, and C. Kuenzer, "Object detection and image segmentation with deep learning on earth observation data: A review-Part II: Applications," *Remote Sens.*, vol. 12, no. 18, 2020, Art. no. 3053.
- [29] T. Kattenborn, J. Leitloff, F. Schiefer, and S. Hinz, "Review on convolutional neural networks (CNN) in vegetation remote sensing," *ISPRS J. Photogrammetry Remote Sens.*, vol. 173, pp. 24–49, 2021.
- [30] Y. LeCun et al., "Backpropagation applied to handwritten zip code recognition," *Neural Comput.*, vol. 1, no. 4, pp. 541–551, 1989.
- [31] I. Ülkü and E. Akagündüz, "A survey on deep learning-based architectures for semantic segmentation on 2D images," *Appl. Artif. Intell.*, vol. 36, no. 1, pp. 1–45, 2022.
- [32] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 3431–3440.
- [33] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1520–1528.
- [34] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Image Comput. Comput. Assist. Interv.*, 2015, pp. 234–241.
- [35] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 12, pp. 2481–2495, Dec. 2017.
- [36] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 801–818.
- [37] L. Zhou, C. Zhang, and M. Wu, "D-LinkNet: LinkNet with pretrained encoder and dilated convolution for high resolution satellite imagery road extraction," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2018, pp. 192–1924.
- [38] H. Li, P. Xiong, H. Fan, and J. Sun, "DFANet: Deep feature aggregation for real-time semantic segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9514–9523.
- [39] I. Ulku, P. Barmpoutis, T. Stathaki, and E. Akagunduz, "Comparison of single channel indices for U-Net based segmentation of vegetation in satellite images," in *Proc. 12th Int. Conf. Mach. Vis.*, 2020, pp. 338–345. [Online]. Available: <https://doi.org/10.1117/12.2556374>
- [40] B. Baheti, S. Innani, S. Gajre, and S. Talbar, "Semantic scene segmentation in unstructured environment with modified DeepLabv3," *Pattern Recognit. Lett.*, vol. 138, pp. 223–229, 2020.
- [41] Y. Wang, B. Liang, M. Ding, and J. Li, "Dense semantic labeling with atrous spatial pyramid pooling and decoder for high-resolution remote sensing imagery," *Remote Sens.*, vol. 11, no. 1, 2019, Art. no. 20.
- [42] S. Du, S. Du, B. Liu, and X. Zhang, "Incorporating DeepLabv3 and object-based image analysis for semantic segmentation of very high resolution remote sensing images," *Int. J. Digit. Earth*, vol. 14, pp. 357–378, 2021.
- [43] S. Hartling, V. Sagan, P. Sidike, M. Maimaitijiang, and J. Carron, "Urban tree species classification using a WorldView-2/3 and LiDAR data fusion approach and deep learning," *Sensors*, vol. 19, no. 6, 2019, Art. no. 1284.
- [44] K. A. Korznikov, D. E. Kislov, J. Altman, J. Doležal, A. S. Vozmishcheva, and P. V. Krestov, "Using U-net-like deep convolutional neural networks for precise tree recognition in very high resolution RGB (red, green, blue) satellite images," *Forests*, vol. 12, no. 1, 2021, Art. no. 66.
- [45] L. Ding, H. Tang, and L. Bruzzone, "LANet: Local attention embedding to improve the semantic segmentation of remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 1, pp. 426–435, Jan. 2021.
- [46] S. Solberg, E. Naeset, and O. M. Bollandsas, "Single tree segmentation using airborne laser scanner data in a structurally heterogeneous spruce forest," *Photogrammetric Eng. Remote Sens.*, vol. 72, no. 12, pp. 1369–1378, 2006.
- [47] Q. Li, P. Yuan, X. Liu, and H. Zhou, "Street tree segmentation from mobile laser scanning data," *Int. J. Remote Sens.*, vol. 41, no. 18, pp. 7145–7162, 2020.
- [48] W. Yan, H. Guan, L. Cao, Y. Yu, C. Li, and J. Lu, "A self-adaptive mean shift tree-segmentation method using UAV LiDAR data," *Remote Sens.*, vol. 12, no. 3, 2020, Art. no. 515.
- [49] J. Yang, Z. Kang, S. Cheng, Z. Yang, and P. H. Akwensi, "An individual tree segmentation method based on watershed algorithm and three-dimensional spatial distribution analysis from airborne LiDAR point clouds," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 13, pp. 1055–1067, 2020.
- [50] C. Zhang, Y. Zhou, and F. Qiu, "Individual tree segmentation from LiDAR point clouds for urban forest inventory," *Remote Sens.*, vol. 7, no. 6, pp. 7892–7913, 2015.
- [51] H. Hamraz, M. A. Contreras, and J. Zhang, "A robust approach for tree segmentation in deciduous forests using small-footprint airborne LiDAR data," *Int. J. Appl. Earth Observ. Geoinformation*, vol. 52, pp. 532–541, 2016.
- [52] F. Pirotti, M. Kobal, and J. Roussel, "A comparison of tree segmentation methods using very high density airborne laser scanner data," *Int. Arch. Photogrammetry Remote Sens. Spatial Inf. Sci.*, vol. 42, pp. 285–290, 2017.
- [53] S. Dersch, M. Heurich, N. Krueger, and P. Krzystek, "Combining graph-cut clustering with object-based stem detection for tree segmentation in highly dense airborne LiDAR point clouds," *ISPRS J. Photogrammetry Remote Sens.*, vol. 172, pp. 207–222, 2021.
- [54] L. Comesaña-Cebral, J. Martínez-Sánchez, H. Lorenzo, and P. Arias, "Individual tree segmentation method based on mobile backpack LiDAR point clouds," *Sensors*, vol. 21, no. 18, 2021, Art. no. 6007.
- [55] K. Itakura, S. Miyatani, and F. Hosoi, "Estimating tree structural parameters via automatic tree segmentation from LiDAR point cloud data," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 15, pp. 555–564, 2022.
- [56] J. Hyypä, O. Kelle, M. Lehtikoinen, and M. Inkinen, "A segmentation-based method to retrieve stem volume estimates from 3-D tree height models produced by laser scanners," *IEEE Trans. Geosci. Remote Sens.*, vol. 39, no. 5, pp. 969–975, May 2001.

- [57] J. B. Roerdink and A. Meijster, "The watershed transform: Definitions, algorithms and parallelization strategies," *Fundamenta Informaticae*, vol. 41, no. 12, pp. 187–228, 2000.
- [58] E. Ayrey et al., "Layer stacking: A novel algorithm for individual forest tree segmentation from LiDAR point clouds," *Can. J. Remote Sens.*, vol. 43, no. 1, pp. 16–27, 2017.
- [59] C. Dechesne, C. Mallet, A. L. Bris, and V. Gouet-Brunet, "Semantic segmentation of forest stands of pure species combining airborne LiDAR data and very high resolution multispectral imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 126, pp. 129–145, 2017.
- [60] P. Krzystek, A. Serebryanyk, C. Schnörr, J. Červenka, and M. Heurich, "Large-scale mapping of tree species and dead trees in Šumava national park and Bavarian forest national park using LiDAR and multispectral imagery," *Remote Sens.*, vol. 12, no. 4, 2020, Art. no. 661.
- [61] D. Pulido, J. Salas, M. Rös, K. Puettmann, and S. Karaman, "Assessment of tree detection methods in multispectral aerial images," *Remote Sens.*, vol. 12, no. 15, 2020, Art. no. 2379.
- [62] S. T. Digumarti et al., "An approach for semantic segmentation of tree-like vegetation," in *Proc. Int. Conf. Robot. Autom.*, 2019, pp. 1801–1807.
- [63] M. Brandt et al., "An unexpectedly large count of trees in the West African Sahara and Sahel," *Nature*, vol. 587, no. 7832, pp. 78–82, 2020.
- [64] A. Bannari, D. Morin, F. Bonn, and A. Huete, "A review of vegetation indices," *Remote Sens. Rev.*, vol. 13, no. 1/2, pp. 95–120, 1995.
- [65] A. R. Huete, "Vegetation indices, remote sensing and forest monitoring," *Geography Compass*, vol. 6, no. 9, pp. 513–532, 2012.
- [66] D. R. Peddle, S. Brunke, and F. G. Hall, "A comparison of spectral mixture analysis and ten vegetation indices for estimating boreal forest biophysical information from airborne data," *Can. J. Remote Sens.*, vol. 27, no. 6, pp. 627–635, 2001.
- [67] S. P. Mohanty et al., "Deep learning for understanding satellite imagery: An experimental survey," *Front. Artif. Intell.*, vol. 3, 2020, Art. no. 85. [Online]. Available: <https://www.frontiersin.org/article/10.3389/fraci.2020.534696>
- [68] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [69] S. Madonsela et al., "Multi-phenology worldview-2 imagery improves remote sensing of savannah tree species," *Int. J. Appl. Earth Observ. Geoinformation*, vol. 58, pp. 65–73, 2017.
- [70] A. Darshik, A. Dev, M. Bharath, B. A. Nair, and G. Gopakumar, "Semantic segmentation of spectral images: A comparative study using FCN8s and U-NET on RIT-18 dataset," in *Proc. 11th Int. Conf. Comput. Commun. Netw. Technol.*, 2020, pp. 1–6.
- [71] H. Hao et al., "Improving building segmentation using uncertainty modeling and metadata injection," in *Proc. 29th Int. Conf. Adv. Geographic Inf. Syst.*, 2021, pp. 117–120.
- [72] A. Song and Y. Kim, "Semantic segmentation of remote-sensing imagery using heterogeneous big data: International society for photogrammetry and remote sensing potsdam and cityscape datasets," *ISPRS Int. J. Geo-Inf.*, vol. 9, no. 10, 2020, Art. no. 601.
- [73] R. Kemker, C. Salvaggio, and C. Kanan, "Algorithms for semantic segmentation of multispectral remote sensing imagery using deep learning," *ISPRS J. Photogrammetry Remote Sens.*, vol. 145, pp. 60–77, 2018.
- [74] L. Mou, Y. Hua, and X. X. Zhu, "A relation-augmented fully convolutional network for semantic segmentation in aerial scenes," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 12408–12417.
- [75] R. Kemker, C. Salvaggio, and C. Kanan, "High-resolution multispectral dataset for semantic segmentation," 2017, *arXiv:1703.01918*.
- [76] Mar. 2017. Accessed: Sep. 9, 2022. [Online]. Available: <https://www.kaggle.com/c/dstl-satellite-imageryfeature-detection>
- [77] Mar. 2002. Accessed: Sep. 9, 2022. [Online]. Available: <http://www.satimagingcorp.com/>
- [78] Z. Wojna et al., "The devil is in the decoder: Classification, regression and GANs," *Int. J. Comput. Vis.*, vol. 127, no. 11, pp. 1694–1706, 2019.
- [79] Oct. 2014. Accessed: Sep. 9, 2022. [Online]. Available: [http://host.robots.ox.ac.uk:8080/leaderboard/displaylb\\_main.php?challengeid=11&compid=6](http://host.robots.ox.ac.uk:8080/leaderboard/displaylb_main.php?challengeid=11&compid=6)
- [80] Z. Tian, T. He, C. Shen, and Y. Yan, "Decoders matter for semantic segmentation: Data-dependent decoding enables flexible feature aggregation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 3121–3130.
- [81] D. L. Torres et al., "Applying fully convolutional architectures for semantic segmentation of a single tree species in urban environment on high resolution UAV optical imagery," *Sensors*, vol. 20, no. 2, 2020, Art. no. 563.
- [82] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," 2017, *arXiv:1706.05587*.
- [83] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected CRFs," 2014, *arXiv:1412.7062*.
- [84] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2018.
- [85] T.-Y. Lin et al., "Microsoft COCO: Common objects in context," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 740–755.
- [86] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [87] C. Szegedy et al., "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 1–9.
- [88] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 1800–1807.
- [89] A. Chaurasia and E. Culurciello, "LinkNet: Exploiting encoder representations for efficient semantic segmentation," in *Proc. IEEE Vis. Commun. Image Process.*, 2017, pp. 1–4.
- [90] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, "ENet: A deep neural network architecture for real-time semantic segmentation," Jun. 2016, *arXiv:1606.02147*.
- [91] S. Mehta, M. Rastegari, A. Caspi, L. Shapiro, and H. Hajishirzi, "ESPNet: Efficient spatial pyramid of dilated convolutions for semantic segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 552–568.
- [92] H. Zhao, X. Qi, X. Shen, J. Shi, and J. Jia, "ICNet for real-time semantic segmentation on high-resolution images," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 405–420.
- [93] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang, "BiSeNet: Bilateral segmentation network for real-time semantic segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 325–341.
- [94] J. Su et al., "Aerial visual perception in smart farming: Field study of wheat yellow rust monitoring," *IEEE Trans. Ind. Informat.*, vol. 17, no. 3, pp. 2242–2249, Mar. 2021.
- [95] R. M. Rustowicz, R. Cheong, L. Wang, S. Ermon, M. Burke, and D. Lobell, "Semantic segmentation of crop type in Africa: A novel dataset and analysis of deep learning methods," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, 2019, pp. 75–82.
- [96] [Online]. Available: [https://earthobservatory.nasa.gov/features/MeasuringVegetation/measuring\\_vegetation\\_2.php](https://earthobservatory.nasa.gov/features/MeasuringVegetation/measuring_vegetation_2.php)
- [97] A. Azadeh, P. Dimitrios, and S. Peter, "Forest canopy density assessment using different approaches—review," *J. Forest Sci.*, vol. 63, no. 3, pp. 107–116, 2017.
- [98] J. Rouse et al., "Monitoring vegetation systems in the great plains with ERTS," *NASA Special Publ.*, vol. 351, no. 1974, 1974, Art. no. 309.
- [99] Y. J. Kaufman and D. Tanre, "Atmospherically resistant vegetation index (ARVI) for EOS-MODIS," *IEEE Trans. Geosci. Remote Sens.*, vol. 30, no. 2, pp. 261–270, Mar. 1992.
- [100] A. R. Huete, "A soil-adjusted vegetation index (SAVI)," *Remote Sens. Environ.*, vol. 25, no. 3, pp. 295–309, 1988.
- [101] A. Paszke et al., "PyTorch: An imperative style, high-performance deep learning library," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2019, pp. 8026–8037.
- [102] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.
- [103] G. Csurka, D. Larlus, F. Perronnin, and F. Meylan, "What is a good evaluation measure for semantic segmentation?," in *Proc. 24th Brit. Mach. Vis. Conf.*, Sep. 2013, pp. 32.1–32.11.
- [104] H. Sheng, X. Chen, J. Su, R. Rajagopal, and A. Ng, "Effective data fusion with generalized vegetation index: Evidence from land cover segmentation in agriculture," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops*, 2020, pp. 267–276.
- [105] K. Yuan, X. Zhuang, G. Schaefer, J. Feng, L. Guan, and H. Fang, "Deep-learning-based multispectral satellite image segmentation for water body detection," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 7422–7434, 2021.
- [106] M. A. E. Bhuiyan et al., "Understanding the effects of optimal combination of spectral bands on deep learning model predictions: A case study based on permafrost Tundra landform mapping using high resolution multispectral satellite imagery," *J. Imag.*, vol. 6, no. 9, 2020, Art. no. 97.

- [107] Y. Cai, H. Huang, K. Wang, C. Zhang, L. Fan, and F. Guo, "Selecting optimal combination of data channels for semantic segmentation in city information modelling (CIM)," *Remote Sens.*, vol. 13, no. 7, 2021, Art. no. 1367.
- [108] H. Grybas and R. G. Congalton, "A comparison of multi-temporal RGB and multispectral UAS imagery for tree species classification in heterogeneous new hampshire forests," *Remote Sens.*, vol. 13, no. 13, 2021, Art. no. 2631.
- [109] Z. Wang et al., "Semantic segmentation and analysis on sensitive parameters of forest fire smoke using Smoke-Unet and Landsat-8 imagery," *Remote Sens.*, vol. 14, no. 1, 2022, Art. no. 45.
- [110] W. Kang, Y. Xiang, F. Wang, and H. You, "DO-Net: Dual-output network for land cover classification from optical remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, 2022, Art. no. 8021205.
- [111] T. Yang et al., "Sea-land segmentation using deep learning techniques for Landsat-8 OLI imagery," *Mar. Geodesy*, vol. 43, no. 2, pp. 105–133, 2020.
- [112] X. Yu et al., "Examining the roles of spectral, spatial, and topographic features in improving land-cover and forest classifications in a subtropical region," *Remote Sens.*, vol. 12, no. 18, 2020, Art. no. 2907.
- [113] M. W. Ahmad, M. Mourshed, and Y. Rezgui, "Trees vs neurons: Comparison between random forest and ANN for high-resolution prediction of building energy consumption," *Energy Buildings*, vol. 147, pp. 77–89, 2017.
- [114] M. Sheykhou, M. Mahdianpari, H. Ghanbari, F. Mohammadimanes, P. Ghamisi, and S. Homayouni, "Support vector machine versus random forest for remote sensing image classification: A meta-analysis and systematic review," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 13, pp. 6308–6325, 2020.
- [115] J. Xia, P. Ghamisi, N. Yokoya, and A. Iwasaki, "Random forest ensembles and extended multiextinction profiles for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 1, pp. 202–216, Jan. 2018.
- [116] S. Wang, H. Yang, Q. Wu, Z. Zheng, Y. Wu, and J. Li, "An improved method for road extraction from high-resolution remote-sensing images that enhances boundary information," *Sensors*, vol. 20, no. 7, 2020, Art. no. 2064.
- [117] M. Zhou, H. Sui, S. Chen, J. Wang, and X. Chen, "BT-RoadNet: A boundary and topologically-aware neural network for road extraction from high-resolution remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 168, pp. 288–306, 2020.
- [118] B. Zhang et al., "Dynamic neural representational decoders for high-resolution semantic segmentation," in *Proc. 35th Conf. Neural Inf. Process. Syst.*, 2021, pp. 17388–17399.
- [119] C. Shen, L. Liu, L. Zhu, J. Kang, N. Wang, and L. Shao, "High-throughput in situ root image segmentation based on the improved DeepLabv3 method," *Front. Plant Sci.*, vol. 11, 2020, Art. no. 576791.
- [120] Y. Bao and Y. Zheng, "Based on the improved DeepLabv3 remote sensing image semantic segmentation algorithm," in *Proc. 4th Int. Conf. Adv. Electron. Mater. Comput. Softw. Eng.*, 2021, pp. 717–720.
- [121] Y. Zhang, H. Ren, W. Yang, Y. Wang, K. Ye, and C.-Z. Xu, "The strong substructure and feature attention mechanism for image semantic segmentation," *Concurrency Computation: Pract. Experience*, vol. 34, no. 12, 2022, Art. no. e5920.
- [122] Z. Wang, J. Guo, W. Huang, and S. Zhang, "High-resolution remote sensing image semantic segmentation based on a deep feature aggregation network," *Meas. Sci. Technol.*, vol. 32, no. 9, 2021, Art. no. 095002.
- [123] D. Lunga, H. L. Yang, A. Reith, J. Weaver, J. Yuan, and B. Bhaduri, "Domain-adapted convolutional networks for satellite image classification: A large-scale interactive learning workflow," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 11, no. 3, pp. 962–977, Mar. 2018.
- [124] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," 2020, *arXiv:2010.11929*.
- [125] Z. Xu, W. Zhang, T. Zhang, Z. Yang, and J. Li, "Efficient transformer for remote sensing image segmentation," *Remote Sens.*, vol. 13, no. 18, 2021, Art. no. 3585. [Online]. Available: <https://www.mdpi.com/2072-4292/13/18/3585>
- [126] Y. Bazi, L. Bashmal, M. Alrahhal, R. Dayil, and N. Ajlan, "Vision transformers for remote sensing image classification," *Remote Sens.*, vol. 13, Feb. 2021, Art. no. 516.
- [127] L. Bashmal, Y. Bazi, and M. A. Rahhal, "Deep vision transformers for remote sensing scene classification," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, 2021, pp. 2815–2818.



She was a Research Associate with Imperial College London, United Kingdom. She is currently an Assistant Professor with the Department of Computer Engineering, Ankara University, Ankara, Turkey. Her research interests include multispectral image processing and deep learning-based semantic segmentation.

**Irem Ulku** received the B.Sc. degrees both in electronics and communication engineering and in industrial engineering (as valedictorian) from Çankaya University, Ankara, Turkey, in 2009 and 2010, respectively, and the M.Sc. degree in electrical and electronics engineering from the Middle East Technical University, Ankara, Turkey, in 2013, and the Ph.D. degree in electronics and communication engineering from Çankaya University, Ankara, Turkey, in 2017.

She was an Instructor with Çankaya University between 2017 to 2019 and 2020 to 2021. In 2019,



**Erdem Akagündüz** (Member, IEEE) received the B.Sc., M.Sc., and Ph.D. degrees in electrical and electronics engineering from Middle East Technical University (METU), Ankara, Turkey, in 2001, 2004, and 2011, respectively.

He is currently an Associate Professor with the Graduate School of Informatics, METU. From 2001 to 2008, he was a Research Assistant with the METU Computer Vision and Intelligent Systems Laboratory, and from 2009 to 2016, he was a Computer Vision Scientist with ASELSAN Inc., Ankara, Turkey. He

then became a Research Associate with the University of York, York, U.K., in 2016. Before joining back to METU, he was an Assistant Professor with the Electrical and Electronics Engineering Department, Çankaya University. His research interests include infrared computer vision, object/target/scene recognition, and deep learning.



**Pedram Ghamisi** (Senior Member, IEEE) received the Ph.D. degree in electrical and computer engineering from the University of Iceland, Reykjavik, Iceland, in 2015.

He is currently the Head of the Machine Learning Group, Helmholtz-Zentrum Dresden-Rossendorf (HZDR), Dresden, Germany, and Research Professor and Group Leader of AI4RS, the Institute of Advanced Research in Artificial Intelligence (IARAI), Vienna, Austria. He is a Cofounder of VasoGnosis Inc. with two branches in San Jose, CA, USA, and

Milwaukee, WI, USA. He was the Co-Chair of the IEEE Image Analysis and Data Fusion Committee (IEEE IADF) between 2019 and 2021. His research interests include interdisciplinary research on machine (deep) learning, image and signal processing, and multisensor data fusion.

Dr. Ghamisi was a recipient of the IEEE Mikio Takagi Prize for winning the Student Paper Competition at IEEE International Geoscience and Remote Sensing Symposium (IGARSS) in 2013, the first prize of the data fusion contest organized by the IEEE IADF in 2017, the Best Reviewer Prize of IEEE GEOSCIENCE AND REMOTE SENSING LETTERS in 2017, and the IEEE Geoscience and Remote Sensing Society 2020 Highest-Impact Paper Award. He is currently an Associate Editor for IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING and IEEE GEOSCIENCE AND REMOTE SENSING LETTERS.