

# Review on Convolutional Neural Networks (CNN) in Vegetation Remote Sensing

This is an uncorrected preprint. The peer-reviewed manuscript is published in  
ISPRS Journal of Photogrammetry and Remote Sensing. The Web- and  
pdf-version can be found at: <https://authors.elsevier.com/a/1cQX33I9x1cgg0>.

Teja Kattenborn<sup>1,\*</sup>, Jens Leitloff<sup>2</sup>, Felix Schiefer<sup>3</sup>, and Stefan Hinz<sup>2</sup>

<sup>1</sup>Remote Sensing Centre for Earth System Research, Leipzig University, Talstr.  
35, 04103 Leipzig, Germany

<sup>2</sup>Institute of Photogrammetry and Remote Sensing (IPF), Karlsruher Institute of  
Technology (KIT), Englerstr. 7, 76131 Karlsruhe, Germany

<sup>3</sup>Institute for Geography and Geoecology (IFGG), Karlsruher Institute of  
Technology (KIT), Kaiserstr. 12, 76131 Karlsruhe, Germany

\*Corresponding author: Teja Kattenborn, [teja.kattenborn@uni-leipzig.de](mailto:teja.kattenborn@uni-leipzig.de)

## Abstract

Identifying and characterizing vascular plants in time and space is required in various disciplines, e.g. in forestry, conservation and agriculture. Remote sensing emerged as a key technology revealing both spatial and temporal vegetation patterns. Harnessing the ever growing streams of remote sensing data for the increasing demands on vegetation assessments and monitoring requires efficient, accurate and flexible methods for data analysis. In this respect, the use of deep learning methods is trend-setting, enabling high predictive accuracy, while learning the relevant data features independently in an end-to-end fashion. Very recently, a series of studies have demonstrated that the deep learning method of Convolutional Neural Networks (CNN) is very effective to represent spatial patterns enabling to extract a wide array of vegetation properties from remote sensing imagery. This review introduces the principles of CNN and distills why they are particularly suitable for vegetation remote sensing. The main part synthesizes current trends and developments, including considerations about spectral resolution, spatial grain, different sensors types, modes of reference data generation, sources of existing reference data, as well as CNN approaches and architectures. The literature review showed that CNN can be applied to various problems, including the detection of individual plants or the pixel-wise segmentation of vegetation classes, while numerous studies have evinced that CNN outperform shallow machine learning methods. Several studies suggest that the ability of CNN to exploit spatial patterns particularly facilitates the value of very high

spatial resolution data. The modularity in the common deep learning frameworks allows a high flexibility for the adaptation of architectures, whereby especially multi-modal or multi-temporal applications can benefit. An increasing availability of techniques for visualizing features learned by CNNs will not only contribute to interpret but to learn from such models and improve our understanding of remotely sensed signals of vegetation. Although CNN has not been around for long, it seems obvious that they will usher in a new era of vegetation remote sensing.

**Keywords**— Convolutional Neural Networks (CNN), Deep Learning, Vegetation, Plants, Remote Sensing, Earth Observation

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Principles of CNNs and relevance for vegetation remote sensing</b>	<b>5</b>
2.1	Basic functioning and structure of CNNs . . . . .	5
2.2	Why CNN for vegetation remote sensing? . . . . .	7
2.3	The training process . . . . .	8
2.4	Implementation, libraries and frameworks . . . . .	10
<b>3</b>	<b>Literature review on CNN-based vegetation remote sensing</b>	<b>11</b>
3.1	Reference data . . . . .	12
3.1.1	Reference Data Sources . . . . .	12
3.1.2	Reference data quantity . . . . .	14
3.2	Common CNN approaches and architectures . . . . .	16
3.2.1	Training strategies . . . . .	16
3.2.2	Approaches and architectures . . . . .	21
3.3	Geographic and thematic areas of CNN application . . . . .	25
3.4	Remote sensing platforms . . . . .	26
3.5	Sensors, spatial and spectral resolution . . . . .	27
3.5.1	Passive optical and SAR data analysis . . . . .	28
3.5.2	LiDAR-based point cloud analysis . . . . .	30
3.5.3	Sensor and data fusion . . . . .	31
3.5.4	Multi-temporal analysis . . . . .	33
3.6	CNN model assessment, understanding, and interpretation . . . . .	35
3.6.1	Numeric evaluation of the predictive performance . . . . .	35
3.6.2	Understanding and interpretation: Opening the <i>black box</i> . . . . .	36
<b>4</b>	<b>Concluding remarks and future perspectives</b>	<b>38</b>
<b>5</b>	<b>Additional resources on CNN theory, implementation and data sources</b>	<b>40</b>

# 1 Introduction

2 Locating and characterizing vascular plants in time and space is key to various  
3 tasks: For instance, nature conservation in the context of global change and biodi-  
4 versity decline can only be successfully implemented and supervised with accurate  
5 spatial representations of the state, structure and functioning of ecosystems and  
6 its flora (Nagendra et al. 2013; Pettorelli et al. 2017; Turner et al. 2003). Forestry  
7 requires regular and extensive information on forest stands, including their struc-  
8 ture, timber volume, species composition, and forest damage (Fassnacht et al. 2016;  
9 McRoberts et al. 2007; White et al. 2016). In agriculture, there is a growing demand  
10 for geoinformation that facilitates resource efficiency and a reduction of environ-  
11 mental impacts (cf. precision farming), including fine-scale predictions of yield,  
12 weed infestations, and plant vigor (Atzberger et al. 2013; Mulla 2013). Concerning  
13 all of these tasks and requirements remote sensing continuously establishes as a key  
14 technology.

15 In the last decades, various technological advances resulted in growing avail-  
16 ability of remote sensing data revealing vegetation patterns on both spatial and  
17 temporal domains (Colomina et al. 2014; Toth et al. 2016). Novel remote sensing  
18 platforms, such as swarms of microsatellites, or unmanned aerial vehicles (UAV),  
19 facilitate a bird’s eye view on vegetation canopies with increasing spatial detail.  
20 Synthetic-aperture radar (SAR), and terrestrial or airborne lasers-scanning enable  
21 to capture the three-dimensional structure of multilayered canopies. Additionally,  
22 there is an ongoing trend of data sharing and open access (cf. *OpenAerialMap*,  
23 *NEON* programme of the US National Science Foundation, EU’s and ESA’s *Copen-  
24 nicus Open Access Hub*).

25 These growing opportunities for vegetation remote sensing come hand in hand  
26 with several challenges, including increased data volumes and computational loads  
27 as well as more diverse data structures with increasing dimensions (spatial, tem-  
28 poral, spectral) often featuring complex relationships. Moreover, the various veg-  
29 etation related tasks and applications fields can differ greatly in their inherent  
30 processes and requirements. Hence, harnessing remote sensing data for vegeta-  
31 tion assessments and monitoring requires efficient, accurate, and flexible analytical  
32 methods.

33 In the context of image analysis and computer vision, deep learning is currently  
34 paving new avenues for remote sensing analysis (Chollet 2017; Hoeser et al. 2020;  
35 Huang et al. 2018; Ronneberger et al. 2015; L. Zhang et al. 2019; Zhu et al. 2017).  
36 In contrast to the previous **shallow** neural network approaches that have been  
37 under investigation for decades, **deep** learning is characterized by a significantly  
38 increased number of successively connected neural layers. This increased amount  
39 of layers and transformations can reveal higher-level features and more abstract  
40 concepts uncovering more complex and hierarchical relationships. A series of stud-  
41 ies has demonstrated that this increased depth can indeed enhance the retrieval of  
42 vegetation-related information contained in remote sensing data (cf. section 3.6).  
43 At the same time, increasing transformations and, thus, deeper levels of complex-  
44 ity commonly require more training data and computational loads. Nevertheless,  
45 deep learning became very popular due to several, corresponding technical devel-  
46 opments, including efficient data processing techniques (e.g. data augmentation or  
47 non-linear activation functions, see section 2.3 and 3.2), high-performance graphic  
48 cards, cloud-computing, as well as open data initiatives providing annotated data.

49 These developments enable an efficient calculation of countless non-linear trans-  
50 formations of the respective input data and, thus, form the core for the essential  
51 strength of deep learning - namely the ability of **end-to-end-learning**. Previous  
52 data analysis methods in remote sensing usually require feature engineering, which  
53 is the heuristic selection of appropriate transformations and hand-crafting latent  
54 variables from the input data prior to modelling. Examples in the field of vege-  
55 tation remote sensing are spectral indices (Adam et al. 2010) or texture metrics  
56 (Haralick 1979), whereas the numerous ways to derive such variables make it often  
57 impossible and inefficient to derive the most effective set of predictors. Moreover,  
58 defining the most appropriate predictors for vegetation analysis can be challenging,  
59 as this may not only require knowledge on the biochemical and structural plant  
60 properties but also on how these interact with the electromagnetic signal measured  
61 by the sensor. By contrast, with deep learning, the neural network itself can learn  
62 the data transformations that are best to solve the problem at hand.

63 The class of deep learning algorithms most commonly used for spatial pattern  
64 analysis are convolutional neural networks (CNNs or ConvNets). CNNs are de-  
65 signed to learn the spatial features, e.g. edges, corners, textures, or more abstract  
66 shapes, that best describe the target class or quantity. The core for learning these  
67 features are manifold and successive transformations of the input data (convolu-  
68 tions) on different spatial scales (by pooling operations). This facilitates identifying  
69 and combining both low-level features and high-level concepts. The functioning of  
70 a CNN can, hence, be regarded as a mimicry of the animal cortex (Angermueller  
71 et al. 2016; Cadieu et al. 2014), where analogously numerous visual stimuli at vary-  
72 ing scales are perceived in the field of vision (counterpart of an image) and the  
73 contained spatial features and their spatial context serves to identify objects. For  
74 example, the shape of a leaf does not necessarily indicate the corresponding vege-  
75 tation type, but its close proximity to branches and the tall and bulky canopy suggest  
76 that it belongs to a tree and not to a herb. The effectiveness of deep learning and  
77 particularly CNNs undoubtedly revolutionized our possibilities to analyse spatial  
78 patterns in Earth observation data. Reference is made here to previous and valu-  
79 able comments and reviews, including a review by Zhu et al. (2017) on the general  
80 principles and potentials of deep learning in remote sensing, Hoese et al. (2020)  
81 summarizing common frameworks and an in depth overview on architectures for  
82 Earth observation data analysis, a comment by Brodrick et al. (2019) highlighting  
83 potentials of CNN for segmentation tasks in ecology and Reichstein et al. (2019)  
84 providing perspectives on how deep learning in general can advance earth system  
85 science.

86 The remote sensing of vegetation is characterized by special requirements and  
87 challenges, such as the often complex acquisition of reference data or the under-  
88 standing of the vegetation specific radiative transfer, the resulting sensor-specific  
89 electromagnetic signals and their dynamics across the phenology. The present re-  
90 view therefore concentrates specifically on CNN applications in the field of vege-  
91 tation remote sensing. A series of recent studies have demonstrated that CNNs enable  
92 to reveal accurate spatial representations of vegetation properties, such as detecting  
93 individual plant organs or individuals, classifying species and communities or quan-  
94 tifying plant traits, from all kinds of remote sensing sensors and platforms. Still,  
95 CNN-based vegetation remote sensing is a very topical but young field of research.  
96 People with a background in remote sensing or vegetation science may require  
97 procedural knowledge on the working principles of CNNs and the anticipated po-

98 tentials for vegetation mapping. In contrast, people from computer sciences may  
99 require declarative knowledge on application tasks in vegetation science, on types  
100 and availability of remote sensing data suitable for vegetation analysis, or on the  
101 relationship between remotely sensed signals and vegetation properties. Thus, the  
102 overall aim of this review is to link procedural and declarative knowledge and pro-  
103 vide an introduction and synthesis on the current state of the art on the utility of  
104 CNNs for vegetation remote sensing.

105 The present review is organized into three main sections: Chapter 2 briefly in-  
106 troduces the basic principles and the general functioning of CNNs and deduces why  
107 it is such a promising method for remote sensing of vegetation. Chapter 3 provides  
108 a summary and meta-analysis on the corresponding literature and synthesizes the  
109 current state of the art and challenges, including:

- 110 • common CNN approaches, architectures and strategies for the retrieval of  
111 vegetation properties,
- 112 • an overview of common applications tasks and demonstrated potentials in  
113 the context of agriculture, forestry and conservation,
- 114 • challenges and corresponding solutions regarding reference data quantity and  
115 quality of continuous and discrete vegetation variables,
- 116 • a consideration of spatial and spectral resolution for CNN-based vegetation  
117 remote sensing and considerations towards different sensors, platforms and  
118 combinations thereof.

119 Lastly, chapter 4 gives concluding remarks and discusses possible future direc-  
120 tions and developments.

## 121 **2 Principles of CNNs and relevance for vegetation** 122 **remote sensing**

123 This chapter introduces the basic principles of CNN, including the functioning of  
124 convolutions, features that make convolutions suitable for vegetation analysis, and  
125 how a CNN is commonly trained and implemented.

### 126 **2.1 Basic functioning and structure of CNNs**

127 As any typical neural network-type model, CNNs are based on **neurons** that are  
128 organized in **layers** and can, hence, learn hierarchical representations. The neurons  
129 between layers are connected through weights and biases. The initial layer is the  
130 input layer, e.g. remote sensing data, and the last layer is the output, such as a  
131 predicted classification into plant species. In between are **hidden layers** trans-  
132 forming the feature space of the input in a way that it matches the output. CNNs  
133 include at least one convolutional layer as a hidden layer to exploit patterns (in the  
134 context of this review predominantly spatial patterns).

135 It can also include other non-convolutional layers. Convolutional layers include  
136 multiple optimizable filters (Fig. 1) that transform the input or preceding hidden  
137 layers. The number of filters defines the **depth** of a convolutional layer. The  
138 resulting transformations are aimed to reveal patterns that are decisive for the

139 problem at hand. The decisive patterns are iteratively learned through convolving,  
 140 which is essentially the sliding of the filter over the layer and the calculation of the  
 141 dot-product of the filter and the layer’s values. The result is a new layer of dot-  
 142 products for each filter, also called a **feature map** (Fig. 1). The early feature maps  
 143 in a CNN may include simple and fine scaled patterns, such as corners, circles, or  
 144 edges. The derived feature maps then serve as input for the next layer, e.g. another  
 145 convolutional layer or a final layer that predicts an outcome based on the detected  
 146 features. In deeper layers of a network, convolving usually reveals more abstract  
 147 patterns and higher-level concepts, such as leaf forms, branching patterns or habit.  
 148 During model training, randomly initialized filters will be iteratively optimized to  
 149 detect the relevant image features (the training procedure is described in Section  
 150 2.3). The combination of several successive convolutional layers with their numerous  
 151 filters, hence, enables the network to learn and combine even subtle image features,  
 152 revealing if a class is present in an image or not (see Fig. 1 for a tree-species-specific  
 153 activation of the network, more details on class activation mapping in Section 3.6.2)

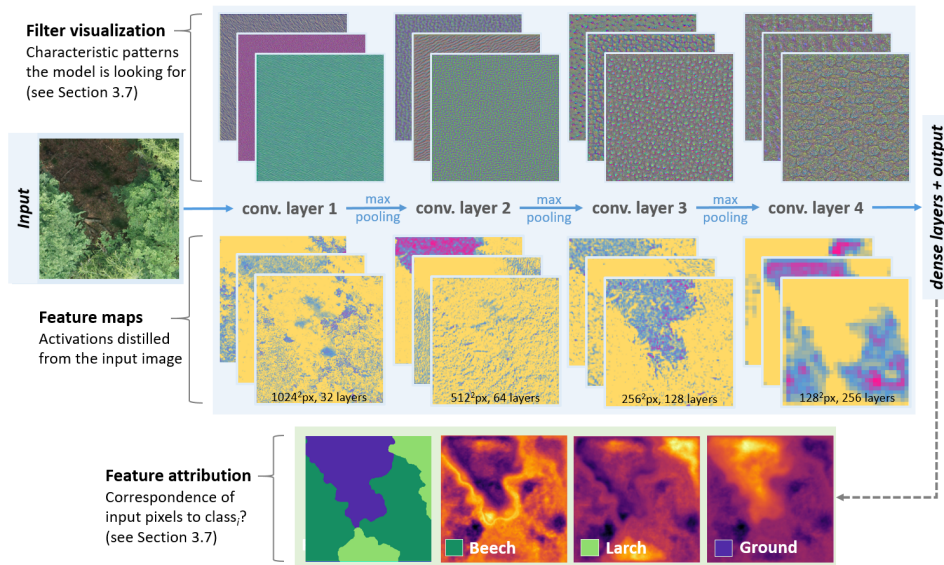


Figure 1: Scheme of a CNN composed of four convolutional layers and subsequent pooling operations trained for tree species classification. The visualization of convolutional filters (top) indicate characteristic patterns the CNN is looking for and were derived by gradient ascent; a technique revealing artificial images maximizing each filter’s activation. The feature maps (center) are the dot-product of the preceding layer and individual filters. Feature attribution maps (bottom) can reveal individual pixels that were decisive for the tree species assignment (details on feature attributions 3.6.2).

154 Between sequences of multiple convolutional layers, the feature maps are com-  
 155 monly spatially down-sampled using spatial **pooling operations** (see Fig. 1).  
 156 Pooling describes the transformation of multiple cells into one cell, similar to re-  
 157 sampling an image to a coarser spatial resolution. Pooling has several advantages:  
 158 It reduces the data size while preserving discriminant information, which in turn de-

159 creases the number of model parameters, thus computational load and the chance  
160 of overfitting; and it enables detecting more abstract features as well as spatial  
161 context across scales and thereby condenses semantic information. Pooling is de-  
162 fined by a filter size, stride (the distance between consecutive pooling operations),  
163 and a reduction operation. The most typical pooling operation is **max-pooling**.  
164 The idea of max-pooling (instead of, for instance, average pooling) is that strong  
165 activations (e.g. edge or line features) are conserved within the network and not  
166 averaged out. A typical max-pooling operation with 2-by-2 filter size and a stride  
167 of 2 reduces the size of the input feature map by a factor of 4, whereas the output  
168 cells contain the maximum value of the 4 input cells within the 2-by-2 filter.

169 The layers of CNNs, e.g. convolutional or pooling layers, can be combined in  
170 very different ways - commonly described as the CNN **architecture**. CNNs can,  
171 hence, have very different architectures, which are basically defined by the task. The  
172 task can be the classification of images, the segmentation of multiple classes, or the  
173 localization of individual objects within a scene (presented in more depth in chapter  
174 3.2.2). The suitability of a CNN architecture largely depends on the complexity  
175 of the task: A more complex problem usually requires a deeper and more sophis-  
176 ticated network. In contrast, limited availability of training data constrains model  
177 complexity due to an increased risk of overfitting. The complexity and general per-  
178 formance of a CNN architecture further depends on the **hyper-parameters**, which  
179 define amongst others the number and characteristics of hidden layers, pooling op-  
180 erations, regularization techniques, or cost-functions. Accordingly, there exists a  
181 wide array of options to implement a CNN towards the specific use case as well  
182 as predefined and established architectures. Examples are given in the literature  
183 review in chapter 3.2. Comprehensive overview of different architectures is given in  
184 Hoeser et al. (2020) and Zhu et al. (2017).

## 185 2.2 Why CNN for vegetation remote sensing?

186 The physiology and morphology of vascular plant canopies is primarily optimized  
187 towards the absorption of solar energy using the photosynthetic machinery and the  
188 corresponding assimilation of carbon for maintenance, further growth, and repro-  
189 duction. Despite these common goals among vascular plants, plant life can differ  
190 greatly on multiple scales, ranging from various morphological features of the in-  
191 dividual, including leaf tissue properties, leaf form, branching patterns, canopy  
192 structure, and the general habitus, to large-scale patterns of vegetation communi-  
193 ties. Furthermore, anthropogenic land use can determine spatial vegetation pat-  
194 terns, either through indirect influences on floral vitality and species composition  
195 or directly through economic activities. Examples include dendritic or fish bone-  
196 like deforestation structures in rain forests, crop rows on plantations, or directed  
197 changes in species composition as a result of gradual nutrient inputs from agricul-  
198 tural land.

199 Remote sensing offers several sensors and acquisitions techniques that are sen-  
200 sitive to physiological and morphological properties of vegetation and, hence, allow  
201 for spatial representations of vegetation patterns at the scale from plant organs  
202 to entire landscapes. This includes close-range observations from terrestrial plat-  
203 forms (e.g., farming robots), fine-resolution data from airborne platforms (UAV or  
204 airplanes) as well as more coarser-resolution satellite-based acquisitions that are  
205 usually focused on large-scale applications.

206 So why are CNNs suitable for vegetation analysis with such remote sensing  
207 data? CNNs are indeed a revolutionary technique but they do not do magic, mean-  
208 ing that they cannot reveal more information than is contained in the data. The  
209 crucial advancement of CNNs is *how* they can extract information from spatial  
210 data. Previous parametric or machine learning methods applied in vegetation re-  
211 mote sensing usually required feature engineering, i.e. the careful screening of  
212 redundancies in the input data and the extraction of latent variables that best de-  
213 scribe the response variable. Simply put, the model needs to be taught how to *see*  
214 the relevant features before it can start solving the problem. Feature engineering  
215 is, hence, based on an understanding of a system and its processes. This enables to  
216 control the model with pre-knowledge but is certainly limited in case of unknown  
217 systems that potentially inherit many dimensions and complex interactions. Espe-  
218 cially for the analysis of 2D or 3D patterns, there are a plethora of transformations  
219 that can be applied to extract spatial features and textures. Examples are Grey  
220 Level Co-Occurrence Matrices (Haralick 1979), Fourier Transformations (Bone et  
221 al. 1986), or 3D multi-scale metrics derived from point clouds (Brodu et al. 2012;  
222 Weinmann et al. 2015). These numerous types of transformations can moreover be  
223 applied with different hyper-parameters (e.g., kernel function or size). The poten-  
224 tial amount of latent variables extracted this way explodes, considering that one  
225 can extract latent variables with such transformations based on different input data  
226 available, e.g., different wavelengths of a multispectral sensors or snapshots from  
227 a time series. Thus, identifying the best combination of possible predictors on a  
228 heuristically basis is often a very inefficient task and often hardly possible.

229 In contrast, a CNN itself learns the ability to *see* by iteratively optimizing the  
230 transformations, i.e., the convolutional layers, during the training process. This  
231 *end-to-end* learning principle can make feature engineering obsolete and, thus, pro-  
232 viding the raw data (e.g. spectral bands or the point cloud) can be already sufficient.  
233 Additional feature engineering, e.g., transformations like vegetation indices or pre-  
234 processing such as speckle reduction, may even introduce an information loss and  
235 decrease the model accuracy (Geng et al. 2017; Hartling et al. 2019; Sothe et al.  
236 2020). In contrast to statistical modeling or machine learning, deep learning, hence,  
237 shifts the focus from *what* a model should learn to *how* a model should learn. The  
238 latter is primarily defined by the model architecture and the optimization of its  
239 hyper-parameters as discussed in the following sections.

### 240 **2.3 The training process**

241 Training a CNN model for vegetation mapping requires the remote sensing data  
242 and matching reference annotations, also called labels or targets. While machine  
243 learning algorithms, such as random forests or support vector machines, require rel-  
244 atively simple array-type data structures, CNN-based training is performed using  
245 more sophisticated data structures called **tensors**. Tensors are essentially stacked  
246 arrays that typically have 4 dimensions, including the individual samples, the spa-  
247 tial dimensions (x, y), a feature dimension (z, e.g. intensity or reflectance), and a  
248 layer dimension (e.g. the corresponding wavelength).

249 During training, the CNN weights are optimized for a certain task, e.g., de-  
250 tecting a certain plant species. This detection is realized by transforming the input  
251 data through convolutional and other hidden layers while being propagated through  
252 the network. The neurons between layers are connected through **activation func-**



253 **tions** determining if a neuron is active – also referred to as firing - or not (ReLU,  
254 the most frequently used activation function, is described in chapter 3.2.1.1). If  
255 activated, the intensity of a neuron’s output is determined by its weights and bi-  
256 ases. The weights and biases are usually optimized using the **gradient descent**  
257 algorithm, which can briefly be described as follows: The term gradient descent  
258 implies the progressing minimization (descent) of errors along a slope (gradient).  
259 Gradient descent is performed in iterations, in which predictions of a model with  
260 momentary parameterization are compared to the annotations of the training data  
261 using a **loss function**. The gradients are derived using the **backpropagation**  
262 algorithm. Given a neural network with an input layer (a tensor), an output layer  
263 (prediction) and  $n$  hidden layers in-between (e.g. convolutional layers), the back  
264 propagation algorithm calculates the gradient of the loss function with respect to  
265 the weights and biases between the hidden layers. This gradient is then used to  
266 evaluate and update the model weights and biases through gradient descent, i.e.  
267 trying to find a global minimum in the high-dimensional feature space. The gradi-  
268 ent descent procedure is performed for multiple samples, followed by averaging the  
269 calculated weights and biases of the hidden layers.

270 Training a CNN is usually computationally very intensive as the explanatory  
271 variables, e.g. image data or point cloud representations, are rich in dimensions  
272 (geolocation + layers) resulting in a myriad of feature maps that depict different  
273 spatial features and context at varying scales. This obviously results in excessive  
274 amounts of data to be processed during CNN training – especially considering that  
275 model training may require many samples to memorize the decisive features of the  
276 target class. These data volumes may, thus, not fit the memory of our system  
277 at once. To overcome this, training is often performed sequentially in **batches**  
278 comprising only a share of the entire dataset. The model weights and biases are  
279 updated based on one average gradient for the entire batch. Separating the dataset  
280 into batches enables to train the model iteratively until it has seen all samples,  
281 which is called an **epoch**. The number of iterations to finish an epoch is, thus, the  
282 total number of observations divided by the batch size.

283 Generally, it is unlikely that a CNN trained in a single epoch already reaches  
284 maximum performance. For instance, observations (in form of batches) that were  
285 shown to the model at the beginning of the training phase may be again useful to  
286 extract more features at a later stage of the training process. Moreover, multiple  
287 steps in the training procedure described above feature stochasticity: The convo-  
288 lutions are based on randomly initialized filters, the assignment of observations  
289 into batches is random, and the gradient descent has a random nature (hence, also  
290 referred to as *stochastic* gradient descent). For this very reason, CNNs are com-  
291 monly optimized within a series of subsequent epochs until the model performance  
292 stops to advance (the model converges) or even decreases (the model overfits). The  
293 number of epochs eventually depends on the complexity of the problem and model  
294 structure.

295 The fact that gradient descent is an iterative algorithm opens several interesting  
296 avenues for CNN-based modelling: Firstly, models can be updated with unseen data  
297 at any time without training the model again from scratch substantially saving  
298 computation loads and processing time. Secondly, models that have seen a lot  
299 of data, e.g. from generic image databases such as *ImageNet*, can be shared and  
300 optimized for a specific problem (further discussed in Section 3.2.1.2). The third  
301 and probably most future-oriented avenue is **federated learning**, which is the

302 training of local models with local data on distributed clients and the simultaneous  
303 sharing of weights coordinated by a central server (Bonawitz et al. 2019). The  
304 server thereby merges the locally derived gradients without ever seeing the data.  
305 Federated Learning follows, thus, the principle of *bringing the code to the data,*  
306 *instead of the data to the code,* which will be inevitable in the geosciences due  
307 to constantly growing data streams. Besides reducing communications costs, this  
308 approach avoids problems related to data access rights, security, or privacy.

## 309 2.4 Implementation, libraries and frameworks

310 Most deep learning frameworks can be used on standard operating system (Linux-  
311 based, Windows, macOS) and provide bindings for different programming lan-  
312 guages. Currently, Python is the most common language in DL research. Training  
313 and inference of deep learning models consist of millions of simple computations,  
314 i.e. multiplications and additions. Thus, it is helpful to use **graphics processing**  
315 **units** (GPU) rather than central processing units (CPU). In contrast to CPUs,  
316 GPUs have rather simple cores but thousands of them, which are optimized to  
317 handle thousands of concurrent operations, leading to a drastic reduction of time  
318 for training and inference. Mostly *NVIDIA* GPU are used, as these feature the  
319 *CUDA Deep Neural Network (cuDNN)* library, which is utilized by common DL  
320 frameworks. The *cuDNN* library provides highly performant primitives for convo-  
321 lutions, pooling operations, normalization and activation functions. Furthermore,  
322 *AMD* provides different tools for deep learning on Linux-based platforms with the  
323 *Radeon Open Compute Platform*. In case of missing hardware it is nowadays pos-  
324 sible, to use (partially free) cloud platforms with GPU support, such as *Alibaba*  
325 *Cloud*, *Amazon Web Services*, *Microsoft Azure* or *Google Cloud Platforms* such as  
326 *Google Earth Engine*. These platforms have completely configured containers for  
327 many frameworks. *Google Colab* <https://colab.research.google.com> even provides  
328 free access to (limited) computing resources including GPUs with no setup.

329 CNNs can be implemented through different **frameworks**. Overviews of for-  
330 mer and current frameworks are given in Hoeser et al. (2020), Nguyen et al.  
331 (2019) and on the corresponding Wikipedia page ([https://en.wikipedia.org/wiki/](https://en.wikipedia.org/wiki/Comparison_of_deep-learning_software)  
332 [Comparison\\_of\\_deep-learning\\_software](https://en.wikipedia.org/wiki/Comparison_of_deep-learning_software)). The currently most prominent deep learn-  
333 ing frameworks are **PyTorch** and **Tensorflow** (Nguyen et al. 2019). Both provide  
334 high-level APIs (e.g. Keras) and various tools for training, data augmentation,  
335 and visualization (e.g. Tensorboard). Furthermore, many vintage and modern DL  
336 architectures can be used directly and with pretrained weights. Extensive documen-  
337 tations, many tutorials, and Jupyter notebooks allow an easy start with both open-  
338 source frameworks. Additionally, the Open Neural Network Exchange (ONNX)  
339 format allows interoperability between many frameworks such as Pytorch, Ten-  
340 sorflow, Keras, mxnet, scikit-learn, Matlab, SAS, and many more. Thus, already  
341 implemented and trained models can be transferred to a favored framework. In  
342 Section 5 links to various tools, models and quick start tutorials are provided.

### 3 Literature review on CNN-based vegetation remote sensing

The literature review was based on a survey on *Google Scholar* and the search terms *CNN, convolutional neural networks, vegetation, plants, forestry, agriculture, land cover, conservation, mapping, Remote Sensing, RGB multispectral, LiDAR TLS, ALS, SAR, RADAR, airborne, satellite, UAV*. The search results were first filtered by the title, by the abstract and then by the content. We only considered primary research articles that underwent a peer-review process. This resulted in a total of 101 research studies considered in the literature review. All studies were published after 2016 and more than 75 % of the studies were published in 2019 or later (see Fig. 2), underlining that CNN-based vegetation remote sensing is a very young but rapidly developing field.

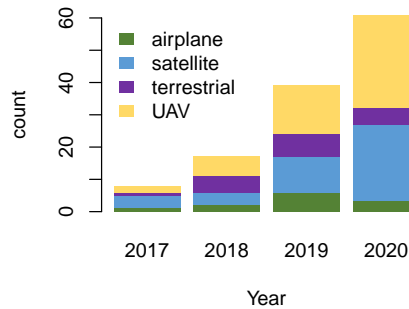


Figure 2: Number of yearly publications based on the literature search indicating a steep increase of studies applying CNNs for vegetation remote sensing. Counts for 2020 were extrapolated based on the number of publications until November.

The resulting literature is very heterogeneous in terms of application areas, vegetation types, target variables, CNN implementations, and remote sensing data (compare Fig. 3). Accordingly, several criteria were defined to structure the literature and identify general trends, including the underlying CNN architecture, remote sensing platform, sensor, spatial resolution of the remote sensing data, mode of reference data acquisition (in-situ or by visual interpretation), number of training and test observations, response type (e.g., object detection or semantic segmentation), geographic location of the study area, accuracy metrics, area of application (agriculture, forestry, conservation or miscellaneous) and specific task (e.g, detecting weed infestation or tree cover mapping). A corresponding spreadsheet including all assessed criteria and studies is available in the Appendix. For the accuracy metrics, we constrained our analysis on the most frequently reported metrics (overall accuracy, precision, recall, F-score and intercept over union). Whenever a study reported multiple accuracy metrics, e.g. when comparing multiple methods, we recorded the best result. The geographic locations of the study areas were derived

370 from place-names using the *Google Geocoding API*, unless the manuscripts explicitly  
371 included the longitude and latitude of the study area.

## 372 3.1 Reference data

### 373 3.1.1 Reference Data Sources

374 As with any supervised modelling approach, training and validating a CNN re-  
375 quires reference observations, also referred to as annotations, labels, or targets.  
376 The large number of parameters in CNNs and the corresponding ability to detect  
377 even subtle patterns are associated with the risk in training a model that is based  
378 on overly-specific details and does not generalize well - it is overfitting. Accordingly,  
379 independent validation of CNNs prior to model deployment is of great importance  
380 to evaluate its robustness and transferability. Ideally, such validation should not  
381 solely involve iteratively shuffling training and validation data, as frequently done  
382 in remote sensing studies (e.g., as with a cross-validation or bootstrapping), but be  
383 based on entirely independent data that the model has never seen before. There-  
384 fore, most CNN-related studies split their reference data in a 1) training data set,  
385 which commonly is split again in **training** and **validation** data during the model  
386 training process, and 2) a **testing** data set used to independently evaluate the  
387 eventual predictive performance of the final model. Typically, a share of 20 to 30  
388 % of the reference data is used for independent testing (median 21 %).

389 In the field of remote sensing of vegetation, reference data was most commonly  
390 acquired in **ground-based surveys** in the form of in-situ plot or point observations  
391 (Fassnacht et al. 2016). The quantity of reference data of ground-based surveys is  
392 generally limited as these involve high logistic efforts and costs for transportation,  
393 equipment, and personnel. In particular for studies in natural environments, lim-  
394 ited accessibility can also greatly hamper the sampling frequency. The effectiveness  
395 of ground-based surveys for CNN modelling may, hence, be limited as the latter  
396 often requires ample reference data. In particular for complex tasks, such as the  
397 differentiation of classes that only differ in subtle features, the quantity of avail-  
398 able reference data can be the critical factor for a successful model training and  
399 convergence. Moreover, tasks as object detection or the segmentation of individ-  
400 ual crown components (3.2.2) require reference data that is spatially explicit and  
401 in exact correspondence with the remote sensing data. Especially for analysis of  
402 very high spatial resolution remote sensing data at centimetre scale, GNSS-coded  
403 reference data acquired in the field is often not directly applicable for two main  
404 reasons: Firstly, geolocation errors of GNSS-measurement typically exceed 0.1-1m;  
405 particularly under dense vegetation canopies (Branson et al. 2018; Kaartinen et al.  
406 2015; Valbuena et al. 2013). Secondly, for practical reasons, field data is usually  
407 measured in form of point observations (e.g., stem position of a tree) or using circ-  
408 ular or rectangular plots, which does commonly not allow for a spatially explicit link  
409 with remote sensing data (Anderson 2018; Kattenborn et al. 2019d; Leitão et al.  
410 2018). Correspondingly, only 14 % of the studies reviewed here used in-situ data  
411 as exclusive reference input.

412 Instead of using in-situ observations, reference data is most often (62 %) di-  
413 rectly acquired in the primary or secondary (e.g., higher resolution) remote sensing  
414 data using **visual interpretation**. In contrast to common in-situ point or plot ob-  
415 servations, reference data acquired by visual interpretation is commonly spatially

416 explicit as it is directly derived from the imagery or point cloud. Furthermore,  
417 there is no position error, as long as the same input data is used for the CNN and  
418 visual interpretation. If secondary data (e.g. higher resolution) is used for visual  
419 interpretation, the geolocation error is relative to the spatial agreement of primary  
420 and secondary data. Visual interpretation provides a very efficient mode of gener-  
421 ating reference data, given that the variable of interest is clearly identifiable in  
422 the imagery. Accordingly, this mode of reference data acquisition is in particular  
423 applicable for discrete classes (e.g. species, plant communities, crop or vegetation  
424 types, individuals). The term visual *interpretation* implies a rather imprecise cap-  
425 ture of the target metric, but it should be noted that in-situ observations do not  
426 necessarily represent (ground) *truth*: As with visual image interpretation, mapping  
427 species in the field is commonly based on visual interpretation and, hence, can also  
428 be prone to errors and bias (Lepš et al. 1992; Lunetta et al. 1991).

429 Annotations from visual interpretation are often derived by delineating tar-  
430 get classes in a GIS environment. This includes the identification of individuals  
431 by points, as often performed for image-based object detection in agricultural en-  
432 vironments (Csillik et al. 2018; Freudenberg et al. 2019), or by delineating the  
433 vegetation components (e.g. in form of polygons) for semantic or instance segmen-  
434 tation (Flood et al. 2019; Kattenborn et al. 2019a). Many studies have also used  
435 special interfaces for an efficient labeling such as *RectLabel*, *LabelMe*, *Labelbox* or  
436 *LableImg* (Russell et al. 2008). Instead of manually labeling the spatial extent of  
437 target classes, a semi-automatic approach using a prior segmentation may be used.  
438 For instance, dos Santos Ferreira et al. (2017) automatically segmented canopy  
439 components in RGB imagery of soybean fields using SLIC (Simple Linear Iterative  
440 Clustering) superpixels (Achanta et al. 2012), and assigned each segment to weeds  
441 or crops by visual interpretation. Natesan et al. (2019) labeled segments derived  
442 from a watershed-based segmentation using a Digital Surface model. In particular,  
443 for LiDAR-based point cloud data, region growing algorithms may be used to effi-  
444 ciently segment points belonging to individual plants (Wang et al. 2019) or plant  
445 components (e.g. stems, branches or foliage; Z. Xi et al. 2018).

446 Despite the above-mentioned advantages, obtaining reference data by visual  
447 interpretation does not rule out misinterpretation. Yet, at the example of mapping  
448 plant species, it has been shown that CNNs can to some extent compensate flawed  
449 or noisy labels (Hamdi et al. 2019; Kattenborn et al. 2020).

450 Although in-situ data may not be the ideal for training and validating CNNs, it  
451 may be an essential requirement in case the target class (e.g. species) is not readily  
452 identifiable in the remote sensing data by means of visual interpretation alone.  
453 According to our review, 22 % of the studies that acquired reference data by visual  
454 interpretation also incorporated in-situ data for training or validation. 84 % of  
455 these studies were either related to forestry or conservation tasks and thus to rather  
456 complex environments, in which visual interpretation alone may not be sufficient.  
457 For instance, Schiefer et al. (2020) and Kattenborn et al. (2019a) used ground-based  
458 full inventory data as a basis to annotate tree species in UAV imagery in temperate  
459 forests forests in Germany, and in highly heterogeneous and complex natural forests  
460 in Waitutu, New Zealand, respectively. Similarly, Sun et al. (2019) used in-situ data  
461 on tree species to map the species diversity in tropical wetlands. Field data may  
462 also provide an independent source to validate CNN-based predictions (Flood et  
463 al. 2019). Especially in cases when a bias by visual interpretation is assumed, a  
464 validation using in-situ reference data is highly recommended.

465 Visual interpretation may be more efficient for data annotation than using  
466 in-situ data alone, but even human labeling through visual interpretation can be  
467 very tedious, especially for large datasets or complex vegetation canopies that re-  
468 quire very detailed annotations. The effort of annotating data may be reduced by  
469 specific training strategies, such as **weakly-** or **semi-supervised learning** (see  
470 section 3.2.1.3), that compensate for few or coarse annotations. Alternatively, if  
471 no knowledge of a vegetation expert is required, crowdsourcing can be used for  
472 labeling. Commercial services are now also available for this purpose. For exam-  
473 ple, Branson et al. (2018) used the service *Amazon Mechanical Turk*<sup>TM</sup> to locate  
474 individual trees in *Google Street View* imagery.

475 Although visual interpretation is an effective labeling approach to many tasks,  
476 it should be noted that there are many vegetation-related applications where it is  
477 not applicable. Particularly, for continuous quantities, such as crop yield or forest  
478 biomass (Ayrey et al. 2018; Castro et al. 2020; Yang et al. 2019), reference data  
479 acquisition is conceptually more difficult as these are often not directly measurable  
480 from the remote sensing data. Here, in-situ measurements or other physically-based  
481 retrieval procedures may often present the only applicable solution. A physically-  
482 based retrieval of reference data was presented by Du et al. (2020), who aimed at  
483 mapping wetland inundation extent in forests on large spatial scales with satel-  
484 lite data (WorldView-2). For parts of their study area, LiDAR data was avail-  
485 able enabling accurate detection of surface waters due to its strong absorption in  
486 near-infrared wavelengths. Reference data acquisition on yield or biomass in an  
487 agricultural context may be automatized by integrating measurement devices on  
488 harvesting machines. For instance, Nevavuori et al. (2019) trained a CNN to pre-  
489 dict wheat and malting barley yield from UAV imagery using training data derived  
490 from a yield measurement device (*John Deere Greenstar 1*) that was coupled with  
491 a GNSS receiver and mounted on a harvester.

492 Concerning biochemical and structural plant traits, an interesting approach is to  
493 train CNNs with simulated data derived from physically-based models. Such hybrid  
494 approaches, i.e. coupling statistical and process-based models, may not only provide  
495 data for training but also enable including priors and realistic constrains in model  
496 training (Reichstein et al. 2019). For instance, Annala et al. (2020) trained a 1D-  
497 CNN with reflectance spectra simulated with the radiative transfer model (RTM)  
498 SLOP (Maier et al. 1999). Although SLOP is a relatively simple leaf reflectance  
499 model, Annala et al. (2020) demonstrated promising tests of this hybrid inversion  
500 method for UAV hyperspectral acquisitions of forest canopies. More sophisticated  
501 RTMs may allow to produce more robust models, e.g. *PROSAIL* (Jacquemoud et  
502 al. 2009) enabling to account for bidirectional reflectance effects in plant canopies,  
503 whereas 3D-RTMs such as *FLIGHT* (North 1996) or *DART* (Gastellu-Etchegorry  
504 et al. 1996) may provide interesting sources for generating synthetic training data  
505 for 2D-CNNs (see Section 3.2 for details on 1D-, 2D- and 3D-CNNs).

### 506 3.1.2 Reference data quantity

507 The quantity of reference data required for the convergence of a CNN depends par-  
508 ticularly on the complexity of the algorithm and most importantly on the contrast  
509 of the features that are decisive for the vegetation property of interest. Fewer refer-  
510 ence data may be required if the vegetation property of interest is easily identifiable  
511 in the remote sensing data (e.g., due to a distinct canopy structure or contrasting

512 flowers). Subtle differences and complex relationships in turn require more com-  
513 plex algorithms and more samples to identify the relevant features. Accordingly,  
514 the effects of varying the training data size cannot be generalized. The results of  
515 Weinstein et al. (2020) suggest that the accuracy first increases rapidly with in-  
516 creasing the reference data quantity and then stagnates. In the context of tree  
517 species mapping in urban environments, Hartling et al. (2019) showed that using  
518 10 % of their available training samples decreased the overall accuracy from 82.58  
519 % to 70.77 %. Using 200-3940 samples and multiple CNN architectures, Fromm  
520 et al. (2019) showed that the reference data quantity can have a large influence on  
521 the overall accuracy for tree seedling mapping (up to 18 %). Using UAV data for  
522 segmenting growth forms in wetlands, T. Liu et al. (2018b) demonstrated that the  
523 effect of sample size (700-3500 samples) can greatly differ across different model  
524 architectures and complexities.

525 Overall, the amount of reference data used in the reviewed studies differed  
526 greatly - most notably between studies using different remote sensing platforms.  
527 Studies based on terrestrial data acquisitions, e.g., terrestrial or mobile LiDAR  
528 scanning, used around 340 reference observations (median). UAV- or airborne-  
529 related studies used a median of 2795 reference observations and studies based on  
530 satellite observations 6001 observations. These large differences may be the result  
531 of two factors: Firstly, studies at the satellite-scale typically cover larger spatial  
532 extents and are, hence, more likely to benefit from previously acquired reference  
533 data sets (cf. Schmitt et al. (2020)), whereas, the coarser spatial resolutions also  
534 allow to incorporate reference data with higher geolocation errors. Secondly, data  
535 acquired at higher resolutions, often TLS or MLS LiDAR data, contains finer infor-  
536 mation on vegetation structures and may thus include more characteristic features.  
537 This may hence facilitate model convergence and decreases the amount of reference  
538 data required.

539 A common training strategy that aims to compensate for few reference data  
540 is **data augmentation**, which inflates the number of reference data by introduc-  
541 ing small manipulations to the existing data or creating synthetic data (see details  
542 section 3.2.1.1). Instead of collecting new reference data, it may be more efficient  
543 to use existing reference data, e.g. from previous research projects or authorities  
544 (e.g. environmental agencies, forestry offices). Accordingly, the establishment of  
545 open access **databases** incorporating labeled remote sensing data is increasingly  
546 demanded but still lacking (Zhu et al. 2017). Such databases would not only fa-  
547 cilitate the efficiency of model training due to ample training data, but would also  
548 allow to assess and improve the extrapolation and transferability of these models  
549 to new domains. This is particularly important as geoscientific models are often  
550 under-constrained due to limited representatives of the training data (Reichstein  
551 et al. 2019). Accordingly, databases can enable to test and improve the model  
552 transferability towards new domains, such as different remote sensing acquisitions  
553 (daytime or sensors), vegetation types, or growth stages. Moreover, databases of  
554 sufficient size could also play an important role to develop backbones that are specifi-  
555 cally oriented to vegetation remote sensing (further discussed in Section 3.2.1.2).  
556 Freely accessible databases can also facilitate more comprehensive and universal  
557 comparisons of algorithms and the identification of improvement opportunities.

558 Despite the described benefits, there exist still only a few databases providing  
559 labeled remote sensing data, which may be explained by the novelty of the scien-  
560 tific field (cf. Fig., 2), associated costs for data storing and sharing (especially in

561 regard to high-resolution data), various fields of application with individual anno-  
562 tation requirements, and lastly the diversity in remote sensing sensors, acquisition  
563 and processing modes. A prime example is the voluntarily organized *ImageCLEF*  
564 initiative (imageclef.org). The latter hosts an evaluation platform and mostly an-  
565 nually recurring competitions for cross-language annotation of images (Kelly et al.  
566 2019). The first competition was hosted in 2003 and aimed at classifications of  
567 generic photograph datasets, whereas in 2011 the first vegetation-specific competi-  
568 tion followed, which was centered on plant species identification from ordinary pho-  
569 tographs. Since 2017, *ImageCLEF* also hosts the *GeoCLEF* competitions, which  
570 focus on plant species identification by means of environmental and remote sensing  
571 data, including high-resolution remote sensing imagery and respective land cover  
572 products. Another example is the *NSF NEON* database (Kampe et al. 2010; Kao et  
573 al. 2012; Marconi et al. 2019) including a wide array of (partly multitemporal) refer-  
574 ence and remote sensing data (most importantly from RGB, LiDAR, hyperspectral  
575 airborne campaigns) on natural and semi-natural ecosystems. This database has  
576 already been proven to be of immense value to train and validate models across  
577 ecosystems and remote sensing acquisitions (Ayrey et al. 2018; Weinstein et al.  
578 2020). For instance, Weinstein et al. (2020) tested cross flight performance of a  
579 CNN for tree crown segmentation in different environments. Their results under-  
580 lined the value of large databases for model training, as the model generalization  
581 with additional datasets greatly improved - even when the target class was not  
582 present in all datasets. Example centered on developing and benchmarking deep  
583 learning towards vegetation types and land-cover mapping with Sentinel-2 imagery  
584 are the *SEN12MS* (Schmitt et al. 2019), *BigEarthNet* (Sumbul et al. 2019) and *EU-*  
585 *ROSAT* (Helber et al. 2019) datasets. In the agricultural context, the *Global Wheat*  
586 *Dataset* (global-wheat.com) includes standardized images on weeds ( $1024 \times 1024$   
587 pixels) with subcentimetre resolution, providing the basis for public challenges, such  
588 as the 2020 challenge to count wheat ears (David et al. 2020).

589 An alternative approach could also be the use of databases that only refer  
590 to vegetation information but can be linked to existing remote sensing data in  
591 other ways, e.g. by taxonomic identities or geo-coordinates. Valuable resources in  
592 this context are the *TRY* database (try-db.org, kattge2020try), which contains a  
593 wealth of morphological, physiological and phenological plant traits, the *opentrees*  
594 database (opentrees.org, providing species and location information of individual  
595 trees in urban areas, or *GBIF* (gbif.org, providing several huge datasets on citizen-  
596 science-based plant photographs together with species names and geo-coordinates,  
597 including the popular *iNaturalist* dataset.

## 598 **3.2 Common CNN approaches and architectures**

### 599 **3.2.1 Training strategies**

600 Training a CNN can be challenging due to a restricted amount of labeled obser-  
601 vations, computation load required for model convergence, and model overfitting.  
602 This chapter lists the most common strategies and methods applied during training  
603 to alleviate these challenges.



### 604 3.2.1.1 Normalization and regularization techniques

605 A famous problem in training artificial neural networks with gradient-based learn-  
606 ing is the **vanishing** or **exploding gradient problem** (Hochreiter 1991, 1998).  
607 During backpropagation, the weights of each node are updated proportionally to  
608 its gradient in respect to the loss. The gradients are derived by calculating the  
609 derivative of an activation function. For a common sigmoid function, this deriva-  
610 tive becomes increasingly small for very low or high values. The derivative of a  
611 layer is calculated by the chain rule and so gradients and corresponding updates of  
612 weights in earlier layers of the network can approach zero (vanish). The opposite  
613 effect, i.e. exploding gradients, can occur for large derivatives. This imbalance in  
614 the network ultimately impairs the network’s ability to find the ideal updates for  
615 the weights.

616 A common counter-measure is **batch normalization**, which is applied in 26  
617 % of the reviewed studies, particularly in networks with many parameters such  
618 as for semantic segmentations (Kattenborn et al. 2019a; Ronneberger et al. 2015;  
619 F. Wagner et al. 2019). Batch normalization normalizes the output of activation  
620 functions to zero-mean and unit variance and thereby prevents the network from  
621 becoming imbalanced due to excessively high or low activations. This smooths the  
622 optimization problem of the gradient descent function and allows for larger ranges  
623 of learning rates and hence facilitates network convergence.

624 The vanishing gradient problem can also be greatly reduced by using the **Rec-**  
625 **tified Linear Unit (ReLU)** activation function. The output weight of the ReLU  
626 function equals the weighted sum of the inputs as long as this sum is  $> 0$  (values  
627  $< 0$  are ignored). For  $> 0$ , ReLU is a simple linear function such that the deriva-  
628 tive is always 1, hence, preventing the vanishing gradient problem. The probably  
629 more important characteristics of ReLU are its non-linearity and its **regulariza-**  
630 **tion** function of the network. The large amount of parameters in deep networks  
631 makes them prone to overfitting and, therefore, regularization aims to facilitate  
632 a network’s ability to generalize. ReLU regularizes the network by reducing the  
633 parameters of the model as it ignores values  $< 0$  - these values are in theory not  
634 activated anyway. The reduction of parameters also greatly decreases the com-  
635 puting time in contrast to conventional hyperbolic tangent functions (Krizhevsky  
636 et al. 2012). Only few studies reported that they used other activation functions  
637 suggesting that in fact most of the studies used ReLU.

638 One of the most common and effective regularization technique is **Dropout**(Srivastava  
639 et al. 2014) (used in at least 31 % of the reviewed studies), and stands for randomly  
640 removing a fraction (typically 50 %) of a layer’s output features during the training  
641 process (these output features are set to zero). The core idea of dropout is to ar-  
642 tificially introduce stochasticity to the training process preventing the model from  
643 learning statistical noise in the data.

644 Still, overfitting does not only depend on the number of parameters in the  
645 model, but also on the representatives of the sampling. Particularly in the context  
646 of vegetation mapping, samples are often taken under limited conditions, while a  
647 model is deployed to further, foreign conditions. The associated risk is therefore  
648 an over-fitting of the model to the situation with limited conditions and repre-  
649 sentatives (e.g, with regard to scene illumination or local vegetation properties).  
650 An obvious solution is a larger amount of training data or covered variation, re-  
651 spectively. To reduce the costs of creating labeled observations, a commonly ap-  
652 plied procedure is to synthetically increase the sample quantity and diversity using

653 **data augmentation** procedures (Chatfield et al. 2014; Krizhevsky et al. 2012).  
654 Data augmentation is the process of producing more samples from existing data  
655 by introducing manipulations them (Shorten et al. 2019). These changes may in-  
656 clude randomly changing the spatial extent of the imagery, e.g., to make a model  
657 more robust for detecting individuals of a plant species with varied sizes. Random  
658 transformations, such as flipping, rotating or translating the imagery, can increase  
659 the generality towards varying sun-azimuth angles and corresponding cast-shadows  
660 (also described as rotational invariance). Random spectral shifts may compensate  
661 for variation in illuminations caused by topography or atmospheric conditions and  
662 may further alleviate data calibration issues or sensor-specific differences. In most  
663 CNN-related studies using LiDAR data, the detection process is not based on the  
664 point cloud, but 2D projections derived from the point cloud (cf. Section 3.5.2).  
665 Here, data augmentation can be performed by varying the viewing geometry prior  
666 to generating the 2D image – also referred to as multi-view-data generation (Jin  
667 et al. 2018; Ko et al. 2018; Su et al. 2015; Zou et al. 2017). The overall effectiveness  
668 of data augmentation is highlighted by the fact that 47 % of the studies used data  
669 augmentation. Fromm et al. (2019) and Safonova et al. (2019) explicitly tested the  
670 effect of data augmentation and found significant improvements for the detection  
671 of tree seedlings and bark beetle-infected trees, respectively.

672 Data augmentation may also be performed by not introducing minor manip-  
673 ulations, but creating new, synthetic observations from the existing data. Gao et  
674 al. (2020) presented an automated procedure for the creation of synthetic images  
675 and labels from original images for detecting weed infestation (*Calystegia sepium*)  
676 in sugar beet fields. Their approach involved the creation of masks for individual  
677 plants from the original images used for cropping and transferring the corresponding  
678 RGB information to other base images. Adding data created from this (very simply  
679 said copy & paste) approach to the original training data indeed increased the pre-  
680 cision from 0.75 to 0.83. For training a CNN for detecting individual tree crowns,  
681 Braga et al. (2020) used the same principle and created synthetic Worldview-3  
682 observations by randomly placing manually-delineated tree crowns on background  
683 tiles.

684 Probably the most elegant framework for generating synthetic data is Genera-  
685 tive Adversarial Networks (GANs). Inspired by game theory, GANs are driven by  
686 the competition of a generator module, creating synthetic data (e.g., images) and  
687 a discriminator module aiming to disambiguate between synthetic and real data  
688 Frid-Adar et al. (2018) and Goodfellow et al. (2014). During training, a GAN,  
689 hence, simultaneously improves on how to synthesize observations from noise and  
690 how to classify them (synthetic vs real data or further classes). At the example of  
691 segmenting weed infestation in crop fields in UAV imagery, Kerdegari et al. (2019)  
692 demonstrated a GAN architecture, composed of a generator and the discriminator  
693 modules with four convolutional layers each. The proposed GAN produced realistic  
694 synthetic visual and near-infrared scenes. Moreover, it was demonstrated that using  
695 the discriminator module for semantic segmentation of unknown images resulted in  
696 comparable accuracy to a pure CNN - even when using only 50 % of the available  
697 labels. The fact that the discriminator was originally trained to detect another  
698 problem, i.e. differentiating synthetic from real data, suggest that applying this  
699 trained discriminator to real world problems could also be considered as a form of  
700 transfer learning - an approach discussed in more detailed in the next chapter.

### 701 3.2.1.2 Transfer learning and backbones

702 As described earlier, training data for vegetation attributes is often limited as its ac-  
703 quisition is commonly costly and limited by accessibility. Furthermore, the training  
704 itself is often associated with high computing costs.

705 A common practice to alleviate this problem is to apply **transfer learning**  
706 during CNN model training. Transfer learning includes **pre-training** of the CNN  
707 model on other, presumably very large and heterogeneous datasets. Such datasets  
708 do not necessarily have to include the target metric or class (e.g. a certain plant  
709 species) and can, for instance, be derived from public and generic databases. Popu-  
710 lar examples are the image databases *MSCOCO* or *ImageNet*, which contain thou-  
711 sands of images from various objects, such as cars, buildings, or people. A very  
712 elegant approach of transfer learning is to built on pre-trained models directly, com-  
713 monly referred to as **pre-trained backbone**, which can potentially reduce data  
714 storage and processing costs.

715 The principle of transfer learning can be transcribed as the process where very  
716 generic images, not necessarily belonging to vegetation-related situations, are used  
717 to teach the CNN the ability to *see* in a general sense. The subsequent step of  
718 adjusting the network can be understood as teaching the CNN how to apply the  
719 ability to see to a very specific problem, such as the differentiation of certain plant  
720 species.

721 There exist various transfer learning approaches (Pires de Lima et al. 2020;  
722 Too et al. 2019; Tuia et al. 2016), which can be roughly grouped into two primary  
723 strategies: The shallow strategy adopts very general, lower-level image features such  
724 as edge detectors from the pre-trained backbone or the generic training dataset.  
725 Only the last layers of the CNN are then fine-tuned for higher level and task-  
726 specific features using imagery corresponding to the specific problem (e.g. plant  
727 species detection). The deep strategy, in contrast, involves fine-tuning the entire  
728 network, i.e. start back-propagation with all layers on the pre-trained network.

729 The use of pre-trained backbones is restricted to available architectures. Yet,  
730 backbones can be customized with output layers (e.g. to apply it on regression  
731 or classification problems), cost functions, and other components or integrated in  
732 existing CNNs. There exist a variety of backbones for popular CNN architectures  
733 (cf. Section 3.2.2), such as *VGG*, *ResNet* or *Inception*. It should be noted that the  
734 popular backbones are usually trained on 3-channel (RGB) data, whereas remote  
735 sensing information often provides more predictors, such as multiple bands, time  
736 steps, or sensor types. In this case, band selection or feature reduction algorithms  
737 provide a promising avenue (Rezaee et al. 2018).

738 According to our review, 30.5 % used pre-trained backbones (e.g., Brahimi et al.  
739 (2018), Branson et al. (2018), Fromm et al. (2019), Gao et al. (2020), Mahdianpari  
740 et al. (2018), and Rezaee et al. (2018)). Mehdipour Ghazi et al. (2017) compared  
741 the utility of three backbones based on *GoogLeNet*, *AlexNet*, *VGGNet*, to identify  
742 plant species in photographs. Brahimi et al. (2018) assessed the value of pre-  
743 training for plant disease recognition based on RGB imagery and multiple CNN  
744 architectures. They showed deep pre-training strategy, i.e. back-propagation on all  
745 layers of the pre-trained model, delivered the highest accuracy. The shallow strategy  
746 was usually worse than training a model from scratch. Fromm et al. (2019) showed  
747 that pre-training not always significantly improved the detection of tree seedlings  
748 and that the value of pre-training depends on the network's complexity, while more  
749 shallow architectures are less likely to benefit from pre-training. Mahdianpari et al.

750 (2018) report that full training resulted in better accuracy than fine-tuning existing  
751 backbones trained on *ImageNet*. This suggests that the detection of vegetation  
752 patterns may not necessarily benefit from features learned on generic datasets.  
753 This also agrees with recent research by He et al. (2018) suggesting that transfer  
754 learning may indeed be useful if training data is scarce and computation power  
755 limited, but otherwise an exhaustive training on task-specific data will result in  
756 higher accuracy than using generic datasets.

### 757 3.2.1.3 Weakly- and semi-supervised learning

758 Besides a lack of reference data, it may occur that reference data already exist,  
759 but do not meet the ideal requirements for the intended application. Accordingly,  
760 several concepts and strategies have evolved to compensate for limited availability  
761 or conceptual incompatibilities of reference data.

762 The aim of **Weakly supervised learning** is to decrease costs for human la-  
763 beling or to make use of existing, lower quality reference data. This concept is  
764 particularly interesting for semantic segmentation tasks, where usually an annota-  
765 tion for each sample (point or pixel) is required. Weakly supervised-learning can,  
766 for instance, involve annotations at an image level instead of at a pixel level, or  
767 sparsely annotated data at a pixel level, such as bounding boxes, lines, or points.  
768 Adhikari et al. (2019) applied weakly supervised learning using the principle of  
769 semantic graphics to map crop rows and individual weed plants in rice paddies.  
770 Semantic graphics defines target objects or concepts through abstract forms. Ac-  
771 cordingly, Adhikari et al. (2019) defined crop rows as line features and weeds as  
772 solid circles and showed that an encoder-decoder CNN is capable of accurately  
773 learning and mapping these concepts. Their findings are particularly interesting  
774 because plant rows are rather fuzzy and not clearly delimitable. The higher-level  
775 concept of a row, however, is clearly definable for humans by abstracting the spatial  
776 context of the individual plants and obviously also reproducible by CNNs. The con-  
777 cept of weakly supervised learning is also applicable when explicit ‘ground truth’  
778 is scarce but frequent datasets from other studies exist that come with their own  
779 errors or lower spatial resolutions. Promising results of this approach were pre-  
780 sented by Schmitt et al. (2020), who predicted vegetation types with Sentinel data  
781 and used training data derived from MODIS land cover maps at 500m resolution  
782 (this dataset is freely available; SEN12MS, Schmitt et al. (2019)). Using a high  
783 resolution imagery, they demonstrated that the Sentinel-based predictions reached  
784 even higher accuracy than the datasets used for training. Another variant of weakly  
785 supervised learning for semantic segmentation is based on saliency maps. The basis  
786 for this approach is a CNN trained for image classification, which can be analyzed  
787 through class activation mapping (cf. Section 3.6.2 and Fig. 1 showing an example  
788 for tree species) to identify those pixels that are decisive for assigning an image<sub>j</sub>  
789 to a class<sub>i</sub>. These pixels are then used to segment the target class<sub>i</sub> based on the  
790 assumption that these pixels highlight the components of the respective class in the  
791 image<sub>j</sub> (e.g., the canopy of a tree species). Although no study has been published  
792 to date that has applied this approach to vegetation remote sensing, the potential  
793 has been demonstrated several times in other disciplines (Lee et al. 2019; K. Li  
794 et al. 2018). This approach could, hence, provide a promising way for an efficient  
795 and automatic segmentation (e.g., of plant species) based on large image databases  
796 without spatially explicit labels, such as the *iNaturalist* data.

797 **Semi-supervised learning** describes the training of a model with only a small  
798 number of reference data and, hence, can be located between supervised and unsu-  
799 pervised learning. Weinstein et al. (2019) applied semi-supervised learning frame-  
800 work for detecting single tree crowns in airborne imagery using a two-step approach:  
801 The first step, which can be considered as unsupervised or weakly-supervised learn-  
802 ing, involved training a CNN with labels (bounding boxes,  $n = 435,551$ ) derived  
803 automatically from LiDAR data and a tree crown segmentation algorithm (Roussel  
804 et al. 2017). In the second step, the CNN was optimized using a few hand-annotated  
805 samples derived from the airborne imagery ( $n = 2,848$ ). Thereby, Weinstein et al.  
806 (2019) demonstrated that only few high-quality samples may be required for train-  
807 ing a robust CNN.

808 However, the number of samples required for a specific task is difficult to es-  
809 timate in advance. In this regard, **Active learning**, which can be considered as  
810 a special case of supervised learning, can be an efficient solution. Active Learning  
811 describes the iterative optimization of a model by repeatedly adding new reference  
812 data until the predictive accuracy saturates or reaches a desired threshold. Ghosal  
813 et al. (2019) exemplified an active learning approach for sorghum head detection  
814 in UAV imagery. Starting point was a single image together bounding boxes of  
815 sorghum heads to train a CNN, which was then applied to another random image.  
816 The image and predictions were afterward fed into an annotation app in which a  
817 human interpreter corrected the predictions before they were added to the training  
818 dataset. The initial model was then optimized using the enlarged training dataset  
819 and the entire procedure was repeated in multiple iterations. In their case study,  
820 the model accuracy already converged between 5-10 iterations, highlighting the ef-  
821 ficiency of active learning for finding the right balance between costs of human  
822 labeling and model performance.

### 823 **3.2.2 Approaches and architectures**

824 Depending on the components and architecture, CNNs can be implemented in many  
825 different ways, which in turn enables a wide range of different applications in the  
826 field of vegetation remote sensing. CNNs can initially be grouped into **1D-**, **2D-**  
827 and **3D-CNNs**, where the number refers to the dimensions of the kernel. 1D-  
828 CNNs are less often used (8 % of the reviewed studies) since they do not explicitly  
829 consider spatial context and are, hence, primarily applied to analyze optical spectra  
830 or multitemporal data (Annala et al. 2020; Guidici et al. 2017; Kussul et al. 2017;  
831 Liao et al. 2020; Y. Xi et al. 2019; Zhong et al. 2019). Most studies applied 2D-CNNs  
832 (88 %), as these readily exploit spatial patterns in common imagery (e.g., RGB or  
833 multispectral imagery, cf. Fromm et al. (2019), Kattenborn et al. (2020), Milioto  
834 et al. (2017), Neupane et al. (2019), F. H. Wagner et al. (2020), and Weinstein  
835 et al. (2019). The added value of spatial patterns, i.e. of 2D versus 1D-CNNs, was  
836 even demonstrated with relatively coarse-resolution Landsat data (Kussul et al.  
837 2017). 3D-CNNs are rarely used (4 %), but are the means of choice when successive  
838 layers have a directional relationship to be considered (e.g. canopy height profiles,  
839 hyperspectral reflectance, or time-series data, e.g., Ayrey et al. (2018), Barbosa  
840 et al. (2020), Jin et al. (2019), Liao et al. (2020), Lottes et al. (2018), Nezami et al.  
841 (2020), and Zhong et al. (2019)). 2D- and 3D-CNNs can be applied to solve different  
842 problems, including assigning values or classes to entire images, detecting individual  
843 objects within images, segmenting the extent of classes, or simultaneously detecting

844 individual objects and segmenting their extent (Fig. 10b). The major differences,  
845 including the required structure of labels and resulting outputs, are described in  
846 the following sections:

### 847 3.2.2.1 Image classification / regression

848 Image classification is the assignment of a class to an entire image (Fig. 9a).  
849 For example, an image may be assigned to the class *shrub* if at least a fraction  
850 is covered with *Ulex europaeus* or *Sambucus nigra*. Training image classification  
851 or regression-based CNNs requires comparably simple annotations in the form of  
852 class correspondences or continuous values, respectively, for each image. Typical  
853 CNN-architectures for image classification and regression include *VGG*, *ResNet*,  
854 *Inception* or *EfficientNet*. *VGG* uses blocks of consecutive convolutions and non-  
855 linear activations. Between those building-blocks max-pooling with stride of 2  
856 reduces the resolution of the layers. The filter size of the convolution is restricted  
857 to 3x3, leading to less parameters and thus more possible layers. The small filter  
858 size is still common in more recent networks. Finally, some fully connected layers  
859 are added for classifying the output of the building-blocks (Fig. 4). *ResNet* also  
860 consists of building-blocks with consecutive convolutions and activations (Fig. 5)  
861 but with some major difference: First, the depth of the layers is drastically reduced  
862 before the 3x3 convolution with a bottleneck 1x1 convolution. Thus, the number  
863 of parameters is much lower compared to *VGG*, even so *ResNet* has up to 10 times  
864 more layers. Second, to compensate for the *vanishing gradient problem* (cf. Section  
865 3.2.1.1) with such a high number of layers (e.g. 152), skip connection with identity  
866 or convolution shortcuts are introduced. Such skip connections are still used in  
867 the current design, allowing very deep networks. Third, *ResNet* only uses one max  
868 pooling layer. Instead, convolution with stride 2 are used for resolution reduction.  
869 Most modern architectures such as *EfficientNet* also dismiss max-pooling operation  
870 to reduce possible information loss during pooling.

871 A typical procedure to map vegetation patterns in remote sensing imagery with  
872 CNN-based image classification or regression is to subset the original imagery into  
873 regular tiles (e.g., 128 x 128 pixels) on which the model is subsequently applied  
874 (details see Section 3.5.1). This procedure was for instance applied to LiDAR and  
875 airborne imagery to map tree species (Sun et al. 2019) or the detection of forest  
876 types using a combination of high-resolution satellite imagery and LiDAR data (C.  
877 Sothe et al. 2020). Image classification or regression may also be applied to segments  
878 derived from previously applied unsupervised image segmentation methods (dos  
879 Santos Ferreira et al. 2017; Hartling et al. 2019; Ko et al. 2018; T. Liu et al.  
880 2018a). Image regression is used when a continuous quantity is assigned to an  
881 entire tile. For example, (Kattenborn et al. 2020) predicted continuous cover values  
882 [%] of plant species and communities in UAV-based tiles (2-5m) along smooth  
883 vegetation gradients. Yang et al. (2019) and Castro et al. (2020) estimated rice  
884 grain yield and forage biomass in pastures, respectively, from UAV-based tiles.  
885 Barbosa et al. (2020) mapped continuous crop yield on coarser scales based on  
886 satellite data. Ayrey et al. (2018) used regression on airborne LiDAR data to  
887 predict forest biomass and tree density.

### 888 3.2.2.2 Object detection

889 Object detection aims at locating individual occurrences of a class (e.g. trees)  
890 within an image (Fig. 9b). The detection typically includes the localization of  
891 the object center and an approximation of its extent using a simple rectangular  
892 bounding box.

893 Widely applied architectures for object detection are region-based CNNs (*R-*  
894 *CNN*, Girshick et al. (2014)), which involve a two-step approach; region proposals  
895 of the object’s location and extent followed by a classification. *R-CNN* was followed  
896 by two successors, i.e. *Fast R-CNN* (Girshick 2015) and the most widely applied  
897 and efficient *Faster R-CNN* (Ren et al. 2017). The more recent *Faster-R-CNN*  
898 forwards feature maps (often derived using a *VGG*-type backbone) to a region  
899 proposal branch that performs an initial prediction on potential object locations  
900 (also referred to as anchors). These rather rough region proposals are then used  
901 to crop areas of the feature maps as input for a fine-scaled object localization and  
902 classification (Fig. 6).

903 Object detection is suitable for countable things with definable spatial extent  
904 within the field of view. Such conditions are often found in agricultural settings  
905 and accordingly 45 % of the studies related to agriculture apply object detection  
906 techniques, such as locating and counting palm or tree individuals in plantations  
907 (Csillik et al. 2018; Freudenberg et al. 2019), individual maize plants in TLS-point-  
908 clouds of crop fields (Jin et al. 2018) or individual strawberry fruits and flowers  
909 in sub-centimeter UAV-imagery (Chen et al. 2019). The application of object de-  
910 tection in natural environments is less frequent, which can be explained by the  
911 presence of continuous gradients and smooth transitions in species cover, traits,  
912 and communities. In forestry or conservation, only 14 % and 10 % of the studies  
913 used object detection. Examples include the localization of fir trees infested by bark  
914 beetle (Safonova et al. 2019), the mapping of individual tree crowns across several  
915 ecosystems (Weinstein et al. 2020) or the detection of *Cactae* (López-Jiménez et al.  
916 2019).

917 Object detection-based CNNs are typically trained using bounding boxes of de-  
918 sired classes as labels. Several tools exists for a fast annotation of bounding boxes  
919 (see Section 3.1.1). However, a problem with bounding boxes in vegetation analysis  
920 is that they often do not explicitly define vegetation boundaries (vegetation is not  
921 rectangular). This in turn can make validation difficult, as inaccurate reference  
922 data do not allow a final assessment of the prediction (Weinstein et al. 2020, 2019).  
923 From this point of view semantic (Section 3.2.2.3) or instance segmentation (Sec-  
924 tion 3.2.2.4) may be more spatially explicit, but also require more sophisticated  
925 annotations.

### 926 3.2.2.3 Semantic segmentation

927 While image classification and object detection aim to detect the presence or lo-  
928 cation of an object, semantic segmentation aims to delineate the explicit spatial  
929 extent of the target class within the image (Fig. 9c). In contrast to object detec-  
930 tion, semantic segmentation assigns all pixels in an image to a class. It is especially  
931 suited to segment uncountable and amorphous *stuff* (frequently used term to il-  
932 lustrate the contrast to countable *things* (cf. Kirillov et al. (2019))). The training  
933 process is typically based on labels in the form of spatially explicit masks to provide  
934 a class assignment for each single pixel (e.g., absence or presence or species a, b, c).

935 The challenge with semantic segmentation is that CNNs usually include mul-  
936 tiple pooling operations to reveal spatial context in the feature maps derived from  
937 the convolutions and, thereby, spatial reference and detail is initially lost. One  
938 solution often referred to as **patch-based**, is to perform a semantic segmentation  
939 by predicting only values for the center pixel of the input image and iteratively  
940 slide the field of view over the image data until every pixel received a label (Baeta  
941 et al. 2017; Fricker et al. 2019; Kussul et al. 2017; Mahdianpari et al. 2018; Rezaee  
942 et al. 2018; M. Zhang et al. 2018). However, this method requires an individual  
943 prediction for each pixel and is rather inefficient considering that the CNN analy-  
944 ses the neighbouring pixels at the same time anyway. A more elegant and effective  
945 way is to build a semantic segmentation on **fully convolutional networks** (FCN)  
946 as first demonstrated by Long et al. (2015). FCN conserve the spatial reference,  
947 by memorizing the pixels that caused activations in earlier stages of the network  
948 and forwarding it to an output segmentation map (see Fig. 7). This way, FCN  
949 do not only allow detecting the presence of a target class within an image (e.g.,  
950 a species) but also the individual pixels that correspond to the target class. A  
951 more recent and frequently applied architecture for semantic segmentation is the  
952 *U-Net* (named after its 'U'-like shape, Ronneberger et al. (2015)). *U-Net* features  
953 encoder-decoder structure, while the spatial scale is subsequently reduced after con-  
954 secutive pooling operations and again increased in a contracting path (see Fig. 8).  
955 The activations from the contracting path are forwarded using skip connections to  
956 the expanding path to reconstruct the spatial identity. Further commonly applied  
957 CNN-architectures for semantic segmentation are *SegNet* (Badrinarayanan et al.  
958 2017) or *FC-DenseNet* (Jégou et al. 2017). Semantic segmentation is widely used  
959 in several contexts, ranging from mapping of plant species (Fricker et al. 2019) and  
960 plant communities (Kattenborn et al. 2019a; F. Wagner et al. 2019), to mapping  
961 deadwood (Fricker et al. 2019; Jiang et al. 2019). Torres et al. (2020) compared  
962 amongst other architectures *U-Net*, *SegNet*, *FC-DenseNet* for mapping *Dipteryx*  
963 *alata* trees in an urban context. Their results suggest that the segmentation ac-  
964 curacy of the three latter algorithms was quite similar, whereas it was found that  
965 more simpler architectures (e.g., *U-net*) require less effort for model training.

#### 966 3.2.2.4 Instance segmentation

967 Instance segmentation aims at detecting individual *things*, such as individual plants  
968 or plant elements, and segmenting their spatial extent. Instance segmentation may,  
969 hence, be considered as a combination of object detection and semantic segmenta-  
970 tion (Fig. 9d). A few studies used CNN-based object detection and subsequently  
971 applied segmentation techniques, such as region growing in the case of point cloud  
972 data, to detect individuals (Wang et al. 2019). However, here we define instance  
973 segmentation as an end-to-end, CNN-based segmentation of individuals. One of  
974 the most popular algorithms for instance segmentation is *Mask-RCNN* (He et al.  
975 2017); a derivative from *R-CNN* described in section 3.2.2.4. Alike *Faster-RCNN*,  
976 it comprises a two-step approach, including an initial region proposal followed by  
977 the localization and classification of the feature maps, while in the case of *Mask-*  
978 *R-CNN*, the proposed region is subject to a segmentation branch (Fig. 6). Similar  
979 to semantic segmentation, fully connected layers are used to create masks at the  
980 original resolution of the input imagery. Despite the potential utility of instance  
981 segmentation, the literature search only comprised few respective studies; Jin et al.



982 (2019) used instance segmentation to map individual leaves and stems in maize  
983 plants, Braga et al. (2020) delineated individual tree crowns in tropical forests and  
984 Chiang et al. (2020) detected individual dead trees. The rare use of instance seg-  
985 mentation could be explained by the more sophisticated collection of reference data,  
986 which involves both the identification of individuals and delineating their explicit  
987 spatial extent. In an agricultural context, the identification of instances of multiple  
988 classes may often not be necessary, as most tasks are situated in mono-cultures.  
989 Instance segmentation in a forestry or conservation context may often not be ap-  
990 plicable because natural canopies often feature smooth transitions or overlapping  
991 crowns.

### 992 3.3 Geographic and thematic areas of CNN application

993 CNN-based vegetation remote sensing has already been applied in many countries  
994 (see Fig. 11), whereas a large amount of studies were carried out in Europe, USA,  
995 Brazil, and China. The pattern suggests that CNN applications are found in many  
996 of the World’s biomes and are hence applicable for a wide range of vegetation types  
997 and applications.

998 Our literature survey revealed that CNN-based vegetation remote sensing is  
999 applied to a wide spectrum of thematic categories (Fig. 12). A classification of  
1000 the studies into broad categories showed that 44 % of the studies are related to  
1001 agriculture, 26 % of the studies have relevance for both conservation and forestry.  
1002 8 % and 22 % exclusively tackled research questions for forestry and conservation,  
1003 respectively. Within these broad categories, the specific tasks are very diverse (the  
1004 interested reader can find the explicit references of each task in the appendix):

1005 Examples in the context of **agriculture** include the mapping of individual crop  
1006 fields at regional scales using medium and high-resolution satellite data, e.g. coffee  
1007 crop fields (Baeta et al. 2017), rice paddies (M. Zhang et al. 2018), safflower, corn,  
1008 alfalfa, tomatoes, and vineyards (Zhong et al. 2019). Several studies used high-  
1009 resolution imagery from airborne and satellite platforms to map individual plants in  
1010 plantations, e.g. citrus trees, palm trees or bananas (Csillik et al. 2018; Freudenberg  
1011 et al. 2019; W. Li et al. 2017; Mubin et al. 2019; Neupane et al. 2019). Besides  
1012 detecting individual citrus trees, Ampatzidis et al. (2019) quantified their crown  
1013 diameter, health status (NDVI-based), and respective canopy gaps in plantation  
1014 rows. A large share of the studies used imagery with milli- or centimeter pixel size  
1015 acquired terrestrially or from UAVs. A prime example of such detailed input data is  
1016 the detection of weed infestations, e.g., in soybean (dos Santos Ferreira et al. 2017)  
1017 or sugar beet fields (Gao et al. 2020; Milioto et al. 2017; Sa et al. 2018)). Lottes  
1018 et al. (2018) presented an automatic approach for mapping weed infestation in  
1019 imagery acquired by a farming robot equipped with a mechanical actuator that can  
1020 stamp detected weeds into the ground. Adhikari et al. (2019) used subcentimeter  
1021 imagery to map crop lines of rice plants in paddy fields to aid navigation of weeding  
1022 robots for the eradication of weeds (*Panicum miliaceum*). Jin et al. (2018) tested  
1023 the detection and height estimation of individual maize plants. Other studies used  
1024 high-resolution imagery for yield estimation, e.g., based on counting individual  
1025 flowers at sub-centimeter resolution as a proxy for strawberries yield (Chen et al.  
1026 2019), segmenting sorghum panicles (Malambo et al. 2019) or applying CNN-based  
1027 regression for rice grain yield estimation (Yang et al. 2019).

1028 In the **forestry** context, most studies use high-resolution data from UAV or

1029 airborne platforms. Ayrey et al. (2018) used airborne LiDAR data to map forest  
1030 biomass and tree density in temperate forests. Weinstein et al. (2020) tested the  
1031 localization of individual tree crowns (object detection) across ecosystems using  
1032 airborne data. Braga et al. (2020) used very high-resolution satellite data to de-  
1033 lineate individual tree crowns (instance segmentation) in tropical forests. A series  
1034 of studies dealt with the mapping of tree species or genera in forests (Fricker et al.  
1035 2019; Kattenborn et al. 2020; Natesan et al. 2019; Nezami et al. 2020; Pinheiro  
1036 et al. 2020; Schiefer et al. 2020; Trier et al. 2018; Zou et al. 2017) and urban areas  
1037 (dos Santos et al. 2019; Hartling et al. 2019; Torres et al. 2020). Fromm et al.  
1038 (2019) tested the detection of individual conifer seedlings in high resolutions air-  
1039 borne imagery for monitoring of tree regeneration. A substantial interest exists  
1040 towards assessments of forest damage, e.g., caused by wind throw (Hamdi et al.  
1041 2019; Korznikov 2020) or bark beetle infestations (Safonova et al. 2019).

1042 Examples in **conservation** with medium resolution data include the mapping  
1043 of wetland types at regional scales with multispectral Landsat and polarimetric  
1044 RADARSAT-2 data (Mahdianpari et al. 2018; Mohammadimanesh et al. 2019;  
1045 Pouliot et al. 2019). de Bem et al. (2020) mapped deforestation in the Amazon us-  
1046 ing stacked pairs of Landsat imagery from consecutive years. In the context of dry-  
1047 land mapping program by *FAO* (Food and Agriculture Organization of the United  
1048 Nations), (Guirado et al. 2020) mapped tree cover (%) using airborne orthoimagery  
1049 and exemplified that CNN-based mapping outperformed previous assessments by  
1050 *FAO* based on photo-interpretation. Examples for mapping at high spatial reso-  
1051 lution include the mapping of rainforest types and disturbance (F. Wagner et al.  
1052 2019), plant succession stages in a glacier-related chronosequence (Kattenborn et  
1053 al. 2019a), herbaceous and woody invasive species species in several environments  
1054 (Kattenborn et al. 2019a; T. Liu et al. 2018c; Qian et al. 2020), shrub cover (Guirado  
1055 et al. 2017), ecosystem structure-relevant plant communities in the Arctic tundra  
1056 (Langford et al. 2019) or the rehabilitation of native tussock grass (*Lomandra longi-*  
1057 *folia*) after weed eradication campaigns (Hamylton et al. 2020).

### 1058 3.4 Remote sensing platforms

1059 Approximately, 17 % of the studies acquired data from the ground or **terrestrial**  
1060 platforms, including stationary photography (Ma et al. 2019), mobile mapping data  
1061 from *Google Street View* (Barbierato et al. 2020; Branson et al. 2018), farming  
1062 robots (Lottes et al. 2018), and terrestrial laser scanning (e.g., Bingxiao et al. (2020)  
1063 and Wang et al. (2019)). The major part of studies using terrestrial platforms took  
1064 place in an agriculture context with a focus on precision farming.

1065 With 36 %, the largest share of studies assessed in this review used data cap-  
1066 tured from **UAV**. This can be explained as UAV feature two important features;  
1067 they enable to autonomously acquire spatially continuous data with automated  
1068 georeferencing - a feature that recently revolutionized possibilities for fast, flex-  
1069 ible, repeated, and cost efficient remote sensing data acquisition for vegetation  
1070 analysis. At the same time, UAV can be operated at low altitudes capturing veg-  
1071 etation canopies with high spatial detail. High-resolution data acquired by UAV  
1072 and CNN-based pattern analysis provide powerful synergies for spatially continu-  
1073 ous vegetation analysis. Due to the inevitable trade-off of spatial resolution and  
1074 image footprint, a drawback of any high-resolution remote sensing is the limited  
1075 area coverage decreasing the efficiency for vegetation assessments on large scales.

1076 One approach to overcome this limitation is the spatial up-scaling of UAV-based  
1077 vegetation maps with satellite data (Kattenborn et al. 2019b), where UAV-based  
1078 maps are used as a reference for coarse-resolution but large-scale satellite-based  
1079 predictions.

1080 Depending on the spatial scale of the vegetation analysis and the size of the de-  
1081 cisive spatial features, **airplanes** may feature a more efficient compromise between  
1082 area coverage and resolution. 11 % of the studies in this review used airborne sen-  
1083 sors. In addition to increased spatial coverage, an advantage of airplane platforms  
1084 is their increased potential payload supporting more sophisticated and high-quality  
1085 sensors. Accordingly, a large proportion of airplane-related studies used LiDAR or  
1086 hyperspectral data or a combination of both.

1087 Aerial data from UAV and airplanes are often generated by matching single  
1088 frames from imaging sensors in concert with photogrammetric processing tech-  
1089 niques. Due to the relatively low height of both platforms, the single image frames  
1090 usually feature a substantial variation in viewing geometry and bidirectional re-  
1091 flectance effects. At first sight, this may challenge the retrieval of vegetation char-  
1092 acteristics, but as T. Liu et al. (2018a,c) have shown, this variation can also be  
1093 a valuable source for increasing the amount of training data and generating more  
1094 robust models. In a case study on mapping vegetation types in UAV imagery,  
1095 they demonstrated increasing model performance when using a multi-view approach  
1096 that combined tiles from orthoimagery and the spatially corresponding single image  
1097 frames.

1098 In total 35 % of the studies used data acquired from **satellites**. The poten-  
1099 tial of CNN-based pattern recognition combined with the unprecedented amount  
1100 of high-resolution satellite data was demonstrated by Brandt et al. (2020) who  
1101 mapped more than 1.8 billion trees across the Sahara and Sahel zone with a mosaic  
1102 of 11,128 satellite scenes (GeoEye-1, WorldView-2, WorldView-3 and QuickBird-  
1103 2). This pioneering study suggest how high resolution data from small satellites  
1104 (weight < 500 kg) and microsattellites (weight < 100 kg) will offer ground braking  
1105 opportunities for CNN-based vegetation analysis. Examples are the *Planet Labs*  
1106 constellation of *PlanetScope* data, which image the entire Earth Surface on a daily  
1107 basis at 3.7 m resolution or *SkySat*, which enable to image targeted areas at 0.72  
1108 m resolution. These satellite constellations may provide sufficient spatial detail for  
1109 various large-scale CNN-based vegetation assessments.

### 1110 **3.5 Sensors, spatial and spectral resolution**

1111 CNN are most frequently applied on passive optical sensors (RGB, multispectral,  
1112 or hyperspectral). Only a few studies (7 %) used products from SAR systems.  
1113 Passive optical and SAR data are commonly analyzed with raster-based methods  
1114 and, hence, discussed together in Section 3.5.1. The second-largest share of studies  
1115 (10 %), incorporated LiDAR data, whereas 3 % used terrestrial LiDAR data, and  
1116 7 % used airborne LiDAR. The common methods for the analysis of LiDAR-based  
1117 point clouds are presented in Section 3.5.2. The fusion of multiple sensor types is  
1118 discussed in Section 3.5.3.

### 1119 3.5.1 Passive optical and SAR data analysis

1120 CNNs involve numerous transformations of the input data and the available (mostly  
1121 GPU-based) memory may, hence, limit the maximum size of the input data. How-  
1122 ever, raster data, such as airborne or spaceborne acquisitions from passive optical or  
1123 SAR-sensors, usually feature multiple layers (e.g., bands of different wavelengths or  
1124 multitemporal data) and can, thus, occupy large data volumes. Moreover, for some  
1125 CNN approaches, e.g. image classification, it would not be meaningful to make  
1126 a single prediction for an entire raster, but, instead, make multiple smaller-scaled  
1127 predictions to reveal the spatial variation within the area covered by the raster.  
1128 For these reasons, CNN training and inference is not performed on entire rasters  
1129 but instead on equally sized **tiles** extracted from a raster. The trained CNN can  
1130 then be used to create spatial maps using a **sliding window** principle. Thereby,  
1131 the CNN is applied to regularly extracted tiles that have the same size as the tiles  
1132 used for training.

1133 The most efficient approach is the seamless extraction of tiles without overlap,  
1134 whereas combining the results of multiple, overlapping tiles may be useful to in-  
1135 crease redundancy and compensate for edge effects (Brandt et al. 2020; Du et al.  
1136 2020). Similarly, Neupane et al. (2019) showed that combining the tiling results  
1137 from different orthophotos acquired at multiple resolutions enhances the detection  
1138 of palm trees. Generally, the tile sized should be maximized as determined by  
1139 memory capacities, as larger sizes increase the CNN’s field of view and, hence, am-  
1140 plifies the available spatial context and thus accuracy of the model. This effect was  
1141 demonstrated in (Kattenborn et al. 2020), where the accuracy in estimating the  
1142 cover of plant species and communities from UAV imagery increased considerably  
1143 from smaller (2m) to larger tile sizes (5m). Likewise, at the example of predicting  
1144 crop yield from UAV imagery, Nevavuori et al. (2019) demonstrated that larger  
1145 tile sizes (10, 20, 40m) resulted in more accurate predictions. Especially for very  
1146 high-resolution data, it should also be considered that increasing the tile size can  
1147 furthermore decrease the effect spatially inaccurate reference data (e.g., geolocation  
1148 errors of in-situ data or inaccurately delineated masks or bounding boxes). How-  
1149 ever, in the case of image regression or classification (Section 3.2.2.1), which results  
1150 in a single prediction per tile, increasing the tile size decreases the spatial grain of  
1151 the mapping output (Kattenborn et al. 2020). For segmentation approaches (Sec-  
1152 tion 3.2.2.3), the spatial extent of the input tiles will have no effect on the output  
1153 resolution. The processing speed of the sliding window approach can be enhanced  
1154 by first pre-filtering areas of the target raster using a region proposal. For instance,  
1155 in the context of shrub cover segmentation in arid areas, Guirado et al. (2017) used  
1156 brightness thresholds and edge-detectors, as these are already a good indicator to  
1157 show the general occurrence of shrubs.

1158 In addition to the spatial context or tile size, the **spatial resolution** is a de-  
1159 cisive factor. The spatial resolution most strongly varies with the remote sensing  
1160 platform (Fig. 13) and additionally depends on operating altitude and sensor prop-  
1161 erties. Although CNN applications are designed for pattern analysis, the highest  
1162 possible resolution will not ultimately be the most operational solution, as higher  
1163 resolution comes with increased storage and computation loads. In addition, data  
1164 acquisition at higher spatial resolution leads to smaller area coverage. The ideal  
1165 spatial resolution is determined by the spatial scale at which the characteristic pat-  
1166 terns of the target class or quantity occur. For instance in the context of tree species  
1167 mapping, Schiefer et al. (2020) showed decreasing the spatial resolution from 2 to 8

1168 cm decreases the accuracy (F-score) by at least 25 %. Fromm et al. (2019) showed  
1169 that the detection accuracy for tree seedlings based on different UAV-image reso-  
1170 lutions (0.3-6.3 cm) can vary up to 20 %. Similarly, Neupane et al. (2019) found  
1171 a 17% decrease in 17 detection accuracy for banana palms in plantations when  
1172 decreasing the pixel size from 40 cm to 60 cm. Weinstein et al. (2020) assessed  
1173 the relationship between object size and spatial resolution the other way around.  
1174 They did not change the spatial resolution of the remote sensing data, but ana-  
1175 lyzed different ecosystems with characteristic tree sizes and concluded that treetop  
1176 detection for small trees (in alpine forests) was the least accurate.

1177 Regarding **spectral resolution** of passive optical sensors, the literature search  
1178 revealed that with 52 % the largest share of studies used RGB imagery, whereas  
1179 only 31 % used multispectral (defined as RGB and at least one additional band)  
1180 and only 9 % used hyperspectral data (defined here as > 20 spectral bands). The  
1181 fact that multispectral, and hyperspectral data are less frequently used is not a  
1182 surprise; multispectral and hyperspectral sensors feature larger pixel sizes as nar-  
1183 rower spectral bands receive less radiation and given that the amount of radiation  
1184 received by the sensor must clearly surpass its signal to noise ratio. Accordingly,  
1185 everything else being equal, multispectral and hyperspectral sensors have a lower  
1186 spatial resolution than RGB data. As CNNs are particularly designed for pattern  
1187 analysis RGB data may often be preferred.

1188 Accordingly, the results of several studies suggest, that for many tasks no high-  
1189 spectral-resolution information may be needed: For instance, Osco et al. (2020)  
1190 found that counting citrus trees did not clearly improve when combining multi-  
1191 spectral with RGB data. Zhao et al. (2019) found no improvement in using mul-  
1192 tispectral over RGB data for rice damage assessments (rice lodging). Yang et al.  
1193 (2019) showed that the added value of multispectral on top of RGB information  
1194 only slightly improved the estimation accuracy of rice grain yield. In the context of  
1195 tree species classification, Nezami et al. (2020) did not report clear improvements  
1196 in using UAV-based hyperspectral data over RGB data. Similarly, Kattenborn et  
1197 al. (2019a) showed that CNN-based species identification is more accurate than a  
1198 pixel-based hyperspectral classification of plant species (Kattenborn et al. 2019c;  
1199 Lopatin et al. 2019).

1200 Yet, for several fields of application spectral data may be absolutely necessary.  
1201 For instance, analysis related to chemical constituents in plant tissue, e.g. as a  
1202 proxy for plant health status or plant diseases (Zarco-Tejada et al. 2019, 2018) may  
1203 not be possible without sufficient spectral information as biochemistry particularly  
1204 changes absorption properties and not patterns.

1205 Finally, it should be noted that high spectral and spatial resolution can also be  
1206 combined. For example, pan-sharpening algorithms, such as *local-mean variance*  
1207 *matching* or *Gramm-Schmidt spectral sharpening*, can be used to sharpen coarser  
1208 multi-spectral bands with spatially high-resolution imagery. Such pan-sharpening  
1209 algorithms are often applied to imagery from very high-resolution satellite sen-  
1210 sors that feature a panchromatic band, as for instance WorldView, QuickBird, or  
1211 Pleiades (cf. Braga et al. (2020), Hartling et al. (2019), Korznikov (2020), and W.  
1212 Li et al. (2017)). Recently, more sophisticated pan-sharpening algorithms based on  
1213 CNNs were proposed (Masi et al. (2016) and Yuan et al. (2018), see also Section  
1214 3.5.3).

1215 SAR backscatter is known to be particularly sensitive to vegetation 3D-structure  
1216 and therefore has a great potential for differentiating vegetation types and growth

1217 forms. The fact that microwaves penetrate clouds makes it especially suitable for  
1218 extracting continuous temporal features and large scale assessments and. Accord-  
1219 ingly, SAR data was most frequently used as input for CNN for land cover and  
1220 vegetation type mapping “Comparing Deep Learning and Shallow Learning for  
1221 Large-Scale Wetland Classification in Alberta, Canada” (2019), Liao et al. (2020),  
1222 and Mohammadimanesh et al. (2019). Although SAR data have been used overall  
1223 relatively rarely so far in combination with CNNs, it can be assumed that CNNs  
1224 are excellently suited to unravel the relatively complex SAR signals and will thus  
1225 play a major role in Earth observation in the long term (see Zhu et al. (2020) for a  
1226 review on analyzing SAR data with deep learning).

### 1227 3.5.2 LiDAR-based point cloud analysis

1228 The analysis of spatial point clouds is basically more computationally intensive than  
1229 for raster data since there is no spatial discretization (and thus no normalization)  
1230 in cells, which often results in larger data sets and more complex spatial repre-  
1231 sentations. A strategy to increase the processing speed is to run the analysis on  
1232 subsets of point clouds, for instance, by detecting key features of the target plant,  
1233 which are then used as seeds to apply region growing algorithms. This approach  
1234 was for instance applied for detecting individual maize plants *Zea parviglumis* (Jin  
1235 et al. 2018) and rubber plants *Hevea brasiliensis* (Wang et al. 2019). The most  
1236 frequently applied strategy to handle point clouds (mostly terrestrial LiDAR) is  
1237 the conversion to simpler and discrete feature representations prior to the CNN  
1238 analysis, including 3D voxels or 2D projections (e.g. depth maps) (Jin et al. 2018;  
1239 Ko et al. 2018; Windrim et al. 2020; Zou et al. 2017).

1240 **Voxels** are volumetric representations of the point cloud that are defined by  
1241 regular and non-overlapping 3D cube-like cells. During the conversion of point  
1242 clouds to voxel datasets, a voxel is created in a delimitable area (x,y,z) if it contains  
1243 one or a minimum number of points. Voxels can be analyzed in a similar way  
1244 as multi-layered rasters, where a layer corresponds to an elevation section of the  
1245 original point cloud. Jin et al. (2019) used a 0.4 cm voxel space with terrestrial  
1246 LiDAR data to separate leaves and stems from individual maize plants. Ayrey et al.  
1247 (2018) used  $25 \times 25 \times 33cm$  voxels created from airborne LiDAR-based point clouds  
1248 to map forest properties, whereas each voxel was assigned the number of points it  
1249 included.

1250 *Projections* are 2D representations of the point cloud from a certain position  
1251 (x,y,z) and viewing angle (azimuth, zenith). The projections can be created by  
1252 different spatial or spectral criteria, e.g. as depth maps, prior extracted 3D-metrics  
1253 describing the local neighbourhood, intensity, or color information (Jin et al. 2018;  
1254 Ko et al. 2018; Zou et al. 2017). For airborne LiDAR data, projections are com-  
1255 monly created using nadir view, for instance, to extract digital height models or to  
1256 extract height percentiles. For terrestrial LiDAR, projections are typically created  
1257 using oblique viewing angles. The transformation from TLS-based point clouds to  
1258 depth images (2D) is usually applied multiple times using different viewing geome-  
1259 tries, which can be considered as a form of data augmentation (see section 3.2.1.1).  
1260 For instance, to train a CNN for detecting individual maize plants in TLS point  
1261 clouds, Jin et al. (2018) created 32 2D-projections with varying oblique angles.

1262 Despite such possibilities to decrease computation load, it has to be considered  
1263 that projections or voxel representations of the point cloud will result in a loss of

1264 the original spatial detail. Therefore, it may be desirable to use end-to-end learn-  
1265 ing directly with the raw point cloud data as input. Using the raw point clouds  
1266 instead of voxel or projections may be more computationally demanding but it can  
1267 be assumed that ongoing developments in processing and algorithms will advance  
1268 capabilities to harness point clouds directly. Another challenge is that point clouds  
1269 are unordered sets of vectors (in contrast to elements in raster layers) and their  
1270 analysis requires a spatial invariance with respect to rotations and translations. A  
1271 well-known CNN architecture that considers these challenges is *PointNet*, which,  
1272 hence, enables efficient end-to-end learning on point clouds. The foundation of  
1273 *PointNet* are symmetric functions to ensure permutation invariance with regard  
1274 to the unordered input and transforms the data into a canonical feature space to  
1275 ensure spatial invariance. Even though *PointNet* or similar algorithms have been  
1276 used comparatively rarely so far, the results are very promising: Jin et al. (2020)  
1277 applied *PointNet* to detect ground points under dense forest canopies and found  
1278 greater accuracy than for traditional non-deep learning methods. Briechele et al.  
1279 (2020) tested *PointNet* to classify temperate tree species in UAV LiDAR data and  
1280 reported an overall accuracy of up to 90 %. Bingxiao et al. (2020) and Windrim  
1281 et al. (2020) used modified versions of *PointNet*, which besides point coordinates  
1282 also considers the LiDAR return intensity, and demonstrated high accuracy in dif-  
1283 ferentiating woody elements and foliage for multiple coniferous and deciduous tree  
1284 species (up to 93-96 % overall accuracy). The results of the aforementioned studies  
1285 are especially remarkable, considering that this approach performs a classification  
1286 at the highest possible detail, i.e. at the level of individual points.

### 1287 3.5.3 Sensor and data fusion

1288 Multimodal remote sensing analysis or data fusion is the combination of acquisitions  
1289 of different sensors types (LiDAR, SAR, passive optical). The different character-  
1290 istics of the sensor types result in different sensitivities towards plant properties:  
1291 Passive optical data is largely shaped by absorption and scattering properties at  
1292 the top of the canopy. SAR signals are composed of directional scattering processes  
1293 originating in a few centimeters or even meters depth in the canopy (depending on  
1294 the wavelength). LiDAR measures backscattered radiation of commonly very small  
1295 footprints enabling to look deep into plant canopies. These different sensing modes  
1296 can hence reveal different plant characteristics and their synergistic use can be used  
1297 to harness complementary information.

1298 A conceptually rather simple fusion approach is to merge the resulting predic-  
1299 tions of multiple, dataset-specific CNNs. This can, for instance, be done by  
1300 majority voting (Baeta et al. 2017) or by probabilistic approaches, such as Condi-  
1301 tional Random Fields (Branson et al. 2018). However, this way only the output  
1302 space is combined, but not the features contained in the different data sources so  
1303 that their synergies cannot be directly integrated and exploited. Therefore it is  
1304 usually more expedient to simultaneously integrate the different data sources in a  
1305 single neural network - also known as **feature level fusion**. Feature level fusion  
1306 requires either preprocessing of the data or an adaption of the CNN architectures  
1307 to comply with different data structures (e.g. point cloud vs. raster data), sensing  
1308 modalities such (e.g., viewing angles from oblique SAR vs. nadir passive optical  
1309 acquisitions).

1310 A frequently used approach for feature level fusion is converting and normalizing

1311 the spatial dimensions of the different sensor products and a subsequent **stacking**  
1312 **to a common tensor**. Based on this tensor, a CNN can be applied to simultane-  
1313 ously extract features from both data sources. This approach is easy to implement  
1314 and most frequently applied. For instance, Trier et al. (2018) stacked hyperspec-  
1315 tral data with normalized Digital Surface Models (also referred to as canopy height  
1316 model) for classifying tree species. Hartling et al. (2019) stacked LiDAR intensities,  
1317 hyperspectral and panchromatic bands for tree species classification in urban ar-  
1318 eas. Prior to applying the CNN, they also used the LiDAR data extract tree crown  
1319 segments by height. In the context of large scale mapping of vegetation types in  
1320 the arctic, Langford et al. (2019) stacked multiple satellite products, including high  
1321 spatial resolution SPOT data, high spectral resolution EO1-Hyperion data and a  
1322 height model derived from SAR-interferometry. In the context of mapping crop  
1323 cover types, Liao et al. (2020) stacked multi-temporal polarimetric RADARSAT-2  
1324 SAR data with VEN $\mu$ S multispectral data using a 1D-CNN. While multispectral  
1325 was superior to SAR data, combining multi-temporal SAR data with multispectral  
1326 data increased the model performance. Kattenborn et al. (2019a, 2020), Nezami et  
1327 al. (2020), and Sothe et al. (2020) used UAV imagery for mapping plant species and  
1328 stacked RGB orthoimagery and canopy height models (CHM) derived from pho-  
1329 togrammetric processing pipelines. Interestingly, Nezami et al. (2020) found minor  
1330 improvements when using CHM information for UAV-based tree species classifica-  
1331 tion. (Kattenborn et al. 2020; Sothe et al. 2020) found that CHM information does  
1332 not significantly improve the accuracy, whereas Kattenborn et al. (2020) suggested  
1333 that at these high spatial resolutions the information represented by the CHM is  
1334 already indirectly visible in the orthoimagery itself through shadows and illumina-  
1335 tion differences. In contrast, at the example of coarser-resolution satellite imagery  
1336 and forest type classification, C. Sothe et al. (2020) reported that stacking LiDAR-  
1337 derived canopy height information with pan-sharpened Worldview-2 contributed  
1338 important information.

1339 Overall, these studies demonstrated that merging the different data sources into  
1340 a single tensor can potentially facilitate the extraction of complementary signals  
1341 through convolutions. This approach is easy to implement as it does not require  
1342 manipulating common CNN structures. However, stacking datasets may not be  
1343 ideal as the normalization to a common tensor may introduce a critical loss of the  
1344 original the information, e.g., by converting point clouds to coarse voxels or depth  
1345 maps (cf. Section 3.5.2), or the viewing geometries and acquisitions modes may  
1346 not be directly compatible, e.g., oblique SAR vs. nadir optical data. Instead of  
1347 fusing datasets through a common tensor, it may, therefore, be more advantageous  
1348 to process the different data sources in parallel branches and perform a **feature**  
1349 **concatenation** at a later stage in the network; that is linking the activations or  
1350 feature maps derived from multiple, sensor- or data-specific CNN. These networks  
1351 are also referred to as **multi-stream networks**. At the example of mapping  
1352 rice grain yield from UAV imagery, (Yang et al. 2019) applied a concatenation  
1353 of feature maps resulting from two CNN branches, namely RGB imagery with  
1354 high and multispectral imagery with low spatial resolution, respectively. A prime  
1355 example on how feature concatenation enables to integrate different data types and  
1356 structures was presented by Branson et al. (2018), who classified tree species in an  
1357 urban environment by concatenating a branch fed with nadir airborne RGB imagery  
1358 and a branch fed with multiple *Google Street View* scenes extracted with varying  
1359 viewing angles and zoom levels. Lottes et al. (2018) used feature concatenations for



1360 detecting crop plants and weed infestations in image sequences taken by a farming  
1361 robot. Their approach takes into account that planting patterns in agricultural  
1362 fields (e.g. row structures) provide additional spatial information for differentiating  
1363 crops from weeds. Accordingly, their approach included the parallel segmentation of  
1364 successive image frames using encoder-decoder CNN structures and the subsequent  
1365 concatenation of the resulting feature maps.

1366 Barbosa et al. (2020) compared data fusion based on both stacking datasets and  
1367 feature concatenation for crop yield mapping based on heterogeneous input data,  
1368 including remote sensing reflectance and elevation data and in-situ maps on nitro-  
1369 gen, seed rate, and soil electroconductivity. They tested multi-stream approaches  
1370 with branches being concatenated at an early stage and a later stage in the network,  
1371 that is before and after applying fully connected layers, respectively. The best per-  
1372 formance was achieved with a concatenation after fully connected layers, followed  
1373 by a feature concatenation at an earlier stage in the network. The worst perfor-  
1374 mance was found when stacking all predictors before applying the CNN, which was  
1375 attributed to a sometimes complex relationship among different input datasets.

1376 Another noteworthy application of multi-stream networks is CNN-based pan-  
1377 sharpening, i.e. the process of fusing high spectral information from the coarser-  
1378 resolution bands with high spatial resolution information. Pan-sharpening is fre-  
1379 quently applied to data from very high-resolution satellites as these are often  
1380 equipped with pan-chromatic bands that have wider spectral bandwidths enabling  
1381 an increased sensitivity for incoming radiance and thus higher spatial resolution  
1382 than the other bands with narrower bandwidths. The fusion of spatial and spectral  
1383 information requires the representation of highly complex and non-linear relation-  
1384 ships - an application for which CNN are ideally suited (C. Dong et al. 2016; Yuan  
1385 et al. 2018). A case study on this seminal technique was presented by Brook et al.  
1386 (2020), who used a multi-scale pan-sharpening algorithm (Yuan et al. 2018) to fuse  
1387 both multispectral and -temporal information from Sentinel-2 satellite data with  
1388 the high spatial information from UAV-imagery at the centimetre scale. The cor-  
1389 responding case study demonstrated that this approach can reveal the temporal  
1390 variation of leaf biochemical status of individual vineyard rows.

1391 It should be noted that multitemporal analysis (e.g., change detection, time  
1392 series analysis) can also be considered as feature level fusion. As discussed in more  
1393 detail in Section 3.5.4, multitemporal analysis can be performed using both of the  
1394 above presented modes, that is **stacking** multirate inputs (de Bem et al. 2020)  
1395 or **concatenating** them in multiple CNN branches operating in parallel (Branson  
1396 et al. 2018; Mazzia et al. 2019).

### 1397 3.5.4 Multi-temporal analysis

1398 Almost all plant life is subject to seasonal variation as a consequence of reoccurring  
1399 changes of abiotic factors, such as radiation driving photosynthesis, temperature  
1400 controlling its efficiency or water input providing the primary oxidation source.  
1401 The seasonal phases or dynamics, also known as phenology, of plants is expressed  
1402 through biochemical and structural properties which in turn determine how plants  
1403 are represented in remote sensing data. This implies that temporal variation in  
1404 plant traits can limit the transferability of our models through time. At the same  
1405 time, temporal dynamics can also provide essential information for plant character-  
1406 ization, e.g. phenological features such as flowers revealing the taxonomic identity

1407 and or the length of the growing season as an essential factor for productivity and  
1408 yield.

1409 A few studies assessed model performances based on comparing or combining  
1410 multitemporal datasets. For instance, Ma et al. (2019) assessed the biomass esti-  
1411 mation with subcentrimetre imagery in wheat crops across 17 acquisition dates and  
1412 found a strong variation in accuracy ( $R^2$  0.60-0.89) highlighting that timing can  
1413 play an important role. Rezaee et al. (2018) successfully tested the transferability  
1414 of a CNN for wetland segmentation on a *RapidEye* scene that was not included in  
1415 the training process. Yang et al. (2019) tested the transferability of CNN models  
1416 across time for rice grain yield estimation, in terms of how good a CNN trained on  
1417 one or multiple phenological phases is applicable to a phenological phase it has not  
1418 seen before. As expected, the models became better the more times were consid-  
1419 ered in the training process. Similarly, M. Zhang et al. (2018) showed that stacking  
1420 multivariate Landsat scenes increased the accuracy of segmenting rice paddies.

1421 In the context of satellite-based land cover classification, (Mazzia et al. 2019)  
1422 incorporated spatial patterns of temporal dynamics by concatenating the pixel-wise  
1423 branches of **recurrent neural networks (RNNs)**, followed by the subsequent  
1424 application of a CNN. RNNs are a type of deep learning approach to analyse recur-  
1425 ring patterns and are therefore perfectly suitable for multitemporal remote sensing  
1426 analysis (Zhong et al. 2019; Zhu et al. 2017). A primary strength of RNNs is their  
1427 ability to resemble temporal patterns despite the presence of data gaps introduced  
1428 by missing scenes, cloud cover, snow, or artefacts. Similar to CNN for spatial pat-  
1429 terns, RNNs allow for end-to-end analysis of temporal signals and therefore makes  
1430 a heuristic definition and engineering of temporal or phenological metrics obsolete.  
1431 Thus, combining CNNs with RNNs enables an end-to-end processing scheme in  
1432 both the spatial and temporal domain. It can, hence, be assumed that the com-  
1433 bination of RNNs and CNNs will be a milestone for vegetation analysis with time  
1434 series data as for instance derived from satellite constellations (Reichstein et al.  
1435 2019).

1436 In contrast to recurring phases, natural disturbances or anthropogenic impacts  
1437 can also cause acute or gradual, directed changes. Such anomalies in temporal veg-  
1438 etation dynamics may be tracked with **change detection** of remote sensing data.  
1439 de Bem et al. (2020) stacked pairs of Landsat imagery to track deforestation in the  
1440 Amazon rainforest. Compared to earlier change detection approaches, which were  
1441 mostly based on metrics for temporal comparison (e.g., NDVI), the approach used  
1442 here is simple and flexible as it does not require sophisticated pre-processing, such  
1443 as the radiometric cross-calibration of the raw data. A disadvantage is the require-  
1444 ment of training data, such as binary classification of changed and stable areas.  
1445 However, the required number of reference data is not very high as deforestation  
1446 is typically clearly visible in remote sensing imagery, and often institutional data  
1447 can be accessed. de Bem et al. (2020). Another change detection approach was  
1448 presented by Branson et al. (2018), who used multi-date *Google Street View* im-  
1449 agery to detect changes of urban trees. As the viewing geometries are not steady in  
1450 street view imagery, a pixel exact stacking is not possible and accordingly, they con-  
1451 catenated **Siamese CNNs** fed with images from the different time steps. Siamese  
1452 CNNs include identical CNNs that operate in parallel branches (Daudt et al. 2018).  
1453 During training, the weights are shared between the branches, which reduces the  
1454 number of learnable parameters but most importantly secures that both branches  
1455 have the same statistics so that their outputs are comparable. The outputs are

1456 then concatenated into fully connected layers to classify similarity.

## 1457 **3.6 CNN model assessment, understanding, and interpreta-** 1458 **tion**

### 1459 **3.6.1 Numeric evaluation of the predictive performance**

1460 The performance of a CNN model can be determined by different metrics that  
1461 are primarily determined by the model approach (cf. 3.2.2): For CNN-based re-  
1462 gressions, the **coefficient of determination** (R<sup>2</sup>) and the **Root Mean Squared**  
1463 **Error** (RMSE) are the means of choice to quantify the correspondence between pre-  
1464 dictions and reference observations. The majority (91 %) of the studies reviewed  
1465 here performed classification tasks, which can be evaluated with several metrics  
1466 (see Tab.1 for the most ones). The most used and intuitive metric is the **overall**  
1467 **accuracy** (used in 71 % of the reviewed studies), which quantifies the proportion  
1468 of correct predictions.

1469 However, the overall accuracy is prone to bias introduced by class imbalance  
1470 and in such case an accuracy assessment based on **precision**, indicates the perfor-  
1471 mance regarding false positives, and **recall**, sensitive to false negatives, should be  
1472 preferred. The **F-score** is the harmonic mean of precision and recall and provides  
1473 a single metric for the overall model performance that is robust for unsymmetrical  
1474 datasets.

1475 For object detection and instance segmentation, the question is not how well  
1476 is the average agreement of all predicted pixels, but how accurately are individual  
1477 objects or segments detected. Here, an F-score may be strongly biased by object  
1478 size. A metric that is robust against size variation of objects is the **Intersect**  
1479 **over Union (IoU)**, which is the ratio of correctly classified pixels and the total  
1480 amount of pixels per segment. Note that recall is also known as producer’s accuracy  
1481 or sensitivity, precision as user’s accuracy, F-score as dice coefficient, and IoU as  
1482 Jaccard-index.

1483 Despite the standardization of accuracy measures, there are several issues that  
1484 constrain a direct comparison between studies. Firstly, it is hard to compare the  
1485 different approaches, i.e. object detection, semantic segmentation, and instance  
1486 segmentation, as these differ in dimensions and thematic complexity. Secondly, the  
1487 mode of reference data acquisition and quality may greatly constrain the informa-  
1488 tive value of accuracy assessments (cf. in-situ vs. visual interpretation in Section  
1489 3.2.2). Thirdly, the remote sensing data and the site characteristics may differ con-  
1490 siderably among studies. For instance, (Weinstein et al. 2020) demonstrated with  
1491 multiple datasets from the *NEON* project that the detection accuracy of individual  
1492 tree crowns in airborne imagery greatly depends on the site conditions, such as tree  
1493 species composition or crown size distribution. Lastly, albeit a common application  
1494 task (e.g. tree species classification), the definition of the classification problem  
1495 and presence of classes among studies may differ, which in turn greatly limits com-  
1496 parison of different mapping methods. For example, the present literature search  
1497 comprises nine studies on tree species classification, none of which examined the  
1498 same composition of tree species. Clearly, these challenges for comparing different  
1499 studies, e.g., in terms of CNN architectures, highlights the need for free accessible  
1500 datasets for comparative studies (cf. section 4).

1501 Despite the challenges related to comparing the different studies, the literature

1502 review revealed unprecedented predictive accuracy of CNN-based vegetation remote  
1503 sensing approaches (see Fig. 14). For instance, studies that targeted the classifica-  
1504 tion of tree species reported at average an overall accuracy of 89 %. In comparison,  
1505 a review on tree species classification with a focus on shallower machine learning  
1506 methods (e.g. Random Forest or Support Vector Machines) by Fassnacht et al.  
1507 (2016) reported an overall accuracy of 83.5 %. This is particularly interesting, as  
1508 the reviewed studies in Fassnacht et al. (2016) primarily used sophisticated sensors  
1509 (e.g., hyperspectral or LiDAR data or their combination), while a large share (43  
1510 %) of the CNN-based studies assessed here used merely RGB data. The overall  
1511 superior performance of CNNs compared to shallower machine learning algorithms  
1512 was demonstrated in several studies and applications tasks (Ayrey et al. 2018; Bar-  
1513 bosa et al. 2020; Briechle et al. 2020; de Bem et al. 2020; L. Dong et al. 2020; dos  
1514 Santos Ferreira et al. 2017; Guidici et al. 2017; Hartling et al. 2019; Knauer et al.  
1515 2019; Liao et al. 2020; T. Liu et al. 2018a,b; Mazzia et al. 2019; Mohammadimanesh  
1516 et al. 2019; Rezaee et al. 2018; Y. Xi et al. 2019; M. Zhang et al. 2018; Zhong et al.  
1517 2019)

### 1518 3.6.2 Understanding and interpretation: Opening the *black box*

1519 Assessing the functioning of a model is important to compare and improve algo-  
1520 rithms, to test causal or physical consistency as well as to trust in and learn from  
1521 models. Transferred to CNNs, this may involve the identification and visualization  
1522 of individual pixels, patterns, or even higher-level concepts that contribute to the  
1523 decision-making process. It is often claimed that deep learning and especially CNN  
1524 models are a black box and it is difficult to grasp the basis on which a CNN makes a  
1525 decision (Reichstein et al. 2019). This can be explained as on one hand, many peo-  
1526 ple are not yet familiar with the principle of the still quite new CNN algorithms and  
1527 on the other hand by the incomparable depth and number of parameters of these  
1528 models. However, most CNNs have a linear and clear structure (mostly consecutive  
1529 sequences of repetitive structures) and the basic operations, such as pooling or acti-  
1530 vation functions, are relatively simple. Despite the abundance of parameters, these  
1531 properties facilitate a converting of abstract vectors into interpretable information  
1532 and understanding of CNN internal processes. CNN interpretation can be grouped  
1533 into two branches, i.e. feature visualisation and feature attribution. **Feature vi-**  
1534 **sualization** is centered on the model and aims to reveal what the network or parts  
1535 of it are looking for by simulating synthetic outputs. **Feature attribution** is cen-  
1536 tered on input data and aims to identify which features in the data activate the  
1537 network in a particular way.

1538 An example of **feature visualization** for tree species mapping is given in Fig-  
1539 ure 1, where the functioning of individual convolutions was visualized using gradient  
1540 ascent-based approach. This technique starts by manipulating a blank image (or  
1541 any other input format) using the gradient ascent, a function that identifies local  
1542 maxima so that the values assigned to the output pixels maximizes the activation  
1543 of the network or a particular layer. The resulting layers, therefore, reflect the pat-  
1544 terns that the network has learned as decisive patterns in the training process (see  
1545 also Schiefer et al. 2020). Feature visualization can hence inform about the general  
1546 behaviour of the model, whereas this branch of understanding CNNs already offers  
1547 a variety of different approaches (cf. Olah et al. (2017) for a comprehensive and  
1548 interactive summary on feature visualization techniques).

1549 A limitation of feature visualization is that the synthetic outputs are often  
1550 unnatural and abstract and it can be very challenging to link these outputs to real-  
1551 world features such as plant organs or canopy forms as seen in remote sensing data.  
1552 Moreover, feature visualization primarily focuses to reveal the general behaviour  
1553 model of a model, e.g., *what are relevant patterns for separating tree species?*, but a  
1554 question at hand could be much more specific, such as *On the basis of which plant*  
1555 *characteristics visible in the image, did the model distinguish the fir tree from the*  
1556 *surrounding spruce trees?*

1557 In this regard, **feature attribution** may enable to analyze CNN models in a  
1558 more intuitive and traceable way as it is directly based on the input data.

1559 The common products of feature attribution are so-called **activation maps**,  
1560 also known as sensitivity, saliency, or pixel attribution maps, which typically rep-  
1561 resent how the input data activates individual feature layers within the network in  
1562 form of heatmaps (see Fig. 1). Activation maps are obtained by forward propa-  
1563 gating individual input images (e.g. through a trained CNN (similar procedures  
1564 are also applied for point cloud data, cf. (B. Zhang et al. 2019))). Mohammadi-  
1565 manesh et al. (2019) for instance derived activation maps of a CNN for classifying  
1566 wetland types in order to visualize characteristic backscatter features of different  
1567 SAR polarization. Moreover, they applied the Uniform Manifold Approximation  
1568 and Projection (UMAP, McInnes et al. (2018)) algorithm, a non-linear dimension  
1569 reduction technique, on the activation maps derived from the last layers of multi-  
1570 ple CNN architectures to compare their ability to discriminate the wetland types.  
1571 Despite their demonstrated value, activation maps in their simplest form are only  
1572 input-specific and not output-specific, so they do not inform how an activation  
1573 contributes to a decision (e.g. predicting a class affiliation).

1574 An output-specific procedure is given by **gradient weighted class activation**  
1575 **mapping (Grad-CAM)**, which distils class-specific gradients to coarsely localizes  
1576 the spatial regions of the last convolutional layer that are discriminative towards  
1577 the network output (Selvaraju et al. 2019). However, tracing class activations to  
1578 input features can be limited, since common CNNs usually involve several pooling  
1579 operations so that the last convolutional layer of a network and corresponding  
1580 activation maps have a much lower spatial resolution than the original input data. A  
1581 fine-grained representation of decisive image features can be obtained by combining  
1582 **Grad-CAM** with guided backpropagation, known as **guided Grad-CAM** in case  
1583 of classifications (Selvaraju et al. 2019), which allows tracing the activation of the  
1584 last convolutional layer to the individual pixels of the input image (see Fig. 1 for an  
1585 example on tree species). The feature attribution at the pixel-level can be further  
1586 enhanced by averaging multiple activation maps generated with stochastic noise, as  
1587 proposed in the **SmoothGrad** approach (Smilkov et al. 2017). Most approaches  
1588 for feature attributions target on classification problems, but similar principles were  
1589 also tested for regression problems, such as regression activation mapping (RAM,  
1590 Z. Wang et al. (2017)).

1591 Although the above-mentioned methods for CNN interpretation are already es-  
1592 tablished in other scientific fields, their application in vegetation remote sensing  
1593 seems to be still in its infancy (but see Castro et al. 2020; Schiefer et al. 2020).  
1594 Nevertheless, according to the demonstrated potential in other disciplines, it can be  
1595 assumed that feature attribution will play an important role in the future: Feature  
1596 attribution can be harnessed to test for model shortcomings, such as non-causal  
1597 relationships and artifacts and as a basis for optimizing CNN architectures and

1598 training processes. Moreover, feature attribution provides an interesting avenue  
1599 for weakly-supervised learning (cf. Section 3.2.1.3), where class activation maps  
1600 derived from a CNN trained with coarse training data (e.g., presence and absence  
1601 instead of detailed masks) can be used as a proxy to segment classes at the pixel  
1602 level (Lee et al. 2019; K. Li et al. 2018). Lastly, it stands to reason that the  
1603 extraction and preparation of insights from artificial intelligence will increase our  
1604 knowledge and capabilities towards technical aspects ranging from sensor develop-  
1605 ment and data acquisition, biophysical and ecological understanding, as well as the  
1606 interrelationship of remote sensing signals and vegetation properties.

## 1607 4 Concluding remarks and future perspectives

1608 The primary findings of the present review can be summarised as follows:

- 1609 • The reviewed literature revealed that CNN can greatly advance our capabili-  
1610 ties for remote sensing-based vegetation mapping in conservation, agriculture,  
1611 and forestry sectors. A series of studies reported an increased performance  
1612 of CNNs over shallower machine learning methods. In addition to high ac-  
1613 curacy, CNNs are readily implemented as they support end-to-end learning,  
1614 enabling immediate use of raw data and, hence, making feature engineer-  
1615 ing and pre-processing in many cases obsolete. This will greatly facilitate  
1616 vegetation mapping in the era of Big Data, as the self-learning capabilities  
1617 will allow to more effectively harness the ever growing data streams across  
1618 temporal and spatial scales.
- 1619 • CNNs can be customized for various mapping operations, such as image-  
1620 or tile-based regression and classification (e.g, yield estimation or absence or  
1621 presence of a class), segmenting classes (e.g., a plant species or communities),  
1622 or identifying individual objects and their extents (e.g., single tree of a specific  
1623 species). Due to phenology and the biochemical and structural diversity  
1624 of plant life, remote sensing of vegetation benefits from multitemporal and  
1625 multimodal remote sensing like no other land cover. Combining multiple  
1626 sensors, perspectives or acquisition dates has often been a technical challenge,  
1627 whereas the modularity of deep learning frameworks facilitates to combine  
1628 data with varying dimensions and will, hence, enable to further exploit the  
1629 diversity of earth observation data.
- 1630 • The challenges of machine learning were in particular focused on feature en-  
1631 gineering (*what should a model see*). The new challenge is to design the  
1632 learning procedure (*how should a model learn to see*). Designing and imple-  
1633 menting an effective CNN architecture requires both technical knowledge on  
1634 deep learning principles in concert with process-understanding of the system  
1635 - here, the remotely sensed vegetation signal.
- 1636 • The core of deep learning, gradient descent is an iterative optimization algo-  
1637 rithm and thereby opens efficient, sustainable and elegant ways for model  
1638 training and exchange, including the subsequent optimization of existing  
1639 models with new samples instead of training a new model from scratch, the  
1640 use of backbones to incorporate and channel big data, or federated learn-  
1641 ing, i.e. the distributed training on multiple clients, to combine computing

- 1642 resources and minimize communication costs (*bringing the code to the data,*  
1643 *instead of the data to the code*).
- 1644 • Exposing CNNs to representative and ample reference data is often a bottle-  
1645 neck for achieving high predictive accuracy and generalization. For reasons of  
1646 efficiency and data compatibility, ground-based reference data is rarely used,  
1647 whereas most studies use visual interpretation or the combination of both.  
1648 Various tools and concepts have been developed to efficiently label remote  
1649 sensing data using visual interpretation or ancillary data, while concepts  
1650 such data augmentation, generation of synthetic training data or semi- and  
1651 weakly supervised learning enable to harness even small quantities or inaccur-  
1652 ate training data. It seems obvious that the success of further capturing the  
1653 seemingly infinite variation of the plant world using deep learning and specif-  
1654 ically CNN techniques will be stimulated by free access to remote sensing  
1655 and reference data and the establishment of corresponding open databases.  
1656 Pooling resources in joint databases will foster a sustainable and effective  
1657 benchmarking of CNN algorithms and building transferable and accurate  
1658 models.
  - 1659 • Most studies reviewed here were related to classification problems, such as  
1660 mapping taxonomic identities, land cover types or functional groups. How-  
1661 ever, many vegetation-related properties are of a continuous nature, for which  
1662 reference data acquisition is usually quite expensive (e.g., biochemical or  
1663 structural plant traits). For many tasks, effective CNN-based vegetation re-  
1664 mote sensing will require creative approaches that go beyond traditional su-  
1665 pervised modelling procedures, including weakly- and semi-supervised learn-  
1666 ing approaches that link remote sensing observations with non-remote sensing  
1667 databases (e.g., plant trait observations or forestry variables), with process-  
1668 based models (e.g., radiative transfer models or forest growth simulators) or  
1669 incorporate citizen science data (e.g., plant photographs).
  - 1670 • For several vegetation-related applications fields, CNN's strength in exploit-  
1671 ing spatial patterns could foster paradigm shifts in the utility of remote sens-  
1672 ing sensors and platforms. A series of studies reported success in locating  
1673 and identifying plant species or individuals by means of simple RGB informa-  
1674 tion and, therefore, highlighted that for a variety of vegetation assessments,  
1675 where previously expensive and complex sensors seemed necessary (e.g. hy-  
1676 perspectral data), more easily available data can now be sufficient. CNN  
1677 techniques are, hence, likely to facilitate the realization of cost-efficient and  
1678 powerful remote sensing solutions for a wide range of users. At the same  
1679 time, the hunger of CNN for spatial detail is likely to catalyse the utility  
1680 of high-resolution remote sensing data, in particular microsatellites, off-the-  
1681 shelf rotary or fixed-wing UAVs as well as terrestrial and airborne LiDAR  
1682 data.
  - 1683 • Contrary to common preconceptions that CNN models are a *black box*, multi-  
1684 ple approaches enable a representation and visualization of a trained model,  
1685 including its behaviour and the key patterns that contribute to decision mak-  
1686 ing process. The respective feature visualization and attribution methods are  
1687 essential to understand CNN models and trust them. The greatest chance of  
1688 these methods, however, lies in distilling new knowledge with regard to the

1689 interaction of vegetation and its relationship with remote sensing signals, but  
1690 particularly towards the diversity of plant form and function.

## 1691 **5 Additional resources on CNN theory, implemen-** 1692 **tation and data sources**

### 1693 **Acquire new reference data**

- 1694 • with geocoding in a GIS-environment: *QGIS* (open source, <https://qgis.org/>) or
- 1695 *ArcGIS* (commercial). ArcGIS supports advanced feature for creating polygons,
- 1696 such as easy tablet and styles support and autocompletion functions.
- 1697 • without geocoding using annotation tools: *LabelMe* ([http://labelme.csail.mit.](http://labelme.csail.mit.edu/Release3.0/)
- 1698 [edu/Release3.0/](http://labelme.csail.mit.edu/Release3.0/)), *LabelImg* (<https://github.com/tzutalin/labelImg>), *Labelbox* (<https://github.com/labelbox/labelbox>)
- 1699 [//github.com/labelbox/labelbox](https://github.com/labelbox/labelbox))
- 1700 • *cleanlab*: Machine learning-oriented *Python* package for identifying erroneous la-
- 1701 bels in datasets and learning with noisy labels ([https://github.com/cgnorthcutt/](https://github.com/cgnorthcutt/cleanlab)
- 1702 [cleanlab](https://github.com/cgnorthcutt/cleanlab))

### 1703 **Use existing reference data**

- 1704 • *NEON*: Partly multitemporal airborne LiDAR, RGB, multi- and hyperspectral
- 1705 acquisitions with in-situ reference data on various ecosystems in the US ([https://](https://data.neonscience.org/)
- 1706 [data.neonscience.org/](https://data.neonscience.org/)).
- 1707 • *EuroSat*: Image patches (64x64 @ 10m resolution) from Sentinel-2 radiance data
- 1708 labelled with vegetation types and land cover classes ([https://github.com/phelber/](https://github.com/phelber/eurosat)
- 1709 [eurosat](https://github.com/phelber/eurosat)).
- 1710 • *BigEarth*: Atmospherically corrected Sentinel-2 patches (120x120 @ 10 m resolu-
- 1711 tion) labelled with CORINE land-cover information (<http://bigearth.net/>).
- 1712 • *SEN12MS*: Sentinel-1 and -2 data (256x256 @ 10m resolution) labelled with
- 1713 MODIS-based land-cover information (<https://dataserv.ub.tum.de/s/m1474000>).
- 1714 • *Awesome Public Datasets*: List of topic-centric public data sources from the
- 1715 fields of biology, earth sciences, agriculture. [https://github.com/awesomedata/](https://github.com/awesomedata/awesome-public-datasets)
- 1716 [awesome-public-datasets](https://github.com/awesomedata/awesome-public-datasets)

### 1717 **Compensate for few reference data or missing computational** 1718 **ressources**

- 1719 • Use pre-trained backbones: Many predefined architectures with trained weights
- 1720 (e.g., derived from *ImageNet*, *MSCOCO*) can be loaded directly. A tutorial for using
- 1721 pre-trained backbones with *Keras* can be found at [https://keras.io/guides/transfer\\_](https://keras.io/guides/transfer_learning/)
- 1722 [learning/](https://keras.io/guides/transfer_learning/) and for *PyTorch* at [https://pytorch.org/tutorials/beginner/transfer\\_](https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html)
- 1723 [learning\\_](https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html)
- 1724 [tutorial.html](https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html)
- 1724 • Weakly supervised learning using self organizing maps (SOM, Riese et al. 2020,
- 1725 <https://doi.org/10.3390/rs12010007> and code: <https://doi.org/10.5281/zenodo.2609130>).
- 1726 • Semi-supervised learning with partially unlabelled datasets presented by *Facebook*
- 1727 *AI* in a Pytorch tutorial: [https://pytorch.org/hub/facebookresearch\\_semi-supervised-ImageNet1K-models\\_](https://pytorch.org/hub/facebookresearch_semi-supervised-ImageNet1K-models_resnext/)
- 1728 [resnext/](https://pytorch.org/hub/facebookresearch_semi-supervised-ImageNet1K-models_resnext/)



## 1729 **First steps to CNN implementation**

- 1730 • *FastAI*: Initiative aiming at introducing AI principles to a wide audience (slogan:  
1731 'Making neural nets uncool again') by maintaining a own *Python*-based library  
1732 designed for easy implementation and a wide range of material, courses and tutorials  
1733 (<https://www.fast.ai>)
- 1734 • *Keras* Developer Guides, including help and tutorials on the Keras API and  
1735 getting started with CNN (<https://keras.io/guides/>).
- 1736 • The textbooks *Deep Learning with R* and *Deep Learning with Python* by F.  
1737 Chollet and J.J. Allaire offer a didactically high-quality, catchy and application-  
1738 oriented introduction to *Keras*, including many hands-on sections and sample codes  
1739 (ISBN: 9781617295546 and 9781617294433).
- 1740 • *Deep Learning with Pytorch*: Introduction to the Pytorch framework including a  
1741 CNN-based image classification example ([https://pytorch.org/tutorials/beginner/  
1742 deep\\_learning\\_60min\\_blitz.html](https://pytorch.org/tutorials/beginner/deep_learning_60min_blitz.html))
- 1743 • Documentation on CNN-based land-cover classification of Sentinel-2 satellite data,  
1744 including different training strategies such as fine-tuning and pre-trained networks:  
1745 <https://github.com/jensleitloff/CNN-Sentinel>

## 1746 **Discover CNN architectures**

- 1747 • *Model Zoo*: Documentation and tutorials on various CNN implementations for  
1748 various frameworks (<https://modelzoo.co>).
- 1749 • *Papers With Code*: Database on scientific publications together with correspond-  
1750 ing data and executable code (<https://paperswithcode.com/>).
- 1751 • *Keras* examples for CNN: <https://keras.io/examples/vision/>
- 1752 • *Segmentation Models library*: High-level *Python* API including multiple seg-  
1753 mentation model architectures and backbones for Keras and Tensorflow ([https://  
1754 github.com/qubvel/segmentation\\_models/](https://github.com/qubvel/segmentation_models/)).
- 1755 • *Awesome Semantic Segmentation*: Links list for the most frequently used segmen-  
1756 tation (e.g. U-net) and instance segmentation models (e.g. Mask-R-CNN) for var-  
1757 ious frameworks. The linklist also includes several annotations tools, datasets and  
1758 additional resources (<https://github.com/mrgloom/awesome-semantic-segmentation/>).
- 1759 • *PyTorch Hub*: Out-of-box models with pretrained weights for PyTorch ([https://  
1760 pytorch.org/hub/](https://pytorch.org/hub/)).
- 1761 • *PyTorch Ecosystem Tools*: Tools, libraries, and more for PyTorch, such as fast.ai  
1762 or Detectron2 (<https://pytorch.org/ecosystem/>).
- 1763 • *TensorFlow Hub* (<https://tfhub.dev/>) and *TensorFlow Model Garden* ([https://  
1764 github.com/tensorflow/models](https://github.com/tensorflow/models)) with hundreds of different (pretrained) models  
1765 .

## 1766 **Feature visualization and attribution (What did the CNN 1767 learn?)**

- 1768 • Comprehensive and interactive resource on principles and approaches for CNN  
1769 feature visualizations of imagery <https://distill.pub/2017/feature-visualization/>
- 1770 • *Interpretable Machine Learning* (Molnar 2019): Constantly updated online book  
1771 providing background and guides for making machine learning decisions inter-  
1772 pretable, including a chapter on CNN-based feature visualization ([https://christophm.  
1773 github.io/interpretable-ml-book/](https://christophm.github.io/interpretable-ml-book/)).

- 1774 • Tutorial on visualizing activation maps with Keras: [https://keras.io/examples/vision/visualizing\\_what\\_convnets\\_learn/](https://keras.io/examples/vision/visualizing_what_convnets_learn/)
- 1775
- 1776 • Tutorial on creating saliency maps with the Grad-CAM approach: [https://keras.io/examples/vision/grad\\_cam/](https://keras.io/examples/vision/grad_cam/)
- 1777
- 1778 • *Uniform Manifold Approximation and Projection* (UMAP): A dimension reduction technique useful for deriving abstract representations of feature maps of a CNN
- 1779 to visualize the input data structure or exploring classification and regression performance. <https://umap-learn.readthedocs.io/en/latest/>
- 1780
- 1781
- 1782 • *The What-If Tool* (WIT): Provides an plugins and web interfaces for expanding
- 1783 understanding of a machine learning models allowing the interactive manipulation
- 1784 of labels and models and comparing resulting outcomes (<https://github.com/pair-code/what-if-tool>).
- 1785

## 1786 Acknowledgements

1787 The authors thank Etienne Laliberté, Fabian Fassnacht, Felix Leidinger, and Marco  
1788 Körner for valuable comments and discussions regarding the manuscript.

## 1789 Appendix

### 1790 Methodology of the cluster analysis of terms found in the 1791 review literature using *VOSviewer*

1792 The cluster analysis was performed using VOSviewer (Van Eck et al. (2010), ver-  
1793 sion 1.6.14) and based on the frequency of terms contained in title and abstracts.  
1794 Terms similar in content, synonyms and generic terms to be excluded that are not  
1795 specifically related to the topic were defined in a thesaurus file. The remaining  
1796 terms were included in the cluster analysis if they occurred at least five times. As  
1797 normalization method the *LinLog modularity* was used. The minimum cluster size  
1798 was set to 10.

### 1799 Data on the reviewed literature

1800 The data extracted from the reviewed literature is available as spreadsheet under  
1801 the following URL:  
1802 <https://tinyurl.com/kattenborn-cnn-meta>  
1803 (link to Google Drive; the host/URL will be changed in case of acceptance)

### 1804 Commonly used accuracy metrics for classification and object 1805 detection purposes.

### 1806 Information on the inception module

## 1807 References

1808 Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., & Süsstrunk, S.  
1809 (2012). SLIC superpixels compared to state-of-the-art superpixel

Table 1: Overview and brief introduction of the most frequently used accuracy metrics for classification and object detection purposes.

Metric	Unit	Description / formula
<b>Overall Accuracy (OA)</b>	[0 – 1]	The overall accuracy is the ratio of true predictions (positive and negative) and the total number of observations $OA = \frac{TP + TN}{TP + TN + FP + FN}$
<b>Precision</b> (also known as user's accuracy)	[0 – 1]	Ratio of true presences classified correctly and the number of all positive predictions. Precision assesses how many of the predicted presences are actually true. $precision_i = \frac{TP_i}{TP_i + FP_i}$
<b>Recall</b> (also known as producer's accuracy or sensitivity)	[0 – 1]	Ratio of true presences classified correctly as i and the total number of instances belonging to class i (true positive and false negative). Recall assess how many of the actual presences were classified as true. $recall_i = \frac{TP_i}{TP_i + FN_i}$
<b>F-score</b> (also known as Sørensen-Dice coefficient or Dice similarity coefficient)	[0 – 1]	The F-score is the harmonic mean of recall and precision and, thus, provides a balanced accuracy metric that is sensitive to both under- and overestimation. $F_i = 2 \times \frac{precision_i \times recall_i}{precision_i + recall_i}$
<b>Intersection over Union (IoU)</b> (also known as Jaccard Index)	[0 – 1]	IoU is closely related to the F-score. IoU measures the relative spatial agreement between reference and predicted surfaces (e.g. a segment or bounding box). The intersect is the area shared among both surfaces (Reference AND prediction), whereas the union is the combined area (Reference OR prediction). $IoU_k = \frac{TP_k}{TP_k + FN_k + FP_k}$

- 1810 methods. *IEEE Transactions on Pattern Analysis and Machine In-*  
1811 *telligence*. <https://doi.org/10.1109/TPAMI.2012.120> (cit. on p. 13)
- 1812 Adam, E., Mutanga, O., & Rugege, D. (2010). Multispectral and hyper-  
1813 spectral remote sensing for identification and mapping of wetland  
1814 vegetation: A review. *Wetlands Ecology and Management*, 18(3),  
1815 281–296 (cit. on p. 4).
- 1816 Adhikari, S. P., Yang, H., & Kim, H. (2019). Learning Semantic Graph-  
1817 ics Using Convolutional Encoder-Decoder Network for Autonomous  
1818 Weeding in Paddy. *Frontiers in Plant Science*. [https://doi.org/10.](https://doi.org/10.3389/fpls.2019.01404)  
1819 3389/fpls.2019.01404 (cit. on pp. 20, 25)
- 1820 Ampatzidis, Y., & Partel, V. (2019). UAV-based high throughput phenotyp-  
1821 ing in citrus utilizing multispectral imaging and artificial intelligence.  
1822 *Remote Sensing*, 11(4). <https://doi.org/10.3390/rs11040410> (cit. on  
1823 p. 25)
- 1824 Anderson, C. B. (2018). Biodiversity monitoring, earth observations and the  
1825 ecology of scale. *Ecology Letters*. <https://doi.org/10.1111/ele.13106>  
1826 (cit. on p. 12)
- 1827 Angermueller, C., Pärnamaa, T., Parts, L., & Stegle, O. (2016). Deep learn-  
1828 ing for computational biology. *Molecular Systems Biology*. [https://](https://doi.org/10.15252/msb.20156651)  
1829 [doi.org/10.15252/msb.20156651](https://doi.org/10.15252/msb.20156651) (cit. on p. 4)
- 1830 Annala, L., Honkavaara, E., Tuominen, S., & Pölönen, I. (2020). Chlorophyll  
1831 concentration retrieval by training convolutional neural network for  
1832 stochastic model of leaf optical properties (SLOP) inversion. *Remote*  
1833 *Sensing*, 12(2), 1–22. <https://doi.org/10.3390/rs12020283> (cit. on  
1834 pp. 14, 21)
- 1835 Atzberger, C., Darvishzadeh, R., Schlerf, M., & Le Maire, G. (2013). Suit-  
1836 ability and adaptation of PROSAIL radiative transfer model for hy-  
1837 perspectral grassland studies. *Remote Sensing Letters*, 4(1), 56–65.  
1838 <https://doi.org/10.1080/2150704X.2012.689115> (cit. on p. 3)
- 1839 Ayrey, E., & Hayes, D. J. (2018). The use of three-dimensional convolutional  
1840 neural networks to interpret LiDAR for forest inventory. *Remote*  
1841 *Sensing*, 10(4), 1–16. <https://doi.org/10.3390/rs10040649> (cit. on  
1842 pp. 14, 16, 21, 22, 26, 30, 36)
- 1843 Badrinarayanan, V., Kendall, A., & Cipolla, R. (2017). SegNet: A Deep Con-  
1844 volutional Encoder-Decoder Architecture for Image Segmentation.  
1845 *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
1846 arXiv 1511.00561. <https://doi.org/10.1109/TPAMI.2016.2644615>  
1847 (cit. on p. 24)
- 1848 Baeta, R., Nogueira, K., Menotti, D., & Dos Santos, J. A. (2017). Learning  
1849 Deep Features on Multiple Scales for Coffee Crop Recognition. *Pro-*  
1850 *ceedings - 30th Conference on Graphics, Patterns and Images, SIB-*  
1851 *GRAPI 2017*, 262–268. <https://doi.org/10.1109/SIBGRAPI.2017.41>  
1852 (cit. on pp. 24, 25, 31)

- 1853 Barbierato, E., Bernetti, I., Capecchi, I., & Saragosa, C. (2020). Integrating  
1854 Remote Sensing and Street View Images to Quantify Urban Forest  
1855 Ecosystem Services. *Remote Sensing*, *12*(2), 329. [https://doi.org/](https://doi.org/10.3390/rs12020329)  
1856 [10.3390/rs12020329](https://doi.org/10.3390/rs12020329) (cit. on p. 26)
- 1857 Barbosa, A., Trevisan, R., Hovakimyan, N., & Martin, N. F. (2020). Mod-  
1858 eling yield response to crop management using convolutional neural  
1859 networks. *Computers and Electronics in Agriculture*, *170*, 105197  
1860 (cit. on pp. 21, 22, 33, 36).
- 1861 Bingxiao, W., Wu, B., Zheng, G., & Chen, Y. (2020). An Improved Con-  
1862 volution Neural Network-Based Model for Classifying Foliage and  
1863 Woody Components from Terrestrial Laser Scanning Data. *Remote*  
1864 *Sensing*, *12*, 1010. <https://doi.org/10.3390/rs12061010> (cit. on  
1865 pp. 26, 31)
- 1866 Bonawitz, K., Eichner, H., Grieskamp, W., Huba, D., Ingerman, A., Ivanov,  
1867 V., Kiddon, C., Konecny, J., Mazzocchi, S., McMahan, H. B. Et al.  
1868 (2019). Towards federated learning at scale: System design. *arXiv*  
1869 *preprint arXiv:1902.01046* (cit. on p. 10).
- 1870 Bone, D. J., Bachor, H.-A., & Sandeman, R. J. (1986). Fringe-pattern anal-  
1871 ysis using a 2-D Fourier transform. *Applied Optics*. [https://doi.org/](https://doi.org/10.1364/ao.25.001653)  
1872 [10.1364/ao.25.001653](https://doi.org/10.1364/ao.25.001653) (cit. on p. 8)
- 1873 Braga, J. R. G., Peripato, V., Dalagnol, R., Ferreira, M. P., Tarabalka,  
1874 Y., Aragão, L. E. O. C., & Velho, H. F. D. C. (2020). Tree Crown  
1875 Delineation Algorithm Based on a Convolutional Neural Network.  
1876 *Remote Sensing*, *12*, 1288. <https://doi.org/10.3390/rs12081288> (cit.  
1877 on pp. 18, 25, 26, 29)
- 1878 Brahimi, M., Arsenovic, M., Laraba, S., Sladojevic, S., Boukhalfa, K., &  
1879 Moussaoui, A. (2018). *Deep Learning for Plant Diseases: Detection*  
1880 *and Saliency Map Visualisation*. Springer International Publishing.  
1881 <https://doi.org/10.1007/978-3-319-90403-0>. (Cit. on p. 19)
- 1882 Brandt, M., Tucker, C. J., Kariryaa, A., Rasmussen, K., Abel, C., Small, J.,  
1883 Chave, J., Rasmussen, L. V., Hiernaux, P., Diouf, A. A., Kergoat,  
1884 L., Mertz, O., Igel, C., Gieseke, F., Schöning, J., Li, S., Melocik,  
1885 K., Meyer, J., Sinno, S., . . . Fensholt, R. (2020). An unexpectedly  
1886 large count of trees in the West African Sahara and Sahel. *Nature*,  
1887 *556*3(August 2019). <https://doi.org/10.1038/s41586-020-2824-5>  
1888 (cit. on pp. 27, 28)
- 1889 Branson, S., Wegner, J. D., Hall, D., Lang, N., Schindler, K., & Perona, P.  
1890 (2018). From Google Maps to a fine-grained catalog of street trees.  
1891 *ISPRS Journal of Photogrammetry and Remote Sensing*, *135*, 13–30.  
1892 <https://doi.org/10.1016/j.isprsjprs.2017.11.008> (cit. on pp. 12, 14,  
1893 19, 26, 31–34)
- 1894 Briechle, S., Krzystek, P., & Vosselman, G. (2020). Classification of Tree  
1895 Species and Standing Dead Trees by Fusing Uav-Based Lidar Data  
1896 and Multispectral Imagery in the 3D Deep Neural Network Point-

- 1897 net++. *ISPRS Annals of the Photogrammetry, Remote Sensing and*  
1898 *Spatial Information Sciences*, 5(2), 203–210. [https://doi.org/10.](https://doi.org/10.5194/isprs-annals-V-2-2020-203-2020)  
1899 5194/isprs-annals-V-2-2020-203-2020 (cit. on pp. 31, 36)
- 1900 Brodrick, P. G., Davies, A. B., & Asner, G. P. (2019). Uncovering Ecological  
1901 Patterns with Convolutional Neural Networks. *Trends in Ecology and*  
1902 *Evolution*, 20, 1–12. <https://doi.org/10.1016/j.tree.2019.03.006> (cit.  
1903 on p. 4)
- 1904 Brodu, N., & Lague, D. (2012). 3D terrestrial lidar data classification of  
1905 complex natural scenes using a multi-scale dimensionality criterion:  
1906 Applications in geomorphology. *ISPRS Journal of Photogrammetry*  
1907 *and Remote Sensing*, 68, 121–134. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.isprsjprs.2012.01.006)  
1908 isprsjprs.2012.01.006 (cit. on p. 8)
- 1909 Brook, A., De Micco, V., Battipaglia, G., Erbaggio, A., Ludeno, G., Cata-  
1910 pano, I., & Bonfante, A. (2020). A smart multiple spatial and tempo-  
1911 ral resolution system to support precision agriculture from satellite  
1912 images: Proof of concept on Aglianico vineyard. *Remote Sensing of*  
1913 *Environment*, 240(January), 111679. [https://doi.org/10.1016/j.rse.](https://doi.org/10.1016/j.rse.2020.111679)  
1914 2020.111679 (cit. on p. 33)
- 1915 Cadieu, C. F., Hong, H., Yamins, D. L., Pinto, N., Ardila, D., Solomon, E. A.,  
1916 Majaj, N. J., & DiCarlo, J. J. (2014). Deep Neural Networks Rival  
1917 the Representation of Primate IT Cortex for Core Visual Object  
1918 Recognition. *PLoS Computational Biology*. [https://doi.org/10.1371/](https://doi.org/10.1371/journal.pcbi.1003963)  
1919 journal.pcbi.1003963 (cit. on p. 4)
- 1920 Castro, W., Junior, J. M., Polidoro, C., Osco, L. P., Gonçalves, W., Ro-  
1921 drrigues, L., Santos, M., Jank, L., Barrios, S., Valle, C., Simeão,  
1922 R., Carromeu, C., Silveira, E., Jorge, L. A. d. C., & Matsubara,  
1923 E. (2020). Deep learning applied to phenotyping of biomass in for-  
1924 ages with uav-based rgb imagery. *Sensors (Switzerland)*, 20(17), 1–  
1925 18. <https://doi.org/10.3390/s20174802> (cit. on pp. 14, 22, 37)
- 1926 Chatfield, K., Simonyan, K., Vedaldi, A., & Zisserman, A. (2014). Return  
1927 of the devil in the details: Delving deep into convolutional nets, In  
1928 *Bmvc 2014 - proceedings of the british machine vision conference*  
1929 *2014*. <https://doi.org/10.5244/c.28.6>. (Cit. on p. 18)
- 1930 Chen, Y., Lee, W. S., Gan, H., Peres, N., Fraisse, C., Zhang, Y., & He, Y.  
1931 (2019). Strawberry yield prediction based on a deep neural network  
1932 using high-resolution aerial orthoimages. *Remote Sensing*, 11(13),  
1933 1–21. <https://doi.org/10.3390/rs11131584> (cit. on pp. 23, 25)
- 1934 Chiang, C. Y., Barnes, C., Angelov, P., & Jiang, R. (2020). Deep Learning-  
1935 Based Automated Forest Health Diagnosis from Aerial Images. *IEEE*  
1936 *Access*, 8arXiv 2010.08437, 144064–144076. [https://doi.org/10.](https://doi.org/10.1109/ACCESS.2020.3012417)  
1937 1109/ACCESS.2020.3012417 (cit. on p. 25)
- 1938 Chollet, F. (2017). Xception: Deep Learning with Depthwise Separable Con-  
1939 volutions. *Proceedings of the IEEE conference on computer vision*

- 1940           *and pattern recognition*, 7(3), 1251–1258. [https://doi.org/10.4271/](https://doi.org/10.4271/2014-01-0975)  
1941           2014-01-0975 (cit. on p. 3)
- 1942 Colomina, I., & Molina, P. (2014). Unmanned aerial systems for photogram-  
1943           metry and remote sensing: A review. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.isprsjprs.2014.02.013)  
1944           isprsjprs.2014.02.013. (Cit. on p. 3)
- 1945 Comparing Deep Learning and Shallow Learning for Large-Scale Wetland  
1946           Classification in Alberta, Canada. (2019). *Remote Sensing*, 12(1), 2.  
1947           <https://doi.org/10.3390/rs12010002> (cit. on p. 30)
- 1948 Csillik, O., Cherbini, J., Johnson, R., Lyons, A., & Kelly, M. (2018). Identifi-  
1949           cation of Citrus Trees from Unmanned Aerial Vehicle Imagery Using  
1950           Convolutional Neural Networks. *Drones*, 2(4), 39. [https://doi.org/](https://doi.org/10.3390/drones2040039)  
1951           10.3390/drones2040039 (cit. on pp. 13, 23, 25)
- 1952 Daudt, R. C., Le Saux, B., & Boulch, A. (2018). Fully convolutional siamese  
1953           networks for change detection, In *2018 25th IEEE International Con-*  
1954           *ference on Image Processing (ICIP)*. IEEE. (Cit. on p. 34).
- 1955 David, E., Madec, S., Sadeghi-Tehran, P., Aasen, H., Zheng, B., Liu, S.,  
1956           Kirchgeßner, N., Ishikawa, G., Nagasawa, K., Badhon, M. A. Et  
1957           al. (2020). Global wheat head detection (gwhd) dataset: A large  
1958           and diverse dataset of high resolution rgb labelled images to de-  
1959           velop and benchmark wheat head detection methods. *arXiv preprint*  
1960           *arXiv:2005.02162* (cit. on p. 16).
- 1961 de Bem, P. P., de Carvalho Junior, O. A., Fontes Guimarães, R., & Tran-  
1962           coso Gomes, R. A. (2020). Change Detection of Deforestation in the  
1963           Brazilian Amazon Using Landsat Data and Convolutional Neural  
1964           Networks. *Remote Sensing*, 12(6), 901. [https://doi.org/10.3390/](https://doi.org/10.3390/rs12060901)  
1965           rs12060901 (cit. on pp. 26, 33, 34, 36)
- 1966 Dong, C., Loy, C. C., & Tang, X. (2016). Accelerating the super-resolution  
1967           convolutional neural network. *Lecture Notes in Computer Science*  
1968           *(including subseries Lecture Notes in Artificial Intelligence and Lec-*  
1969           *ture Notes in Bioinformatics)*, 9906 LNCSarXiv 1608.00367, 391–  
1970           407. [https://doi.org/10.1007/978-3-319-46475-6\\_25](https://doi.org/10.1007/978-3-319-46475-6_25) (cit. on p. 33)
- 1971 Dong, L., Du, H., Han, N., Li, X., Zhu, D., Mao, F., Zhang, M., Zheng,  
1972           J., Liu, H., Huang, Z., & He, S. (2020). Application of convolu-  
1973           tional neural network on lei bamboo above-ground-biomass (AGB)  
1974           estimation using Worldview-2. *Remote Sensing*, 12(6), 958. <https://doi.org/10.3390/rs12060958> (cit. on p. 36)
- 1975
- 1976 dos Santos Ferreira, A., Matte Freitas, D., Gonçalves da Silva, G., Pis-  
1977           tori, H., & Theophilo Folhes, M. (2017). Weed detection in soybean  
1978           crops using ConvNets. *Computers and Electronics in Agriculture*,  
1979           143(February), 314–324. [https://doi.org/10.1016/j.compag.2017.](https://doi.org/10.1016/j.compag.2017.10.027)  
1980           10.027 (cit. on pp. 13, 22, 25, 36)
- 1981 dos Santos, A. A., Marcato Junior, J., Araújo, M. S., Di Martini, D. R.,  
1982           Tetila, E. C., Siqueira, H. L., Aoki, C., Eltner, A., Matsubara, E. T.,  
1983           Pistori, H., Feitosa, R. Q., Liesenberg, V., & Gonçalves, W. N.

- 1984 (2019). Assessment of CNN-Based Methods for Individual Tree De-  
1985 tection on Images Captured by RGB Cameras Attached to UAVs.  
1986 *Sensors*, 19(16), 3595. <https://doi.org/10.3390/s19163595> (cit. on  
1987 p. 26)
- 1988 Du, L., McCarty, G. W., Zhang, X., Lang, M. W., Vanderhoof, M. K., Li,  
1989 X., Huang, C., Lee, S., & Zou, Z. (2020). Mapping forested wetland  
1990 inundation in the delmarva peninsula, USA using deep convolutional  
1991 neural networks. *Remote Sensing*, 12(4), 644. [https://doi.org/10.](https://doi.org/10.3390/rs12040644)  
1992 [3390/rs12040644](https://doi.org/10.3390/rs12040644) (cit. on pp. 14, 28)
- 1993 Fasnacht, F. E., Latifi, H., Stere??czak, K., Modzelewska, A., Lefsky, M.,  
1994 Waser, L. T., Straub, C., & Ghosh, A. (2016). Review of studies on  
1995 tree species classification from remotely sensed data. *Remote Sensing*  
1996 *of Environment*, 186arXiv arXiv:1011.1669v3, 64–87. [https://doi.](https://doi.org/10.1016/j.rse.2016.08.013)  
1997 [org/10.1016/j.rse.2016.08.013](https://doi.org/10.1016/j.rse.2016.08.013) (cit. on pp. 3, 12, 36)
- 1998 Flood, N., Watson, F., & Collett, L. (2019). Using a U-net convolutional  
1999 neural network to map woody vegetation extent from high resolution  
2000 satellite imagery across Queensland, Australia. *International Jour-*  
2001 *nal of Applied Earth Observation and Geoinformation*, 82(June),  
2002 101897. <https://doi.org/10.1016/j.jag.2019.101897> (cit. on p. 13)
- 2003 Freudenberg, M., Nölke, N., Agostini, A., Urban, K., Wörgötter, F., &  
2004 Kleinn, C. (2019). Large scale palm tree detection in high resolution  
2005 satellite images using U-Net. *Remote Sensing*, 11(3), 1–18. [https:](https://doi.org/10.3390/rs11030312)  
2006 [//doi.org/10.3390/rs11030312](https://doi.org/10.3390/rs11030312) (cit. on pp. 13, 23, 25)
- 2007 Fricker, G. A., Ventura, J. D., Wolf, J. A., North, M. P., Davis, F. W., &  
2008 Franklin, J. (2019). A convolutional neural network classifier identi-  
2009 fies tree species in mixed-conifer forest from hyperspectral imagery.  
2010 *Remote Sensing*, 11(19), 2326 (cit. on pp. 24, 26).
- 2011 Frid-Adar, M., Klang, E., Amitai, M., Goldberger, J., & Greenspan, H.  
2012 (2018). Synthetic data augmentation using gan for improved liver  
2013 lesion classification, In *2018 ieee 15th international symposium on*  
2014 *biomedical imaging (isbi 2018)*. IEEE. (Cit. on p. 18).
- 2015 Fromm, M., Schubert, M., Castilla, G., Linke, J., & McDermid, G. (2019).  
2016 Automated detection of conifer seedlings in drone imagery using con-  
2017 volutional neural networks. *Remote Sensing*, 11(21). [https://doi.](https://doi.org/10.3390/rs11212585)  
2018 [org/10.3390/rs11212585](https://doi.org/10.3390/rs11212585) (cit. on pp. 15, 18, 19, 21, 26, 29)
- 2019 Gao, J., French, A. P., Pound, M. P., He, Y., Pridmore, T. P., & Pieters,  
2020 J. G. (2020). Deep convolutional neural networks for image-based  
2021 *Convolvulus sepium* detection in sugar beet fields. *Plant Methods*,  
2022 16(1), 1–12. <https://doi.org/10.1186/s13007-020-00570-z> (cit. on  
2023 pp. 18, 19, 25)
- 2024 Gastellu-Etchegorry, J.-P., Demarez, V., Pinel, V., & Zagolski, F. (1996).  
2025 Modeling radiative transfer in heterogeneous 3-d vegetation canopies.  
2026 *Remote sensing of environment*, 58(2), 131–156 (cit. on p. 14).



- 2027 Geng, J., Wang, H., Fan, J., & Ma, X. (2017). Deep supervised and contrac-  
 2028 tive neural network for sar image classification. *IEEE Transactions*  
 2029 *on Geoscience and Remote Sensing*, 55(4), 2442–2459 (cit. on p. 8).
- 2030 Ghosal, S., Zheng, B., Chapman, S. C., Potgieter, A. B., Jordan, D. R.,  
 2031 Wang, X., Singh, A. K., Singh, A., Hirafuji, M., Ninomiya, S., Gana-  
 2032 pathysubramanian, B., Sarkar, S., & Guo, W. (2019). A Weakly Su-  
 2033 pervised Deep Learning Framework for Sorghum Head Detection and  
 2034 Counting. *Plant Phenomics*, 2019, 1–14. [https://doi.org/10.34133/](https://doi.org/10.34133/2019/1525874)  
 2035 2019/1525874 (cit. on p. 21)
- 2036 Girshick, R. (2015). Fast R-CNN, In *Proceedings of the ieee international*  
 2037 *conference on computer vision*. [https://doi.org/10.1109/ICCV.2015.](https://doi.org/10.1109/ICCV.2015.169)  
 2038 169. (Cit. on p. 23)
- 2039 Girshick, R., Donahue, J., Darrell, T., Berkeley, U. C., & Malik, J. (2014).  
 2040 R-CNN. *1311.2524v5*, arXiv 1311.2524. [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR.2014.81)  
 2041 CVPR.2014.81 (cit. on p. 23)
- 2042 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D.,  
 2043 Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial  
 2044 nets, In *Advances in neural information processing systems*. (Cit. on  
 2045 p. 18).
- 2046 Guidici, D., & Clark, M. L. (2017). One-Dimensional Convolutional Neural  
 2047 Network Land-Cover Classification of Multi-Seasonal Hyperspectral  
 2048 Imagery in the San Francisco Bay Area , California. *Remote Sens*,  
 2049 9, 629. <https://doi.org/10.3390/rs9060629> (cit. on pp. 21, 36)
- 2050 Guirado, E., Alcaraz-Segura, D., Cabello, J., Puertas-Ruiz, S., Herrera, F.,  
 2051 & Tabik, S. (2020). Tree Cover Estimation in Global Drylands from  
 2052 Space Using Deep Learning. *Remote Sensing*, 12(3), 343. [https://](https://doi.org/10.3390/rs12030343)  
 2053 [doi.org/10.3390/rs12030343](https://doi.org/10.3390/rs12030343) (cit. on p. 26)
- 2054 Guirado, E., Tabik, S., Alcaraz-Segura, D., Cabello, J., & Herrera, F. (2017).  
 2055 Deep-learning Versus OBIA for scattered shrub detection with Google  
 2056 Earth Imagery: Ziziphus lotus as case study. *Remote Sensing*, 9(12),  
 2057 1–22. <https://doi.org/10.3390/rs9121220> (cit. on pp. 26, 28)
- 2058 Hamdi, Z. M., Brandmeier, M., & Straub, C. (2019). Forest Damage Assess-  
 2059 ment Using Deep Learning on High Resolution Remote Sensing Data.  
 2060 *Remote Sensing*, 11(17), 1976. <https://doi.org/10.3390/rs11171976>  
 2061 (cit. on pp. 13, 26)
- 2062 Hamylton, S., Morris, R., Carvalho, R., Roder, N., Barlow, P., Mills, K.,  
 2063 & Wang, L. (2020). Evaluating techniques for mapping island vege-  
 2064 tation from unmanned aerial vehicle (UAV) images: Pixel classifica-  
 2065 tion, visual interpretation and machine learning approaches. *Inter-  
 2066 national Journal of Applied Earth Observation and Geoinformation*,  
 2067 89(March), 102085. <https://doi.org/10.1016/j.jag.2020.102085>  
 2068 (cit. on p. 26)

- 2069 Haralick, R. M. (1979). Statistical and structural approaches to texture.  
 2070 *Proceedings of the IEEE*, 67(5), 786–804. [https://doi.org/10.1109/](https://doi.org/10.1109/PROC.1979.11328)  
 2071 [PROC.1979.11328](https://doi.org/10.1109/PROC.1979.11328) (cit. on pp. 4, 8)
- 2072 Hartling, S., Sagan, V., Sidike, P., Maimaitijiang, M., & Carron, J. (2019).  
 2073 Urban tree species classification using a worldview-2/3 and liDAR  
 2074 data fusion approach and deep learning. *Sensors (Switzerland)*, 19(6),  
 2075 1–23. <https://doi.org/10.3390/s19061284> (cit. on pp. 8, 15, 22, 26,  
 2076 29, 32, 36)
- 2077 He, K., Girshick, R., & Dollár, P. (2018). Rethinking ImageNet Pre-training.  
 2078 *arXiv preprint*, arXiv 1811.08883, 1–10 (cit. on p. 20).
- 2079 He, K., Gkioxari, G., Dollar, P., & Girshick, R. (2017). Mask R-CNN, In  
 2080 *Proceedings of the IEEE international conference on computer vision*.  
 2081 <https://doi.org/10.1109/ICCV.2017.322>. (Cit. on p. 24)
- 2082 Helber, P., Bischke, B., Dengel, A., & Borth, D. (2019). Eurosat: A novel  
 2083 dataset and deep learning benchmark for land use and land cover  
 2084 classification. *IEEE Journal of Selected Topics in Applied Earth Ob-*  
 2085 *servations and Remote Sensing*, 12(7), 2217–2226 (cit. on p. 16).
- 2086 Hochreiter, S. (1991). Untersuchungen zu dynamischen neuronalen netzen.  
 2087 *Diploma, Technische Universität München*, 91(1) (cit. on p. 17).
- 2088 Hochreiter, S. (1998). The vanishing gradient problem during learning re-  
 2089 current neural nets and problem solutions. *International Journal of*  
 2090 *Uncertainty, Fuzziness and Knowledge-Based Systems*, 6(02), 107–  
 2091 116 (cit. on p. 17).
- 2092 Hoeser, T., & Kuenzer, C. (2020). Object Detection and Image Segmentation  
 2093 with Deep Learning on Earth Observation Data : A Review-Part I  
 2094 : Evolution and Recent Trends. *Remote Sensing*, 12(May), 1667.  
 2095 <https://doi.org/10.3390/rs12101667> (cit. on pp. 3, 4, 7, 10)
- 2096 Huang, B., Lu, K., Audebert, N., Khalel, A., Tarabalka, Y., Malof, J.,  
 2097 Boulch, A., Saux, B. L., Collins, L., Bradbury, K., Lefevre, S., &  
 2098 El-Saban, M. (2018). Large-scale semantic classification: Outcome of  
 2099 the first year of inria aerial image labeling benchmark. *International*  
 2100 *Geoscience and Remote Sensing Symposium (IGARSS), 2018-July*,  
 2101 6947–6950. <https://doi.org/10.1109/IGARSS.2018.8518525> (cit. on  
 2102 p. 3)
- 2103 Jacquemoud, S., Verhoef, W., Baret, F., Bacour, C., Zarco-Tejada, P. J., As-  
 2104 ner, G. P., François, C., & Ustin, S. L. (2009). Prospect+ sail models:  
 2105 A review of use for vegetation characterization. *Remote sensing of*  
 2106 *environment*, 113, S56–S66 (cit. on p. 14).
- 2107 Jégou, S., Drozdal, M., Vazquez, D., Romero, A., & Bengio, Y. (2017).  
 2108 The one hundred layers tiramisu: Fully convolutional densenets for  
 2109 semantic segmentation, In *Proceedings of the IEEE conference on com-*  
 2110 *puter vision and pattern recognition workshops*. (Cit. on p. 24).
- 2111 Jiang, S., Yao, W., & Heurich, M. (2019). Dead Wood Detection Based on  
 2112 Semantic Segmentation of Vhr Aerial Cir Imagery Using Optimized

- 2113 Fcn-Densenet. *ISPRS - International Archives of the Photogramme-*  
 2114 *try, Remote Sensing and Spatial Information Sciences, XLII-2/W16* (September),  
 2115 127–133. [https://doi.org/10.5194/isprs-archives-xlii-2-w16-127-](https://doi.org/10.5194/isprs-archives-xlii-2-w16-127-2019)  
 2116 2019 (cit. on p. 24)  
 2117 test123
- 2118 Jin, S., Su, Y., Zhao, X., Hu, T., & Guo, Q. (2020). A Point-Based Fully  
 2119 Convolutional Neural Network for Airborne LiDAR Ground Point  
 2120 Filtering in Forested Environments. *IEEE Journal of Selected Topics*  
 2121 *in Applied Earth Observations and Remote Sensing*, 13, 3958–3974.  
 2122 <https://doi.org/10.1109/JSTARS.2020.3008477> (cit. on p. 31)
- 2123 Jin, S., Su, Y., Gao, S., Wu, F., Hu, T., Liu, J., Li, W., Wang, D., Chen,  
 2124 S., Jiang, Y., Pang, S., & Guo, Q. (2018). Deep learning: Individual  
 2125 maize segmentation from terrestrial lidar data using faster R-CNN  
 2126 and regional growth algorithms. *Frontiers in Plant Science*, 9(June),  
 2127 1–10. <https://doi.org/10.3389/fpls.2018.00866> (cit. on pp. 18, 23,  
 2128 25, 30)
- 2129 Jin, S., Guan, H., Zhang, J., Guo, Q., Su, Y., Gao, S., Wu, F., Xu, K., Ma, Q.,  
 2130 Hu, T., Liu, J., & Pang, S. (2019). Separating the Structural Com-  
 2131 ponents of Maize for Field Phenotyping Using Terrestrial LiDAR  
 2132 Data and Deep Convolutional Neural Networks. *IEEE Transactions*  
 2133 *on Geoscience and Remote Sensing*, PP, 1–15. [https://doi.org/10.](https://doi.org/10.1109/tgrs.2019.2953092)  
 2134 1109/tgrs.2019.2953092 (cit. on pp. 21, 24, 30)
- 2135 Kaartinen, H., Hyyppä, J., Vastaranta, M., Kukko, A., Jaakkola, A., Yu,  
 2136 X., Pyörälä, J., Liang, X., Liu, J., Wang, Y., Kaijaluoto, R., Melkas,  
 2137 T., Holopainen, M., & Hyyppä, H. (2015). Accuracy of kinematic  
 2138 positioning using global satellite navigation systems under forest  
 2139 canopies. *Forests*. <https://doi.org/10.3390/f6093218> (cit. on p. 12)
- 2140 Kampe, T. U., Johnson, B. R., Kuester, M. A., & Keller, M. (2010). Neon:  
 2141 The first continental-scale ecological observatory with airborne re-  
 2142 mote sensing of vegetation canopy biochemistry and structure. *Jour-*  
 2143 *nal of Applied Remote Sensing*, 4(1), 043510 (cit. on p. 16).
- 2144 Kao, R. H., Gibson, C. M., Gallery, R. E., Meier, C. L., Barnett, D. T.,  
 2145 Docherty, K. M., Blevins, K. K., Travers, P. D., Azuaje, E., Springer,  
 2146 Y. P. Et al. (2012). Neon terrestrial field observations: Design-  
 2147 ing continental-scale, standardized sampling. *Ecosphere*, 3(12), 1–  
 2148 17 (cit. on p. 16).
- 2149 Kattenborn, T., Eichel, J., & Fassnacht, F. E. (2019a). Convolutional Neural  
 2150 Networks enable efficient, accurate and fine-grained segmentation of  
 2151 plant species and communities from high-resolution UAV imagery.  
 2152 *Scientific Reports*, 9(1), 17656. [https://doi.org/10.1038/s41598-](https://doi.org/10.1038/s41598-019-53797-9)  
 2153 019-53797-9 (cit. on pp. 13, 17, 24, 26, 29, 32)
- 2154 Kattenborn, T., Eichel, J., Schmidtlein, S., Wisser, S., Burrows, L., & Fass-  
 2155 nacht, F. E. (2020). Convolutional Neural Networks accurately pre-  
 2156 dict cover fractions of plant species and communities in Unmanned

- 2157 Aerial Vehicle imagery. *Remote Sensing in Ecology and Conserva-*  
2158 *tion*, 1–15. <https://doi.org/10.1002/rse2.146> (cit. on pp. 13, 21, 22,  
2159 26, 28, 32)
- 2160 Kattenborn, T., Lopatin, J., Förster, M., Braun, A. C., & Fassnacht, F. E.  
2161 (2019b). UAV data as alternative to field sampling to map woody  
2162 invasive species based on combined Sentinel-1 and Sentinel-2 data.  
2163 *Remote Sensing of Environment*, 227, 61–73. [https://doi.org/10.](https://doi.org/10.1016/j.rse.2019.03.025)  
2164 [1016/j.rse.2019.03.025](https://doi.org/10.1016/j.rse.2019.03.025) (cit. on p. 27)
- 2165 Kattenborn, T., Lopatin, J., Förster, M., Braun, A. C., & Fassnacht, F. E.  
2166 (2019c). UAV data as alternative to field sampling to map woody  
2167 invasive species based on combined Sentinel-1 and Sentinel-2 data.  
2168 *Remote Sensing of Environment*, 227(January), 61–73. [https://doi.](https://doi.org/10.1016/j.rse.2019.03.025)  
2169 [org/10.1016/j.rse.2019.03.025](https://doi.org/10.1016/j.rse.2019.03.025) (cit. on p. 29)
- 2170 Kattenborn, T., & Schmidtlein, S. (2019d). Radiative transfer modelling  
2171 reveals why canopy reflectance follows function. *Scientific Reports*,  
2172 9(1), 6541. <https://doi.org/10.1038/s41598-019-43011-1> (cit. on  
2173 p. 12)
- 2174 Kelly, L., Suominen, H., Goeuriot, L., Neves, M., Kanoulas, E., Li, D.,  
2175 Azzopardi, L., Spijker, R., Zuccon, G., Scells, H. Et al. (2019).  
2176 Overview of the clef ehealth evaluation lab 2019, In *International*  
2177 *conference of the cross-language evaluation forum for european lan-*  
2178 *guages*. Springer. (Cit. on p. 16).
- 2179 Kerdegari, H., Razaak, M., Argyriou, V., & Remagnino, P. (2019). Smart  
2180 Monitoring of Crops Using Generative Adversarial Networks. *Lecture*  
2181 *Notes in Computer Science (including subseries Lecture Notes in*  
2182 *Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11678  
2183 *LNCS*, 554–563. [https://doi.org/10.1007/978-3-030-29888-3\\_45](https://doi.org/10.1007/978-3-030-29888-3_45)  
2184 (cit. on p. 18)
- 2185 Kirillov, A., He, K., Girshick, R., Rother, C., & Dollar, P. (2019). Panoptic  
2186 segmentation. *Proceedings of the IEEE Computer Society Confer-*  
2187 *ence on Computer Vision and Pattern Recognition, 2019-June*arXiv  
2188 1801.00868, 9396–9405. <https://doi.org/10.1109/CVPR.2019.00963>  
2189 (cit. on p. 23)
- 2190 Knauer, U., von Rekowski, C. S., Stecklina, M., Krokotsch, T., Pham Minh,  
2191 T., Hauffe, V., Kiliass, D., Ehrhardt, I., Sagischewski, H., Chmara,  
2192 S., & Seiffert, U. (2019). Tree Species Classification Based on Hybrid  
2193 Ensembles of a Convolutional Neural Network (CNN) and Random  
2194 Forest Classifiers. *Remote Sensing*, 11(23), 2788. [https://doi.org/](https://doi.org/10.3390/rs11232788)  
2195 [10.3390/rs11232788](https://doi.org/10.3390/rs11232788) (cit. on p. 36)
- 2196 Ko, C., Kang, J., & Sohn, G. (2018). DEEP MULTI-TASK LEARNING  
2197 for TREE GENERA CLASSIFICATION. *ISPRS Annals of the Pho-*  
2198 *togrammetry, Remote Sensing and Spatial Information Sciences*, 4(2),  
2199 153–159. <https://doi.org/10.5194/isprs-annals-IV-2-153-2018> (cit.  
2200 on pp. 18, 22, 30)

- 2201 Korznikov, K. (2020). Automatic Windthrow Detection Using Very-High-  
 2202 Resolution Satellite Imagery and Deep Learning. *Remote Sensing*,  
 2203 12(April), 1145. <https://doi.org/10.3390/rs12071145> (cit. on pp. 26,  
 2204 29)
- 2205 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classi-  
 2206 fication with deep convolutional neural networks, In *Advances in*  
 2207 *neural information processing systems*. [https://doi.org/10.1201/](https://doi.org/10.1201/9781420010749)  
 2208 [9781420010749](https://doi.org/10.1201/9781420010749). (Cit. on pp. 17, 18)
- 2209 Kussul, N., Lavreniuk, M., Skakun, S., & Shelestov, A. (2017). Deep learning  
 2210 classification of land cover and crop types using remote sensing data.  
 2211 *IEEE Geoscience and Remote Sensing Letters*, 14(5), 778–782 (cit.  
 2212 on pp. 21, 24).
- 2213 Langford, Z. L., Kumar, J., Hoffman, F. M., Breen, A. L., & Iversen, C. M.  
 2214 (2019). Arctic vegetation mapping using unsupervised training datasets  
 2215 and convolutional neural networks. *Remote Sensing*, 11(1), 1–23.  
 2216 <https://doi.org/10.3390/rs11010069> (cit. on pp. 26, 32)
- 2217 Lee, J., Kim, E., Lee, S., Lee, J., & Yoon, S. (2019). Ficklenet: Weakly  
 2218 and semi-supervised semantic image segmentation using stochastic  
 2219 inference, In *Proceedings of the ieee conference on computer vision*  
 2220 *and pattern recognition*. (Cit. on pp. 20, 38).
- 2221 Leitão, P. J., Schwieder, M., Pötzschner, F., Pinto, J. R. R., Teixeira, A. M. C.,  
 2222 Pedroni, F., Sanchez, M., Rogass, C., van der Linden, S., Busta-  
 2223 mante, M. M. C., & Hostert, P. (2018). From sample to pixel: multi-  
 2224 scale remote sensing data for upscaling aboveground carbon data in  
 2225 heterogeneous landscapes. *Ecosphere*, 9(8), e02298. [https://doi.org/](https://doi.org/10.1002/ecs2.2298)  
 2226 [10.1002/ecs2.2298](https://doi.org/10.1002/ecs2.2298) (cit. on p. 12)
- 2227 Lepš, J., & Hadincová, V. (1992). *How reliable are our vegetation analyses?*  
 2228 (Tech. rep. No. 1). <https://doi.org/10.2307/3236006>. (Cit. on p. 13)
- 2229 Li, K., Wu, Z., Peng, K.-C., Ernst, J., & Fu, Y. (2018). Tell me where  
 2230 to look: Guided attention inference network, In *Proceedings of the*  
 2231 *ieee conference on computer vision and pattern recognition*. (Cit. on  
 2232 pp. 20, 38).
- 2233 Li, W., Fu, H., Yu, L., & Cracknell, A. (2017). Deep learning based oil  
 2234 palm tree detection and counting for high-resolution remote sensing  
 2235 images. *Remote Sensing*, 9(1). <https://doi.org/10.3390/rs9010022>  
 2236 (cit. on pp. 25, 29)
- 2237 Liao, C., Wang, J., Xie, Q., Al Baz, A., Huang, X., Shang, J., & He,  
 2238 Y. (2020). Synergistic Use of Multi-Temporal RADARSAT-2 and  
 2239 VENUS Data for Crop Classification Based on 1D Convolutional  
 2240 Neural Network CSA SOAR-E View project NSERC discovery View  
 2241 project Synergistic Use of Multi-Temporal RADARSAT-2 and VENUS  
 2242 Data for Crop Classifi. *Remote Sensing*, 12(832), 832. [https://doi.](https://doi.org/10.3390/rs12050832)  
 2243 [org/10.3390/rs12050832](https://doi.org/10.3390/rs12050832) (cit. on pp. 21, 30, 32, 36)

- 2244 Liu, T., & Abd-Elrahman, A. (2018a). Deep convolutional neural network  
 2245 training enrichment using multi-view object-based analysis of Un-  
 2246 manned Aerial systems imagery for wetlands classification. *ISPRS*  
 2247 *Journal of Photogrammetry and Remote Sensing*, *139*, 154–170. <https://doi.org/10.1016/j.isprsjprs.2018.03.006> (cit. on pp. 22, 27, 36)  
 2248
- 2249 Liu, T., Abd-Elrahman, A., Morton, J., & Wilhelm, V. L. (2018b). Compar-  
 2250 ing fully convolutional networks, random forest, support vector  
 2251 machine, and patch-based deep convolutional neural networks for  
 2252 object-based wetland mapping using images from small unmanned  
 2253 aircraft system. *GIScience and Remote Sensing*, *55*(2), 243–264.  
 2254 <https://doi.org/10.1080/15481603.2018.1426091> (cit. on pp. 15,  
 2255 36)
- 2256 Liu, T., Abd-Elrahman, A., Zare, A., Dewitt, B. A., Flory, L., & Smith,  
 2257 S. E. (2018c). A fully learnable context-driven object-based model  
 2258 for mapping land cover using multi-view data from unmanned air-  
 2259 craft systems. *Remote Sensing of Environment*, *216*(June), 328–344.  
 2260 <https://doi.org/10.1016/j.rse.2018.06.031> (cit. on pp. 26, 27)
- 2261 Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks  
 2262 for semantic segmentation, In *Proceedings of the ieee conference on*  
 2263 *computer vision and pattern recognition*. (Cit. on pp. 24, 66).
- 2264 Lopatin, J., Dolos, K., Kattenborn, T., & Fassnacht, F. E. (2019). How  
 2265 canopy shadow affects invasive plant species classification in high  
 2266 spatial resolution remote sensing. *Remote Sensing in Ecology and*  
 2267 *Conservation*, 1–16. <https://doi.org/10.1002/rse2.109> (cit. on p. 29)
- 2268 López-Jiménez, E., Vasquez-Gomez, J. I., Sanchez-Acevedo, M. A., Herrera-  
 2269 Lozada, J. C., & Uriarte-Arcia, A. V. (2019). Columnar cactus recog-  
 2270 nition in aerial images using a deep learning approach. *Ecological In-*  
 2271 *formatics*, *52*, 131–138. <https://doi.org/10.1016/j.ecoinf.2019.05.005>  
 2272 (cit. on p. 23)
- 2273 Lottes, P., Behley, J., Milioto, A., & Stachniss, C. (2018). Fully convolutional  
 2274 networks with sequential information for robust crop and weed de-  
 2275 tection in precision farming. *IEEE Robotics and Automation Letters*,  
 2276 *3*(4), 1806.03412, 2870–2877. [https://doi.org/10.1109/LRA.2018.](https://doi.org/10.1109/LRA.2018.2846289)  
 2277 [2846289](https://doi.org/10.1109/LRA.2018.2846289) (cit. on pp. 21, 25, 26, 32)
- 2278 Lunetta, R. S., Congalton, R. G., Fenstermaker, L. K., Jensen, J. R., McG-  
 2279 wire, K. C., & Tinney, L. R. (1991). Remote sensing and geographic  
 2280 information system data integration: error sources and research is-  
 2281 sues. *Photogrammetric Engineering & Remote Sensing*, *57*(6), 677–  
 2282 687 (cit. on p. 13).
- 2283 Ma, J., Li, Y., Chen, Y., Du, K., Zheng, F., Zhang, L., & Sun, Z. (2019).  
 2284 Estimating above ground biomass of winter wheat at early growth  
 2285 stages using digital images and deep convolutional neural network.  
 2286 *European Journal of Agronomy*, *103*(June 2018), 117–129. <https://doi.org/10.1016/j.eja.2018.12.004> (cit. on pp. 26, 34)  
 2287

- 2288 Mahdianpari, M., Salehi, B., Rezaee, M., Mohammadimanesh, F., & Zhang,  
2289 Y. (2018). Very deep convolutional neural networks for complex land  
2290 cover mapping using multispectral remote sensing imagery. *Remote*  
2291 *Sensing*, 10(7). <https://doi.org/10.3390/rs10071119> (cit. on pp. 19,  
2292 24, 26)
- 2293 Maier, S., Lüdeker, W., & Günther, K. (1999). Slop: A revised version of  
2294 the stochastic model for leaf optical properties. *Remote Sensing of*  
2295 *Environment*, 68(3), 273–280 (cit. on p. 14).
- 2296 Malambo, L., Rooney, W., Zhou, T., Popescu, S., Ku, N.-W., & Moore, S.  
2297 (2019). A deep learning semantic segmentation-based approach for  
2298 field-level sorghum panicle counting. *Remote Sensing*, 11(24). <https://doi.org/10.3390/rs11242939> (cit. on p. 25)
- 2300 Marconi, S., Graves, S. J., Gong, D., Nia, M. S., Le Bras, M., Dorr, B. J.,  
2301 Fontana, P., Gearhart, J., Greenberg, C., Harris, D. J. Et al. (2019).  
2302 A data science challenge for converting airborne remote sensing data  
2303 into ecological information. *PeerJ*, 6, e5843 (cit. on p. 16).
- 2304 Masi, G., Cozzolino, D., Verdoliva, L., & Scarpa, G. (2016). Pansharpening  
2305 by convolutional neural networks. *Remote Sensing*, 8(7). <https://doi.org/10.3390/rs8070594> (cit. on p. 29)
- 2307 Mazzia, V., Khaliq, A., & Chiaberge, M. (2019). Improvement in Land Cover  
2308 and Crop Classification based on Temporal Features Learning from  
2309 Sentinel-2 Data Using Recurrent-Convolutional Neural Network (R-  
2310 CNN). *Applied Sciences*, 10(1), 238. [https://doi.org/10.3390/](https://doi.org/10.3390/app10010238)  
2311 [app10010238](https://doi.org/10.3390/app10010238) (cit. on pp. 33, 34, 36)
- 2312 McInnes, L., Healy, J., & Melville, J. (2018). Umap: Uniform manifold ap-  
2313 proximation and projection for dimension reduction. *arXiv preprint*  
2314 *arXiv:1802.03426* (cit. on p. 37).
- 2315 McRoberts, R. E., & Tomppo, E. O. (2007). Remote sensing support for  
2316 national forest inventories. *Remote Sensing of Environment*. <https://doi.org/10.1016/j.rse.2006.09.034> (cit. on p. 3)
- 2318 Mehdipour Ghazi, M., Yanikoglu, B., & Aptoula, E. (2017). Plant identifi-  
2319 cation using deep neural networks via optimization of transfer learn-  
2320 ing parameters. *Neurocomputing*, 235(April 2016), 228–235. <https://doi.org/10.1016/j.neucom.2017.01.018> (cit. on p. 19)
- 2322 Milioto, A., Lottes, P., & Stachniss, C. (2017). REAL-TIME BLOB-WISE  
2323 SUGAR BEETS VS WEEDS CLASSIFICATION for MONITOR-  
2324 ING FIELDS USING CONVOLUTIONAL NEURAL NETWORKS.  
2325 *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial*  
2326 *Information Sciences*, 4(2W3), 41–48. [https://doi.org/10.5194/](https://doi.org/10.5194/isprs-annals-IV-2-W3-41-2017)  
2327 [isprs-annals-IV-2-W3-41-2017](https://doi.org/10.5194/isprs-annals-IV-2-W3-41-2017) (cit. on pp. 21, 25)
- 2328 Mohammadimanesh, F., Salehi, B., Mahdianpari, M., Gill, E., & Molinier,  
2329 M. (2019). A new fully convolutional neural network for semantic  
2330 segmentation of polarimetric SAR imagery in complex land cover  
2331 ecosystem. *ISPRS Journal of Photogrammetry and Remote Sensing*,

- 2332           151 (March), 223–236. [https://doi.org/10.1016/j.isprsjprs.2019.03.](https://doi.org/10.1016/j.isprsjprs.2019.03.015)  
2333           015 (cit. on pp. 26, 30, 36, 37)
- 2334 Molnar, C. (2019). *Interpretable machine learning: A guide for making black*  
2335           *box models explainable* [[https://christophm.github.io/interpretable-](https://christophm.github.io/interpretable-ml-book/)  
2336           ml-book/]. (Cit. on p. 41).
- 2337 Mubin, N. A., Nadarajoo, E., Shafri, H. Z. M., & Hamedianfar, A. (2019).  
2338           Young and mature oil palm tree detection and counting using convo-  
2339           lutional neural network deep learning method. *International Journal*  
2340           *of Remote Sensing*, 40(19), 7500–7515. [https://doi.org/10.1080/](https://doi.org/10.1080/01431161.2019.1569282)  
2341           01431161.2019.1569282 (cit. on p. 25)
- 2342 Mulla, D. J. (2013). Twenty five years of remote sensing in precision agri-  
2343           culture: Key advances and remaining knowledge gaps. [https://doi.](https://doi.org/10.1016/j.biosystemseng.2012.08.009)  
2344           org/10.1016/j.biosystemseng.2012.08.009. (Cit. on p. 3)
- 2345 Nagendra, H., Lucas, R., Honrado, J. P., Jongman, R. H., Tarantino, C.,  
2346           Adamo, M., & Mairota, P. (2013). Remote sensing for conservation  
2347           monitoring: Assessing protected areas, habitat extent, habitat con-  
2348           dition, species diversity, and threats. *Ecological Indicators*. [https:](https://doi.org/10.1016/j.ecolind.2012.09.014)  
2349           //doi.org/10.1016/j.ecolind.2012.09.014 (cit. on p. 3)
- 2350 Natesan, S., Armenakis, C., & Vepakomma, U. (2019). Resnet-based tree  
2351           species classification using uav images. *International Archives of the*  
2352           *Photogrammetry, Remote Sensing and Spatial Information Sciences*  
2353           - *ISPRS Archives*, 42(2/W13), 475–481. [https://doi.org/10.5194/](https://doi.org/10.5194/isprs-archives-XLII-2-W13-475-2019)  
2354           isprs-archives-XLII-2-W13-475-2019 (cit. on pp. 13, 26)
- 2355 Neupane, B., Horanont, T., & Hung, N. D. (2019). Deep learning based  
2356           banana plant detection and counting using high-resolution red-green-  
2357           blue (RGB) images collected from unmanned aerial vehicle (UAV).  
2358           *PloS one*, 14(10), e0223906. [https://doi.org/10.1371/journal.pone.](https://doi.org/10.1371/journal.pone.0223906)  
2359           0223906 (cit. on pp. 21, 25, 28, 29)
- 2360 Nevavuori, P., Narra, N., & Lipping, T. (2019). Crop yield prediction with  
2361           deep convolutional neural networks. *Computers and Electronics in*  
2362           *Agriculture*, 163(June), 104859. [https://doi.org/10.1016/j.compag.](https://doi.org/10.1016/j.compag.2019.104859)  
2363           2019.104859 (cit. on pp. 14, 28)
- 2364 Nezami, S., Khoramshahi, E., Pölonen, I., Nevalainen, O., Honkavaara, E.,  
2365           Honkavaara@nls, E., & Fi, E. H. (2020). Tree Species Classifica-  
2366           tion of Drone Hyperspectral and RGB Imagery with Deep Learning  
2367           Convolutional Neural Networks Hyperspectral imaging guided skin  
2368           cancer diagnostics View project DroneKnowledge View project So-  
2369           mayeh Nezami Finnish Geodetic Institute Tre. [https://doi.org/10.](https://doi.org/10.20944/preprints202002.0334.v1)  
2370           20944/preprints202002.0334.v1 (cit. on pp. 21, 26, 29, 32)
- 2371 Nguyen, G., Dlugolinsky, S., Bobák, M., Tran, V., Garcia, A. L., Heredia,  
2372           I., Malik, P., & Hluchy, L. (2019). Machine learning and deep learn-  
2373           ing frameworks and libraries for large-scale data mining: A survey.  
2374           *Artificial Intelligence Review*, 52(1), 77–124 (cit. on p. 10).



- 2375 Noack, A. (2007). Energy models for graph clustering. *J. Graph Algorithms*  
2376 *Appl.*, 11(2), 453–480 (cit. on p. 64).
- 2377 North, P. R. (1996). Three-dimensional forest light interaction model using  
2378 a monte carlo method. *IEEE Transactions on geoscience and remote*  
2379 *sensing*, 34(4), 946–956 (cit. on p. 14).
- 2380 Olah, C., Mordvintsev, A., & Schubert, L. (2017). Feature visualization.  
2381 *Distill*, 2(11), e7 (cit. on p. 36).
- 2382 Osco, L. P., de Arruda, M. d. S., Marcato Junior, J., da Silva, N. B., Ramos,  
2383 A. P. M., Moryia, É. A. S., Imai, N. N., Pereira, D. R., Creste, J. E.,  
2384 Matsubara, E. T., Li, J., & Gonçalves, W. N. (2020). A convolutional  
2385 neural network approach for counting and geolocating citrus-trees in  
2386 UAV multispectral imagery. *ISPRS Journal of Photogrammetry and*  
2387 *Remote Sensing*, 160(November 2019), 97–106. [https://doi.org/10.](https://doi.org/10.1016/j.isprsjprs.2019.12.010)  
2388 [1016/j.isprsjprs.2019.12.010](https://doi.org/10.1016/j.isprsjprs.2019.12.010) (cit. on p. 29)
- 2389 Pettorelli, N., Schulte to Bühne, H., Tulloch, A., Dubois, G., Macinnis-  
2390 Ng, C., Queirós, A. M., Keith, D. A., Wegmann, M., Schrodte, F.,  
2391 Stellmes, M., Sonnenschein, R., Geller, G. N., Roy, S., Somers, B.,  
2392 Murray, N., Bland, L., Geijzendorffer, I., Kerr, J. T., Broszeit, S.,  
2393 . . . Nicholson, E. (2017). Satellite remote sensing of ecosystem func-  
2394 tions: opportunities, challenges and way forward. *Remote Sensing in*  
2395 *Ecology and Conservation*, 1–23. <https://doi.org/10.1002/rse2.59>  
2396 (cit. on p. 3)
- 2397 Pinheiro, M., Roberti, D., Almeida, A. D., Almeida, D. D., Baldez, J., Min-  
2398 ervino, S., Franklin, H., Veras, P., Formighieri, A., Alexandre, C.,  
2399 Santos, N., Aurélio, M., & Ferreira, D. (2020). Forest Ecology and  
2400 Management Individual tree detection and species classification of  
2401 Amazonian palms using UAV images and deep learning. *Forest Ecol-*  
2402 *ogy and Management*, 475, 118397. [https://doi.org/10.1016/j.foreco.](https://doi.org/10.1016/j.foreco.2020.118397)  
2403 [2020.118397](https://doi.org/10.1016/j.foreco.2020.118397) (cit. on p. 26)
- 2404 Pires de Lima, R., & Marfurt, K. (2020). Convolutional neural network  
2405 for remote-sensing scene classification: Transfer learning analysis.  
2406 *Remote Sensing*, 12(1), 86 (cit. on p. 19).
- 2407 Pouliot, D., Latifovic, R., Pasher, J., & Duffe, J. (2019). Assessment of  
2408 convolution neural networks for wetland mapping with landsat in  
2409 the central Canadian boreal forest region. *Remote Sensing*, 11(7).  
2410 <https://doi.org/10.3390/rs11070772> (cit. on p. 26)
- 2411 Qian, W., Huang, Y., Liu, Q., Fan, W., Sun, Z., Dong, H., Wan, F., & Qiao,  
2412 X. (2020). UAV and a deep convolutional neural network for moni-  
2413 toring invasive alien plants in the wild. *Computers and Electronics in*  
2414 *Agriculture*, 174(May), 105519. [https://doi.org/10.1016/j.compag.](https://doi.org/10.1016/j.compag.2020.105519)  
2415 [2020.105519](https://doi.org/10.1016/j.compag.2020.105519) (cit. on p. 26)
- 2416 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Car-  
2417 valhais, N., & Prabhat. (2019). Deep learning and process under-  
2418 standing for data-driven Earth system science. *Nature*, 566(7743),

- 2419 195–204. <https://doi.org/10.1038/s41586-019-0912-1> (cit. on pp. 4,  
2420 14, 15, 34, 36)
- 2421 Ren, S., He, K., Girshick, R., & Sun, J. (2017). Faster R-CNN, arXiv  
2422 1506.01497. <https://doi.org/10.1109/TPAMI.2016.2577031> (cit.  
2423 on p. 23)
- 2424 Rezaee, M., Mahdianpari, M., Zhang, Y., & Salehi, B. (2018). Deep Convo-  
2425 lutional Neural Network for Complex Wetland Classification Using  
2426 Optical Remote Sensing Imagery. *IEEE Journal of Selected Topics in  
2427 Applied Earth Observations and Remote Sensing*, 11(9), 3030–3039.  
2428 <https://doi.org/10.1109/JSTARS.2018.2846178> (cit. on pp. 19, 24,  
2429 34, 36)
- 2430 Riese, F. M., Keller, S., & Hinz, S. (2020). Supervised and semi-supervised  
2431 self-organizing maps for regression and classification focusing on hy-  
2432 perspectral data. *Remote Sensing*, 12(1), 7 (cit. on p. 40).
- 2433 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional net-  
2434 works for biomedical image segmentation, In *International confer-  
2435 ence on medical image computing and computer-assisted interven-  
2436 tion*, Cham., Springer. [https://doi.org/10.1007/978-3-319-24574-  
2437 4\\_28](https://doi.org/10.1007/978-3-319-24574-4_28). (Cit. on pp. 3, 17, 24)
- 2438 Roussel, J.-R., Auty, D., De Boissieu, F., & Sánchez Meador, A. (2017).  
2439 Lidr: Airborne lidar data manipulation and visualization for forestry  
2440 applications. r package version 1.2. 0. (Cit. on p. 21).
- 2441 Russell, B. C., Torralba, A., Murphy, K. P., & Freeman, W. T. (2008).  
2442 LabelMe: A database and web-based tool for image annotation. *In-  
2443 ternational Journal of Computer Vision*. [https://doi.org/10.1007/  
2444 s11263-007-0090-8](https://doi.org/10.1007/s11263-007-0090-8) (cit. on p. 13)
- 2445 Sa, I., Popović, M., Khanna, R., Chen, Z., Lottes, P., Liebisch, F., Nieto,  
2446 J., Stachniss, C., Walter, A., & Siegwart, R. (2018). WeedMap: A  
2447 large-scale semantic weed mapping framework using aerial multi-  
2448 spectral imaging and deep neural network for precision farming. *Re-  
2449 mote Sensing*, 10(9), arXiv 1808.00100. [https://doi.org/10.3390/  
2450 rs10091423](https://doi.org/10.3390/rs10091423) (cit. on p. 25)
- 2451 Safonova, A., Tabik, S., Alcaraz-Segura, D., Rubtsov, A., Maglinets, Y., &  
2452 Herrera, F. (2019). Detection of Fir Trees (*Abies sibirica*) Damaged  
2453 by the Bark Beetle in Unmanned Aerial Vehicle Images with Deep  
2454 Learning. *Remote Sensing*, 11(6), 643. [https://doi.org/10.3390/  
2455 rs11060643](https://doi.org/10.3390/rs11060643) (cit. on pp. 18, 23, 26)
- 2456 Schiefer, F., Kattenborn, T., Frick, A., Frey, J., Schall, P., Koch, B., &  
2457 Schmidlein, S. (2020). Mapping forest tree species in high resolution  
2458 uav-based rgb-imagery by means of convolutional neural networks.  
2459 *ISPRS Journal of Photogrammetry and Remote Sensing*, 170, 205–  
2460 215. <https://doi.org/https://doi.org/10.1016/j.isprsjprs.2020.10.015>  
2461 (cit. on pp. 13, 26, 28, 36, 37)

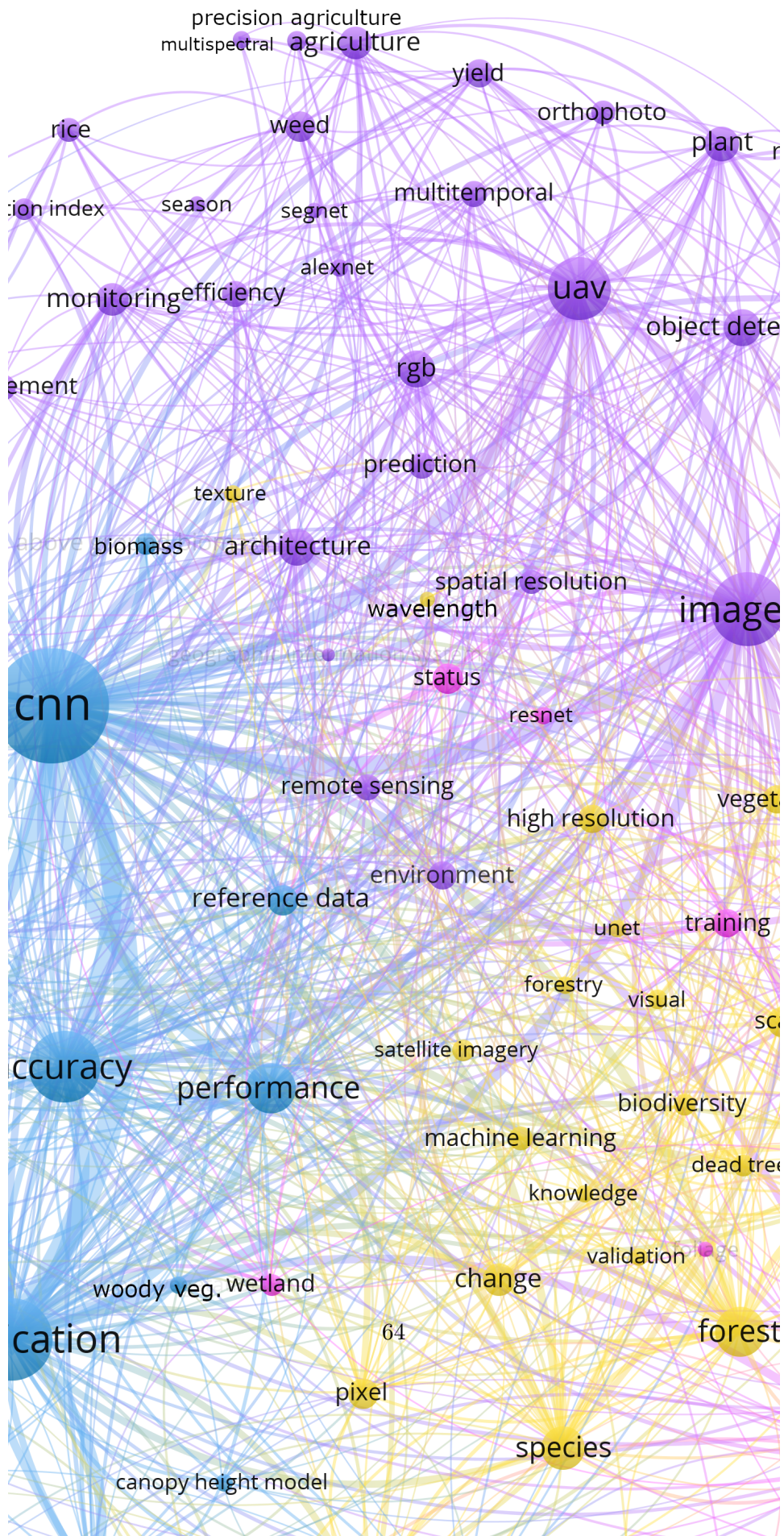
- 2462 Schmitt, M., Hughes, L. H., Qiu, C., & Zhu, X. X. (2019). Sen12ms—a curated  
2463 dataset of georeferenced multi-spectral sentinel-1/2 imagery for deep  
2464 learning and data fusion. *arXiv preprint arXiv:1906.07789* (cit. on  
2465 pp. 16, 20).
- 2466 Schmitt, M., Prexl, J., Ebel, P., Liebel, L., & Zhu, X. X. (2020). Weakly Su-  
2467 pervised Semantic Segmentation of Satellite Images for Land Cover  
2468 Mapping – Challenges and Opportunities. *arXiv preprint*, arXiv 2002.08254.  
2469 <http://arxiv.org/abs/2002.08254> (cit. on pp. 15, 20)
- 2470 Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra,  
2471 D. (2019). Grad-CAM: Visual Explanations from Deep Networks  
2472 via Gradient-Based Localization. *International Journal of Computer*  
2473 *Vision*, 17arXiv 1610.02391, 331–336. <https://doi.org/10.1007/s11263-019-01228-7> (cit. on p. 37)
- 2475 Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on Image Data Aug-  
2476 mentation for Deep Learning. *Journal of Big Data*, 6(1). <https://doi.org/10.1186/s40537-019-0197-0> (cit. on p. 18)
- 2478 Smilkov, D., Thorat, N., Kim, B., Viégas, F., & Wattenberg, M. (2017).  
2479 SmoothGrad: removing noise by adding noise, arXiv 1706.03825.  
2480 <http://arxiv.org/abs/1706.03825> (cit. on p. 37)
- 2481 Sothe, C. [C.], De Almeida, C. M., Schimalski, M. B., Liesenberg, V., La  
2482 Rosa, L. E., Castro, J. D., & Feitosa, R. Q. (2020). A compari-  
2483 son of machine and deep-learning algorithms applied to multisource  
2484 data for a subtropical forest area classification. *International Jour-  
2485 nal of Remote Sensing*, 41(5), 1943–1969. <https://doi.org/10.1080/01431161.2019.1681600> (cit. on pp. 22, 32)
- 2487 Sothe, C., Almeida, C. M. D., Schimalski, M. B., Rosa, L. E. C. L., Castro,  
2488 J. D. B., Feitosa, R. Q., Dalponte, M., Lima, C. L., Liesenberg, V.,  
2489 & Miyoshi, G. T. (2020). Comparative performance of convolutional  
2490 neural network , weighted and conventional support vector machine  
2491 and random forest for classifying tree species using hyperspectral  
2492 and photogrammetric data. *GIScience & Remote Sensing*, 00(00),  
2493 1–26. <https://doi.org/10.1080/15481603.2020.1712102> (cit. on pp. 8,  
2494 32)
- 2495 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov,  
2496 R. (2014). Dropout: A simple way to prevent neural networks from  
2497 overfitting. *Journal of Machine Learning Research* (cit. on p. 17).
- 2498 Su, H., Maji, S., Kalogerakis, E., & Learned-Miller, E. (2015). Multi-view  
2499 convolutional neural networks for 3D shape recognition, In *Proceed-  
2500 ings of the ieee international conference on computer vision*. <https://doi.org/10.1109/ICCV.2015.114>. (Cit. on p. 18)
- 2502 Sumbul, G., Charfuelan, M., Demir, B., & Markl, V. (2019). Bigearthnet: A  
2503 large-scale benchmark archive for remote sensing image understand-  
2504 ing, In *Igarss 2019-2019 ieee international geoscience and remote  
2505 sensing symposium*. IEEE. (Cit. on p. 16).

- 2506 Sun, Y., Huang, J., Ao, Z., Lao, D., & Xin, Q. (2019). Deep Learning Ap-  
2507 proaches for the Mapping of Tree Species Diversity in a Tropical  
2508 Wetland Using Airborne LiDAR and High-Spatial-Resolution Re-  
2509 mote Sensing Images. *Forests*, *10*(11), 1047. [https://doi.org/10.](https://doi.org/10.3390/f10111047)  
2510 [3390/f10111047](https://doi.org/10.3390/f10111047) (cit. on pp. 13, 22)
- 2511 Too, E. C., Yujian, L., Njuki, S., & Yingchun, L. (2019). A comparative  
2512 study of fine-tuning deep learning models for plant disease identifi-  
2513 cation. *Computers and Electronics in Agriculture*, *161*, 272–279 (cit.  
2514 on p. 19).
- 2515 Torres, D. L., Feitosa, R. Q., Happ, P. N., La Rosa, L. E. C., Junior, J. M.,  
2516 Martins, J., Bressan, P. O., Gonçalves, W. N., & Liesenberg, V.  
2517 (2020). Applying fully convolutional architectures for semantic seg-  
2518 mentation of a single tree species in urban environment on high  
2519 resolution UAV optical imagery. *Sensors (Switzerland)*, *20*(2), 1–20.  
2520 <https://doi.org/10.3390/s20020563> (cit. on pp. 24, 26)
- 2521 Toth, C., & Józków, G. (2016). Remote sensing platforms and sensors: A  
2522 survey. <https://doi.org/10.1016/j.isprsjprs.2015.10.004>. (Cit. on  
2523 p. 3)
- 2524 Trier, Ø. D., Salberg, A. B., Kermit, M., Rudjord, Ø., Gobakken, T., Næsset,  
2525 E., & Aarsten, D. (2018). Tree species classification in Norway from  
2526 airborne hyperspectral and airborne laser scanning data. *European*  
2527 *Journal of Remote Sensing*, *51*(1), 336–351. [https://doi.org/10.](https://doi.org/10.1080/22797254.2018.1434424)  
2528 [1080/22797254.2018.1434424](https://doi.org/10.1080/22797254.2018.1434424) (cit. on pp. 26, 32)
- 2529 Tuia, D., Persello, C., & Bruzzone, L. (2016). Domain adaptation for the  
2530 classification of remote sensing data: An overview of recent advances.  
2531 *IEEE geoscience and remote sensing magazine*, *4*(2), 41–57 (cit. on  
2532 p. 19).
- 2533 Turner, W., Spector, S., Gardiner, N., Fladeland, M., Sterling, E., & Steininger,  
2534 M. (2003). Remote sensing for biodiversity science and conservation.  
2535 [https://doi.org/10.1016/S0169-5347\(03\)00070-3](https://doi.org/10.1016/S0169-5347(03)00070-3). (Cit. on p. 3)
- 2536 Valbuena, R., Mauro, F., Rodriguez-Solano, R., & Manzanera, J. A. (2013).  
2537 Accuracy and precision of GPS receivers under forest canopies in a  
2538 mountainous environment. *Spanish Journal of Agricultural Research*.  
2539 <https://doi.org/10.5424/sjar/2010084-1242> (cit. on p. 12)
- 2540 Van Eck, N., & Waltman, L. (2010). Software survey: Vosviewer, a computer  
2541 program for bibliometric mapping. *scientometrics*, *84*(2), 523–538  
2542 (cit. on pp. 42, 64).
- 2543 Wagner, F. H., Sanchez, A., Aidar, M. P. M., Rochelle, A. L. C., Tarabalka,  
2544 Y., Fonseca, M. G., Phillips, O. L., Gloor, E., & Aragão, L. (2020).  
2545 Mapping Atlantic rainforest degradation and regeneration history  
2546 with indicator species using convolutional network. *Plos One*, *15*(2),  
2547 e0229448. <https://doi.org/10.1371/journal.pone.0229448> (cit. on  
2548 p. 21)

- 2549 Wagner, F., Sanchez, A., Tarabalka, Y., Lotte, R. G., Ferreira, M. P., Aidar,  
2550 M. P., Gloor, E., Phillips, O. L., & Aragao, L. (2019). Using the u-  
2551 net convolutional network to map forest types and disturbance in the  
2552 atlantic rainforest with very high resolution images. *Remote Sensing*  
2553 *in Ecology and Conservation*, 5(4), 360–375 (cit. on pp. 17, 24, 26).
- 2554 Wang, Z., & Yang, J. (2017). Diabetic retinopathy detection via deep con-  
2555 volutional networks for discriminative localization and visual expla-  
2556 nation. *arXiv preprint arXiv:1703.10757* (cit. on p. 37).
- 2557 Wang, Chen, Cao, An, Chen, Xue, & Yun. (2019). Individual Rubber Tree  
2558 Segmentation Based on Ground-Based LiDAR Data and Faster R-  
2559 CNN of Deep Learning. *Forests*, 10(9), 793. [https://doi.org/10.](https://doi.org/10.3390/f10090793)  
2560 [3390/f10090793](https://doi.org/10.3390/f10090793) (cit. on pp. 13, 24, 26, 30)
- 2561 Weinmann, M., Jutzi, B., Hinz, S., & Mallet, C. (2015). Semantic point cloud  
2562 interpretation based on optimal neighborhoods, relevant features and  
2563 efficient classifiers. *ISPRS Journal of Photogrammetry and Remote*  
2564 *Sensing*. <https://doi.org/10.1016/j.isprsjprs.2015.01.016> (cit. on  
2565 p. 8)
- 2566 Weinstein, B. G., Marconi, S., Bohlman, S. A., Zare, A., & White, E. P.  
2567 (2020). Cross-site learning in deep learning RGB tree crown de-  
2568 tection. *Ecological Informatics*, 56(December 2019), 101061. [https:](https://doi.org/10.1016/j.ecoinf.2020.101061)  
2569 [//doi.org/10.1016/j.ecoinf.2020.101061](https://doi.org/10.1016/j.ecoinf.2020.101061) (cit. on pp. 15, 16, 23, 26,  
2570 29, 35)
- 2571 Weinstein, B. G., Marconi, S., Bohlman, S., Zare, A., & White, E. (2019).  
2572 Individual tree-crown detection in rgb imagery using semi-supervised  
2573 deep learning neural networks. *Remote Sensing*, 11(11), 1–13. [https:](https://doi.org/10.3390/rs11111309)  
2574 [//doi.org/10.3390/rs11111309](https://doi.org/10.3390/rs11111309) (cit. on pp. 21, 23)
- 2575 White, J. C., Coops, N. C., Wulder, M. A., Vastaranta, M., Hilker, T., &  
2576 Tompalski, P. (2016). Remote Sensing Technologies for Enhancing  
2577 Forest Inventories: A Review. [https://doi.org/10.1080/07038992.](https://doi.org/10.1080/07038992.2016.1207484)  
2578 [2016.1207484](https://doi.org/10.1080/07038992.2016.1207484). (Cit. on p. 3)
- 2579 Windrim, L., & Bryson, M. (2020). Detection, segmentation, and model  
2580 fitting of individual tree stems from airborne laser scanning of forests  
2581 using deep learning. *Remote Sensing*, 12(9). [https://doi.org/10.](https://doi.org/10.3390/RS12091469)  
2582 [3390/RS12091469](https://doi.org/10.3390/RS12091469) (cit. on pp. 30, 31)
- 2583 Xi, Y., Ren, C., Wang, Z., Wei, S., Bai, J., Zhang, B., Xiang, H., & Chen, L.  
2584 (2019). Mapping Tree Species Composition Using OHS-1 Hyperspec-  
2585 tral Data and Deep Learning Algorithms in Changbai Mountains,  
2586 Northeast China. *Forests*, 10(9), 818. [https://doi.org/10.3390/](https://doi.org/10.3390/f10090818)  
2587 [f10090818](https://doi.org/10.3390/f10090818) (cit. on pp. 21, 36)
- 2588 Xi, Z., Hopkinson, C., & Chasmer, L. (2018). Filtering stems and branches  
2589 from terrestrial laser scanning point clouds using deep 3-D fully con-  
2590 volutional networks. *Remote Sensing*, 10(8). [https://doi.org/10.](https://doi.org/10.3390/rs10081215)  
2591 [3390/rs10081215](https://doi.org/10.3390/rs10081215) (cit. on p. 13)

- 2592 Yang, Q., Shi, L., Han, J., Zha, Y., & Zhu, P. (2019). Deep convolutional  
2593 neural networks for rice grain yield estimation at the ripening stage  
2594 using UAV-based remotely sensed images. *Field Crops Research*,  
2595 *235*(February), 142–153. <https://doi.org/10.1016/j.fcr.2019.02.022>  
2596 (cit. on pp. 14, 22, 25, 29, 32, 34)
- 2597 Yuan, Q., Wei, Y., Meng, X., Shen, H., & Zhang, L. (2018). A Multiscale  
2598 and Multidepth Convolutional Neural Network for Remote Sensing  
2599 Imagery, *11*(3), 978–989 (cit. on pp. 29, 33).
- 2600 Zarco-Tejada, P. J., Hornero, A., Beck, P. S., Kattenborn, T., Kempeneers,  
2601 P., & Hernández-Clemente, R. (2019). Chlorophyll content estima-  
2602 tion in an open-canopy conifer forest with Sentinel-2A and hyper-  
2603 spectral imagery in the context of forest decline. *Remote Sensing of*  
2604 *Environment*, *223*, 320–335. [https://doi.org/10.1016/j.rse.2019.01.](https://doi.org/10.1016/j.rse.2019.01.031)  
2605 *031* (cit. on p. 29)
- 2606 Zarco-Tejada, P. J., Camino, C., Beck, P. S., Calderon, R., Hornero, A.,  
2607 Hernández-Clemente, R., Kattenborn, T., Montes-Borrego, M., Susca,  
2608 L., Morelli, M., Gonzalez-Dugo, V., North, P. R., Landa, B. B.,  
2609 Boscia, D., Saponari, M., & Navas-Cortes, J. A. (2018). Previsual  
2610 symptoms of *Xylella fastidiosa* infection revealed in spectral plant-  
2611 trait alterations. *Nature Plants*, *4*(7), 432–439. [https://doi.org/10.](https://doi.org/10.1038/s41477-018-0189-7)  
2612 *1038/s41477-018-0189-7* (cit. on p. 29)
- 2613 Zhang, B., Huang, S., Shen, W., & Wei, Z. (2019). Explaining the point-  
2614 net: What has been learned inside the pointnet?, In *Proceedings of*  
2615 *the ieee conference on computer vision and pattern recognition work-*  
2616 *shops*. (Cit. on p. 37).
- 2617 Zhang, L., Shao, Z., Liu, J., & Cheng, Q. (2019). Deep learning based re-  
2618 trieval of forest aboveground biomass from combined LiDAR and  
2619 landsat 8 data. *Remote Sensing*, *11*(12). [https://doi.org/10.3390/](https://doi.org/10.3390/rs11121459)  
2620 *rs11121459* (cit. on p. 3)
- 2621 Zhang, M., Lin, H., Wang, G., Sun, H., & Fu, J. (2018). Mapping paddy  
2622 rice using a Convolutional Neural Network (CNN) with Landsat 8  
2623 datasets in the Dongting Lake Area, China. *Remote Sensing*, *10*(11).  
2624 <https://doi.org/10.3390/rs10111840> (cit. on pp. 24, 25, 34, 36)
- 2625 Zhao, X., Yuan, Y., Song, M., Ding, Y., Lin, F., Liang, D., & Zhang, D.  
2626 (2019). Use of unmanned aerial vehicle imagery and deep learning  
2627 unet to extract rice lodging. *Sensors (Switzerland)*, *19*(18), 1–13.  
2628 <https://doi.org/10.3390/s19183859> (cit. on p. 29)
- 2629 Zhong, L., Hu, L., & Zhou, H. (2019). Deep learning based multi-temporal  
2630 crop classification. *Remote Sensing of Environment*, *221*, 430–443.  
2631 <https://doi.org/10.1016/j.rse.2018.11.032> (cit. on pp. 21, 25, 34, 36)
- 2632 Zhu, X. X., Montazeri, S., Ali, M., Hua, Y., Wang, Y., Mou, L., Shi, Y., Xu,  
2633 F., & Bamler, R. (2020). Deep learning meets sar. (Cit. on p. 30).
- 2634 Zhu, X. X., Tuia, D., Mou, L., Xia, G. S., Zhang, L., Xu, F., & Fraundorfer,  
2635 F. (2017). Deep Learning in Remote Sensing: A Comprehensive Re-

2636 view and List of Resources. *IEEE Geoscience and Remote Sensing*  
2637 *Magazine*, 5(4), 8–36. <https://doi.org/10.1109/MGRS.2017.2762307>  
2638 (cit. on pp. 3, 4, 7, 15, 34)  
2639 Zou, X., Cheng, M., Wang, C., Xia, Y., & Li, J. (2017). Tree Classification  
2640 in Complex Forest Point Clouds Based on Deep Learning. *IEEE*  
2641 *Geoscience and Remote Sensing Letters*, 14(12), 2360–2364. <https://doi.org/10.1109/LGRS.2017.2764938> (cit. on pp. 18, 26, 30)  
2642





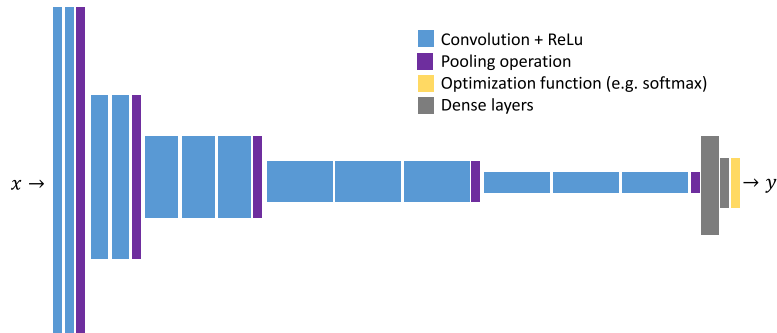


Figure 4: Schematic diagram of the *VGG-16* architecture. The 16 stands for the number of convolutional and dense layers. Frequently used alternatives are *VGG-8* and *VGG-19*.

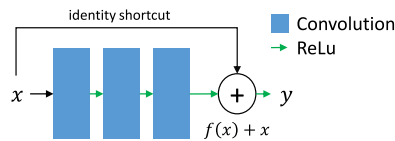


Figure 5: Schematic diagram of a residual building block used in repeated sequence in common *ResNet* architectures.

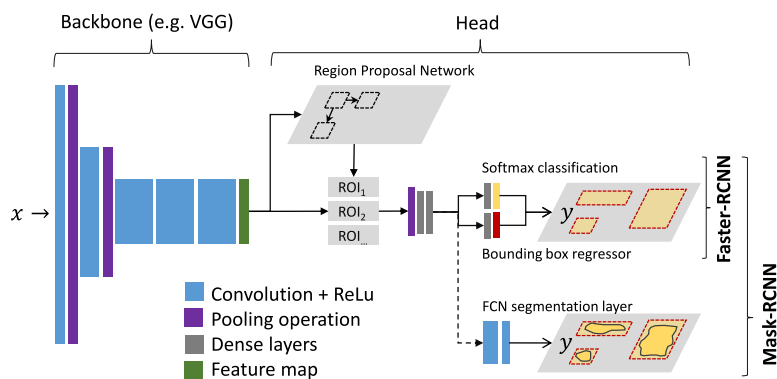


Figure 6: Faster-R-CNN and Mask-RCNN, respectively.

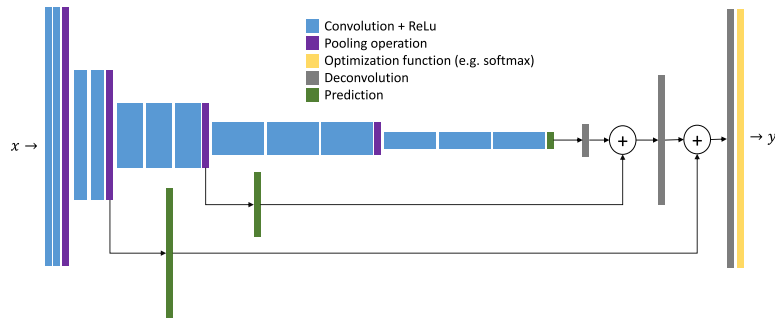


Figure 7: Schematic diagram of the FCN architecture as proposed by Long et al. (2015). Predictions (also referred to as 'scores') within the network are forwarded to deeper layers to relate respective activations to the original spatial resolution.

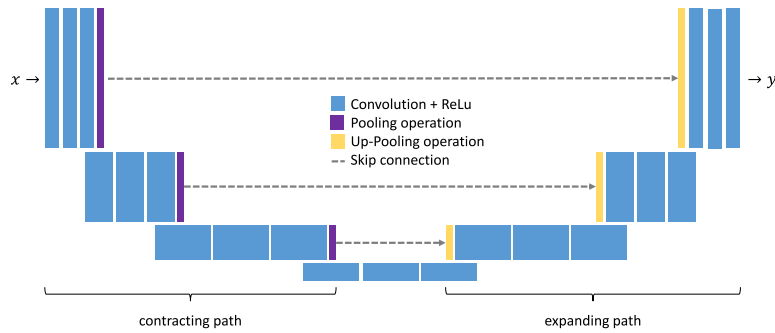


Figure 8: Schematic diagram of the *U-Net* architecture depicting its encoder-decoder structure using an contracting and expanding path.

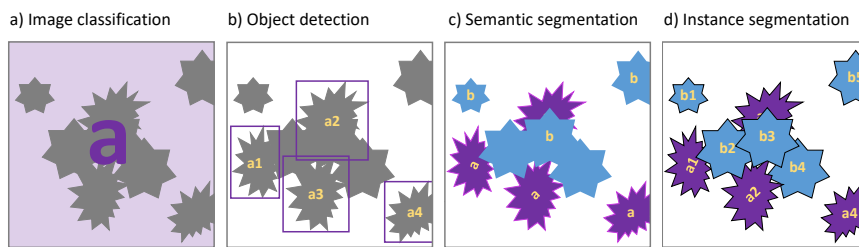


Figure 9: Schemes illustrating the conceptual differences between different CNN approaches, including a) image classification, where the entire image is assigned to a class; b) object detection, where individual occurrences are localized and their extent estimated with bounding boxes; c) semantic segmentation, which assigns each pixel of the input image to the target classes; and d) instance segmentation, where individuals belonging to a class are mapped.

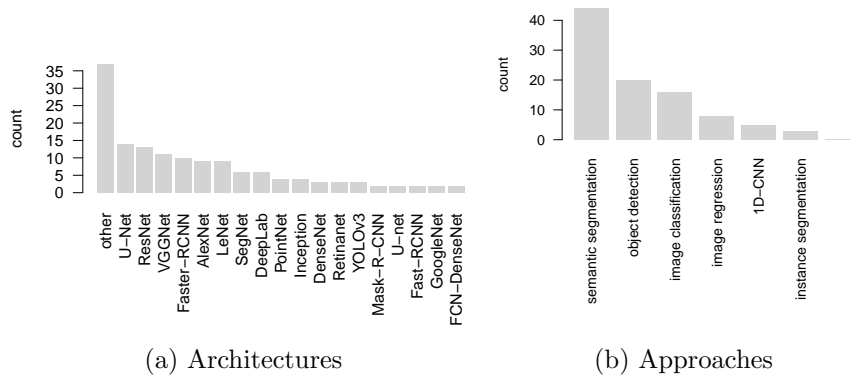


Figure 10: Barplots characterizing the reviewed literature in terms of frequency of a) different architectures, including direct implementations as well as modifications of the original architecture and b) different approaches

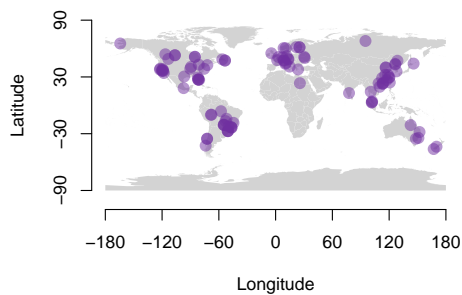


Figure 11: Study areas of the reviewed studies

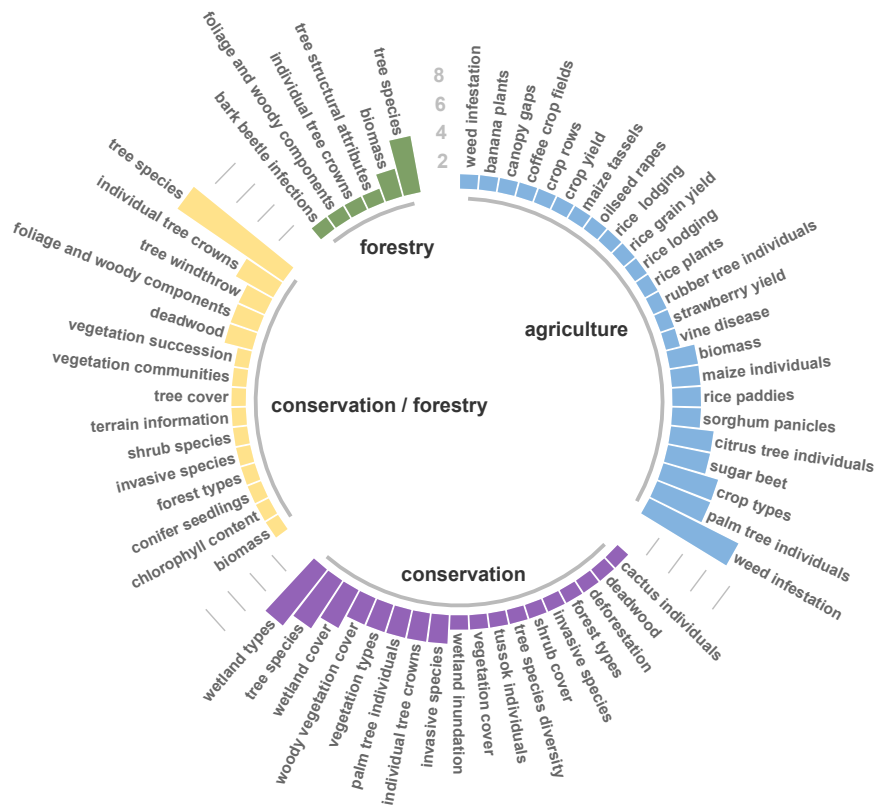


Figure 12: Frequency of studies in the context of agriculture, forestry, and conservation. The class *forestry/conservation* includes studies that are relevant for both fields.

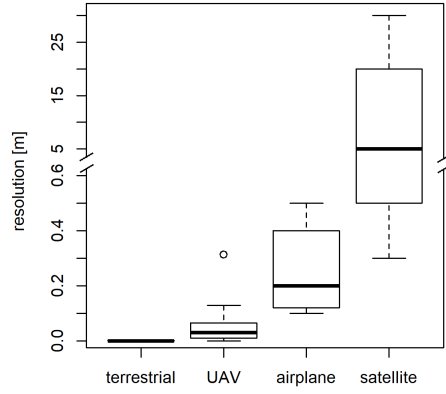


Figure 13: Frequency distribution of spatial resolutions by different remote sensing platforms among the reviewed studies (only raster products considered).

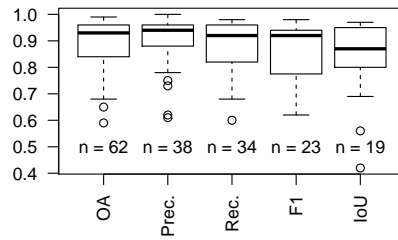


Figure 14: Validation results of the CNN-based predictions derived from the reviewed studies. The studies used different metrics (frequency = n), including Overall Accuracy (OA), Precision (Prec.), Recall (Rec.), F-score (F) and IoU (Intersect over Union).

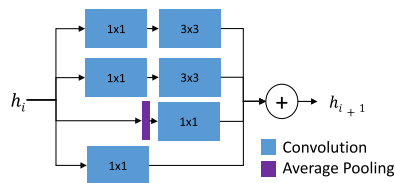


Figure 15: A schematic representation of an Inception-module