

1 Semantic segmentation of water bodies in very high-resolution satellite 2 and aerial images

3

4 **Marc Wieland** ^{1*}, **Sandro Martinis** ¹, **Ralph Kiefl** ¹, **Veronika Gstaiger** ²

5 ¹ German Remote Sensing Data Center (DFD), German Aerospace Center (DLR), Oberpfaffenhofen, D-82234
6 Wessling, Germany

7 ² Remote Sensing Technology Institute (IMF), German Aerospace Center (DLR), Oberpfaffenhofen, D-82234
8 Wessling, Germany

9

10 * Correspondence: marc.wieland@dlr.de

11

12 **Abstract**

13 This study evaluates the performance of convolutional neural networks for semantic segmentation of water bodies
14 in very high-resolution satellite and aerial images from multiple sensors with particular focus on flood emergency
15 response applications. Different model architectures (U-Net and DeepLab-V3+) are combined with encoder
16 backbones (MobileNet-V3, ResNet-50 and EfficientNet-B4) and tested for their ability to delineate inundated
17 areas under varying environmental conditions and data availability scenarios. An unprecedented reference dataset
18 of 1,120 globally sampled images with quality checked binary water masks is introduced and used to train, validate
19 and test the models for water body segmentation. Furthermore, independent test datasets are developed to test the
20 generalization ability of the trained models across regions, sensors (IKONOS, GeoEye-1, WorldView-2,
21 WorldView-3 and four different airborne camera systems) and tasks (normal water and flood water segmentation).
22 Results indicate that across all tested scenarios a U-Net model with Mobilenet-V3 backbone pre-trained on
23 ImageNet performs best. While using R-G-B image bands performs well, adding the near infrared band (if
24 available) slightly improves prediction results. Similarly, adding slope information from an independent digital
25 elevation model increases accuracies. Train-time augmentation and contrast enhancement could improve
26 transferability across sensors and in particular between satellite and aerial images. Moreover, adding noisy training
27 data from freely available online resources could further improve performance with minimal annotation effort.

28

29 **Keywords**

30 Convolutional Neural Networks; Semantic segmentation; Water; Emergency response; Rapid mapping

31 **1 Introduction**

32 Very high-resolution optical satellite and aerial images (in the following defined as having a ground
33 sampling distance < 4 m and spectral bands in the visible and near infrared part of the electromagnetic
34 spectrum) are frequently requested during rapid mapping activations to gain situational awareness about
35 large-scale flood disasters (Martinis et al., 2017; Voigt et al., 2016). Near-real time information about
36 (flood) water extent and affected infrastructure are amongst the most critical components that can be
37 delivered by this kind of images to support prioritization of response actions. The tasking capability of
38 some of the satellite imaging platforms and the fact that airborne sensors can acquire images at any (day)
39 time and below cloud-coverage makes them particularly favourable for time-critical emergency response
40 applications. These can provide a valuable supplement to systematically acquiring satellites with
41 Synthetic Aperture Radar (SAR) or multi-spectral sensors, like Sentinel-1, Sentinel-2 or Landsat. These
42 satellites acquire images at lower spatial resolution (~ 10 -30 m) and on a regular schedule usually every
43 couple of days depending on the area of interest. Hence, very high-resolution optical satellite and aerial
44 images help to fill temporal acquisition gaps and provide more spatial details of an ongoing disaster
45 situation.

46 To assure that geo-information products have the highest possible spatial and temporal resolutions and
47 information content, it is crucial that algorithms used for remote sensing image-based emergency
48 mapping are able to simultaneously use data from a variety of acquisition platforms and sensors.
49 Automated water segmentation in this kind of images, however, is a challenging task because of
50 significant variations in spectral reflectance characteristics, size and shape of water bodies. Compared to
51 normal water, flood water for example may be characterized by higher contents of sediment and debris.
52 Vegetation, infrastructure, boats or other vehicles may further interact with the water surface or sub-
53 surface, and also lighting conditions (e.g., sun reflectance on water surface or shadows by buildings or
54 vegetation) may cause high intra and inter water body variability. The fractal geometry of the land-water
55 border further increases the complexity of the segmentation task, since the definition of a hard class
56 boundary may be subjective depending on the image resolution. Moreover, bottom features may be
57 visible through the water column, especially in clear shallow water bodies, which further increases the
58 complexity of automated water segmentation. As a result, common operational procedures for remote
59 sensing image-based emergency flood mapping rely largely on manual image interpretation.

60 Traditional methods for the segmentation of water bodies mainly focus on thresholding water indices,
61 such as the Normalized Difference Water Index (NDWI) (McFeeters, 1996). These methods exclusively

62 rely on spectral information and require at least one or more image bands in the near- or short-wave
63 infrared, which negatively impacts their transferability to other sensors. When applied to very high-
64 resolution images such methods, moreover, show a lack of generalization ability, because the more
65 spatial details are visible, the higher the spectral variability within the water class becomes. To overcome
66 some of the limitations of simple thresholding methods, supervised machine learning methods such as
67 Support Vector Machine (SVM) that learn a decision boundary between classes from labelled training
68 samples and based on hand-crafted features have been applied for water mapping (Ireland et al., 2015).
69 In recent years, deep learning approaches have shown superior performance on various image analysis
70 tasks across disciplines. Convolutional Neural Networks (CNNs) are widely used for semantic
71 segmentation because of their ability to implement nonlinear decision functions and to learn features
72 directly from raw images by combining convolutional and pooling layers (Ball et al., 2017).

73 Current semantic segmentation algorithms, such as U-Net (Ronneberger et al., 2015), PSP-Net (Zhao et
74 al., 2017) or DeepLab-V3+ (Chen et al., 2017) follow an encoder-decoder architecture, which consists
75 of a down-sampling path to learn dense image features and capture context (encoder) as well as a
76 symmetric up-sampling path that learns an optimal interpolation of the features to recover the original
77 image resolution (decoder). U-Net uses skip connections for feature propagation, whereas PSP-Net and
78 DeepLab-V3+ use spatial pyramid pooling to get contextual features at multiple scales. Dilated
79 convolutions are utilized by the latter to capture more contextual information without increasing the
80 number of trainable parameters. Different encoder networks have been proposed in literature, amongst
81 which the most prominent are VGG-16 (Simonyan and Zisserman, 2015) and ResNet (He et al., 2015a).
82 These networks are originally intended for image classification, but can be used as feature extractors for
83 semantic segmentation. Accuracy improvements between different network versions and architectures
84 have largely been achieved by scaling in depth (adding more layers), width (increasing spatial size of
85 layers) or resolution (adding more channels per layer). This has mostly been done by manual tuning
86 which, given the large number of parameters involved, makes finding an optimal performance point
87 difficult and resource consuming. Tan and Le (2019) showed that scaling all three dimensions uniformly
88 with a fixed set of scaling factors can improve model performance while keeping resource consumption
89 minimal (Tan and Le, 2019). They used neural architecture search to develop a novel baseline model and
90 apply compound scaling to optimize network depth, width, and resolution simultaneously. The main
91 building blocks of their EfficientNet models are mobile inverted bottleneck blocks (MBConv) with
92 squeeze-and-excitation optimization (Hu et al., 2020). EfficientNet models outperform other CNNs on
93 ImageNet (Deng et al., 2009) and several transfer learning benchmark datasets while consistently being

94 smaller and thus more efficient than previous models. Howard et al. (2019) also use hardware-aware
95 neural architecture search to develop Mobilenet-V3, which is particularly light-weight and optimized to
96 run on limited hardware resources like mobile phone CPUs. Despite their promising characteristics,
97 Mobilenet-V3 and EfficientNet encoders are not (yet) widely used in remote sensing studies. Other
98 encoders that have been tested in remote sensing applications include but are not limited to ResNext (Xie
99 et al., 2017), PSPNet (Yuan et al., 2022) and ShuffleNet (Gomes et al., 2021). A new branch of networks
100 that gain popularity recently are transformer networks. Cao et al. (2021) for example report superior
101 performance of their U-Net-like transformer for medical image segmentation compared to fully
102 convolutional networks. First studies that apply transformer networks on remote sensing images also
103 show promising results (Ding et al., 2022; Gu et al., 2022; Xu et al., 2021).

104 Since reference data in remote sensing is expensive to generate and satellite images are often restricted
105 by conservative licenses that prevent re-distribution of purchased imagery, only few publicly available
106 benchmark datasets exist for semantic segmentation of very high-resolution satellite and aerial images,
107 namely ISPRS (Rottensteiner et al., 2013), Dstl (Iglovikov et al., 2017) and Deepglobe (Demir et al.,
108 2018). Of these only the Dstl Kaggle Challenge and Deepglobe Land Cover Classification Challenge
109 datasets contain water classes, whereas only the latter covers samples outside of urban areas. Castillo-
110 Navarro et al. (2021) created the MiniFrance dataset that is specifically designed for semi-supervised
111 segmentation of land-use. It contains over 2,000 labelled and unlabelled aerial images across France,
112 with the water class covering approximately 1 % of the dataset. Due to their general focus on land-cover
113 / land-use classification, the available datasets can cover only limited variations in spectral reflectance
114 characteristics, size and shape of water bodies and none of them considers a flood water class. Recently,
115 several datasets for natural disaster applications have been released (Bonafilia et al., 2020; Gupta et al.,
116 2019). Hänsch et al. (2022) present the SpaceNet 8 dataset for the detection of flooded roads and
117 buildings. Rahnemoonfar et al. (2021) introduce the FloodNet dataset, which is composed of partially
118 annotated UAV images after Hurricane Harvey and aims at post flood scene understanding. The semantic
119 segmentation labels include amongst others a water class and several classes of flooded infrastructure.
120 The dataset is, however, from a single location and covers on a relatively small amount of labelled data.
121 The available disaster related datasets are often limited in geographical coverage and / or have a strong
122 focus on infrastructure.

123 To artificially increase the number of training samples in limited data scenarios, data augmentation is
124 increasingly being used in remote sensing. Also transfer learning (Pan and Yang, 2010) is an approach
125 to ease the problem of scarce training data. CNNs trained on large-scale natural image datasets (e.g.,

126 ImageNet) have successfully been transferred to segment remote sensing images, either by means of fine-
127 tuning the network with data from the target domain or by using the pre-trained network directly as
128 feature extractor. However, in particular deep models with a large number of trainable parameters may
129 require large amounts of labelled samples of the target domain to avoid over-fitting. Semi-supervised
130 learning, which jointly uses labelled and unlabelled data, provides another promising approach to
131 improve segmentation when training data is limited (Li et al., 2019). Few studies also experiment with
132 training on noisy labels. Mnih and Hinton (2012) and Kaiser et al. (2017) report that satisfying
133 performance can be obtained with significantly less manual annotation effort, by exploiting noisy large-
134 scale training data from independent maps or data from crowd-sourcing platforms like OpenStreetMap
135 (Haklay, 2010).

136 While most deep learning studies that use very high-resolution optical satellite and aerial images focus
137 on building segmentation (Neupane et al., 2021), road delineation (Mattyus et al., 2015) or object
138 detection (Azimi et al., 2021), relatively few studies exist for water body segmentation (Huang et al.,
139 2018). Yang et al. (2020) propose a model based on mask R-CNN to detect and segment water bodies in
140 GaoFen-2 images. They train their models on Google Earth image tiles and GaoFen-2 images and test
141 on a hold-out dataset of GaoFen-2 images. Their results indicate superior performance of a ResNet-50-
142 based model compared to a more complex ResNet-100-based model. Guo et al. (2020) introduce a multi-
143 scale water extraction CNN for GaoFen-1 images and compare their results to U-Net and DeepLab-V3+
144 models. Chen et al. (2018) apply a self-adaptive pooling CNN on clusters of image super-pixels to extract
145 urban water bodies in GaoFen-2 and Ziyuan-3 images. They train and test their method on images of
146 Beijing, Chengdu and Tianjin and report superior accuracy compared to SVM and water index
147 thresholding. Gebrehiwot et al. (2019) investigate the potential of a VGG-based fully convolutional
148 network to extract flooded areas from Unmanned Aerial Vehicle (UAV) images acquired over three study
149 areas in North Carolina, United States after Hurricane Matthew. Their results support conclusions of
150 studies from other disciplines that adapting pre-trained models and fine-tuning them can lead to superior
151 performance, especially when domain-specific training data are scarce. Duan and Hu (2020) compare a
152 multi-scale refinement network with U-Net, SegNet and DeepLab-V3+ networks for water segmentation
153 in WorldView-3 and GaoFen-2 satellite images. They report superior performance of the refinement
154 network and DeepLab-V3+ compared to U-Net and SegNet. Feng et al. (2019) propose a modified U-
155 Net architecture and combine it with a super-pixel segmentation method to refine details of the land-
156 water borders. They test their method on four images acquired by WorldView-2 and GaoFen-2 satellites.

157 Existing studies may achieve encouraging results for experimental setups, but largely focus on limited

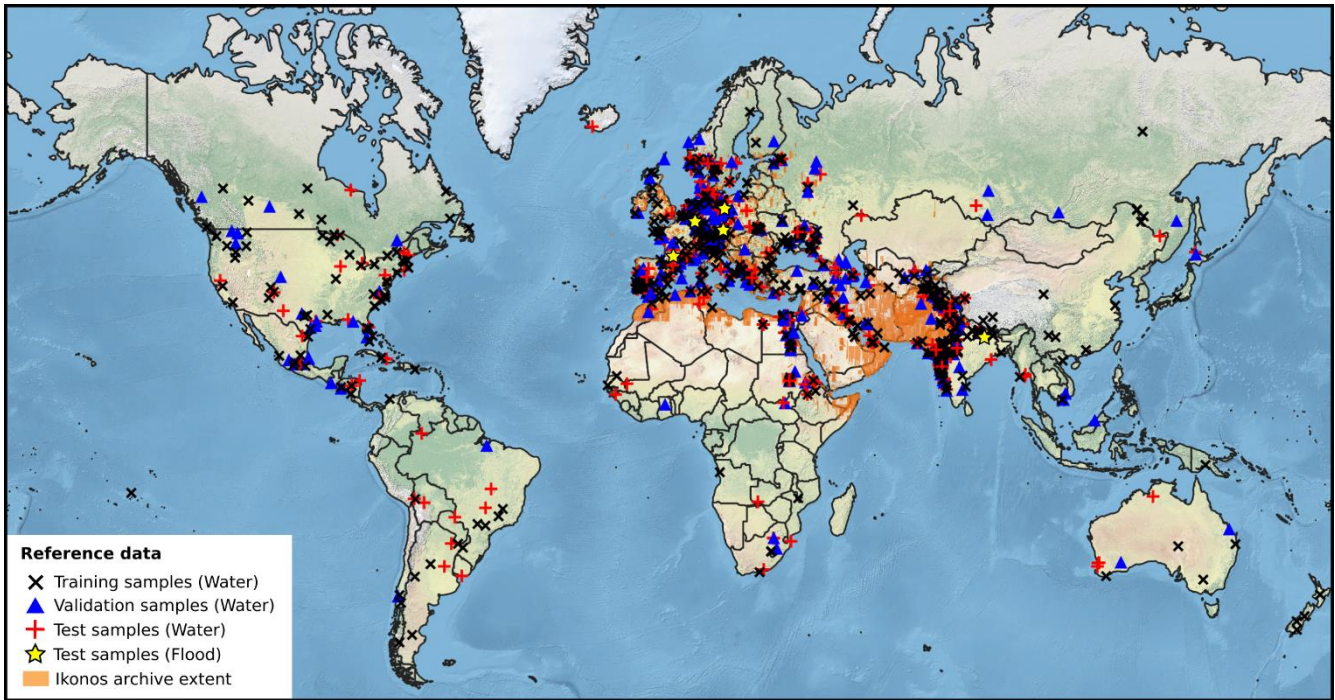
158 geographical coverage, single sensors and propose complex solutions that may not scale well to time-
159 critical emergency response applications. Further research is required to train water segmentation models
160 that are applicable to different satellite sensors and images across varying atmospheric conditions and
161 scene properties. The objective of this study is to evaluate the performance of CNNs for semantic
162 segmentation of surface water bodies in very high-resolution satellite and aerial images. We render this
163 as a binary segmentation task with classes “land” and “water”. Depending on whether a satellite image
164 has been acquired during normal hydrological conditions or during a flood event, we distinguish between
165 scenes that primarily depict “normal water” or “flood water”.

166 Most promising architectures (U-Net and DeepLab-V3+) are combined with different encoders
167 (MobileNet-V3, ResNet-50 and EfficientNet-B4) and tested for their ability to delineate inundated areas
168 under varying environmental conditions and data availability scenarios. For each combination of
169 architecture and encoder we test different weight initializations and evaluate the performance of the
170 trained models to (I) segment normal water in images of the same sensors, (II) segment normal water in
171 images of other sensors, and (III) segment flood water in images of other sensors. We introduce an
172 unprecedented reference dataset of 1,120 images with quality checked binary water masks. We also
173 develop an independent test dataset for four flood events to evaluate the generalization ability of the
174 trained models across regions, sensors and tasks. Compared to other studies in this direction, we consider
175 sensors from multiple sensors and acquisition platforms (IKONOS, GeoEye-1, WorldView-2,
176 WorldView-3, four different airborne camera systems and images from a Web Map Service of Mapbox
177 (“Mapbox Satellite,” 2021)) and train, validate and test considered networks on a globally sampled
178 dataset that covers variations in environmental and atmospheric conditions as well as spectral reflectance
179 characteristics, size and shape of water bodies.

180 **2 Data**

181 We compile a reference dataset based on 1,120 globally distributed very high-resolution satellite and
182 aerial images to train, test and validate the water segmentation models. For the dataset to be representative
183 for a large variety of climatic, atmospheric, and land-cover conditions, we combine several global
184 environmental layers into a unified stratification layer. Each stratum is aligned to a global grid with
185 approximately 1.5 x 1.5 km spacing, contains water bodies (Pekel et al., 2016) and represents a unique
186 combination of biome (Olson et al., 2001) and pre-dominant land-cover (Buchhorn et al., 2020). Samples
187 are selected by means of a stratified random sampling procedure (Figure 1). Data availability constraints
188 are considered, acquisition times cover different seasons and the maximum cloud-cover percentage for

189 image acquisitions is set to 5 %. For the purpose of this study, we acquire Internal Data Access (IDA) to
190 archived IKONOS images received by the satellite down-link for Europe and Middle-East between 2000
191 and 2009 (Schreier et al., 2008). For each sample we acquire the respective images, orthorectify them,
192 convert digital numbers to top of canopy reflectance, pan-sharpen them to 0.8 m spatial resolution, stack
193 the Red (R), Green (G), Blue (B) and Near-Infrared (NIR) image bands together and create a 2,048 x
194 2,048 pixels subset. Thematic masks are delineated into classes “water” and “land” based on semi-
195 automated image analysis. We first apply a simple NDWI thresholding to extract an initial water mask,
196 which we iteratively improve by manual adjustments and quality controls by several experienced
197 operators (Figure 2). Due to the limited geographical coverage of the IDA IKONOS archive, we
198 complement the reference dataset with freely available optical images from a Web Map Service of
199 Mapbox (“Mapbox Satellite,” 2021). Images are acquired as geo-referenced R-G-B composites at zoom
200 levels 17 and 18, which approximately equals a ground sampling distance of 0.6 to 1.2 m depending on
201 the latitude of the sample location. Data at these zoom levels are a combination of satellite and aerial
202 images from Maxar and various open national aerial image archives with JPG compression, reduced
203 radiometric resolution and enhanced contrast compared to the original images. We download thematic
204 water masks for each of these samples from OpenStreetMap (“OpenStreetMap,” 2021) and iteratively
205 improve them by manual quality control and adjustments by experienced operators. Images and masks
206 are shuffled and divided into training (50 %), validation (25 %) and testing (25 %) datasets, which are
207 then split into non-overlapping tiles with 256 x 256 pixels size. The final dataset covers approximately
208 90,000 tiles and 2,800 km².



209

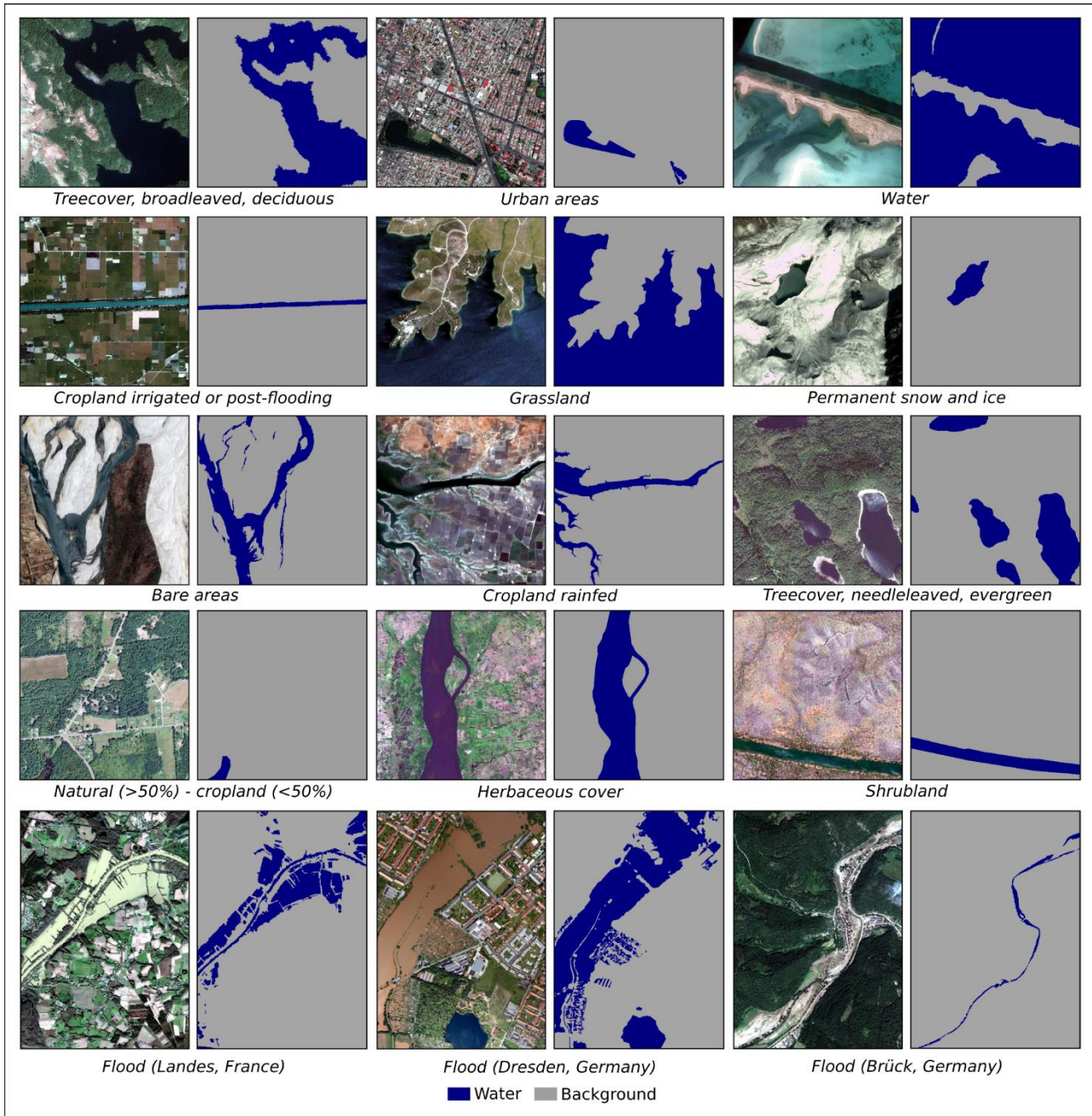
210 *Figure 1: Spatial distribution of the reference data samples used for training, validation and testing of the water*
 211 *segmentation method; Independent test samples for flood water segmentation from multiple sensors (GeoEye-1,*
 212 *WorldView-2, WorldView-3 and four different airborne camera systems).*

213 To independently test the performance of trained models on flood water scenes and to consider additional
 214 platforms and sensors, we select six locations of four flood events (Dresden Germany 2013, Landes
 215 France 2019, Bihar India 2020, Minden, Brück and Altenahr Germany 2021), for which normal water
 216 (pre- or post-event) and flood water (co-event) satellite or aerial images are available. We use normal
 217 water images to test transfer of models to other sensors for water segmentation. To test model transfer to
 218 other sensors and tasks (flood water segmentation), we use flood images. For each test case we manually
 219 delineate water extents by means of a standard rapid mapping workflow and make sure to cross-check
 220 results between different experienced operators (Long et al., 2021). Cloud and cloud shadow pixels are
 221 manually masked in both reference and predicted masks and excluded from accuracy assessments. Table
 222 1 shows an overview of the available images for each case study. The dataset covers approximately
 223 15,000 tiles (256 x 256 pixels) and 204 km².

224 Each sample that we consider in this study is further complemented with slope information derived from
 225 a Digital Elevation Model (DEM) raster of the freely and globally available Copernicus DEM (Fahrland
 226 et al., 2020). The slope raster is up-sampled from its original 30 m resolution to align with the respective
 227 satellite and aerial images.

228 *Table 1: Overview of normal water and flood images for case studies with information about acquisition dates,*

Case study		Normal water images			Flood water images		
		Acquisition	Sensor Resolution Spectral bands	Coverage (km ²)	Acquisition	Sensor Resolution Spectral bands	Coverage (km ²)
Satellite	Landes, France <i>Riverine floods</i>	2019-09-11	GeoEye-1 0.50 m R-G-B-NIR	35	2019-12-19	GeoEye-1 0.50 m R-G-B-NIR	35
	Bihar, India <i>Monsoon rains</i>	2020-02-11	WorldView-2 0.50 m R-G-B-NIR	60	2020-08-24	GeoEye-1 0.50 m R-G-B-NIR	60
	Brück, Germany <i>Riverine floods</i>	2021-09-07	GeoEye-1 0.50 m R-G-B-NIR	4	2021-07-21	Worldview-3 0.50 m R-G-B-NIR	4
Aerial	Altenahr, Germany <i>Riverine floods</i>	2021-10-23	Zenmuse H20T 0.07 m R-G-B	1	2021-07-20	DLR MACS 0.20 m R-G-B	1
	Minden, Germany <i>Riverine floods</i>	2019-06-28	DOP20 0.20 m R-G-B	1	2021-07-16	DLR 4K 0.07 m R-G-B	1
	Dresden, Germany <i>Riverine floods</i>	2011-05-31	DOP20 0.20 m R-G-B	1	2013-06-08	DLR 3K 0.25 m R-G-B	1



231

232 *Figure 2: Examples of reference data samples with pairs of image and corresponding water mask for different*
 233 *land-cover settings. Each sample covers an area of approximately 1.3 km². The last row shows flood water images*
 234 *of the independent test dataset. The area covered by the flood water samples varies and is depicted in Table 1.*

235 **3 Method**

236 In a series of experiments, we aim to answer the following specific research questions.

- 237 1. Baseline model (BM): Which combination of architecture and encoder performs best across
 238 varying data availability scenarios?

- 239 2. Input bands (IB): Can additional spectral band and independent slope information improve the
240 performance?
- 241 3. Data augmentation (DA): Can different data augmentation techniques improve the generalization
242 ability?
- 243 4. Pre-training (PT): What influence do different pre-trained weights have on the performance?
- 244 5. Transfer model – satellite to aerial (TM): Can a model trained on satellite images be transferred
245 to contrast enhanced aerial images with reduced radiometric resolution? Can noisy training data
246 improve the performance?

247 In the following we describe the setup for each experiment and general considerations that apply for all
248 experiments.

249 *Baseline model (BM)*

250 Several studies have shown that the U-Net architecture (Ronneberger et al., 2015) is able to deliver state-
251 of-the-art results in water segmentation tasks in high-resolution satellite images while keeping
252 computational complexity low (Wieland and Martinis, 2019). For comparison we choose the
253 DeepLabV3+ architecture (Chen et al., 2017), which extracts information on multiple scales
254 simultaneously with atrous spatial pyramid pooling using diluted convolution kernels with different rates.
255 Depth-wise separable convolution is applied to decrease computation complexity. The hypothesis is, that
256 these architectural differences may allow the model to have a higher level of context-awareness, which
257 may positively impact on the segmentation performance. In combination with these architectures we
258 compare promising encoders, namely MobileNet-V3, ResNet-50 and EfficientNet-B4. We selected these
259 for our experiments since they show a good trade-off between number of model parameters and ImageNet
260 Top-1 accuracy (Tan and Le, 2019).

261 *Input bands (IB)*

262 Water shows low reflectance in the NIR wavelengths as it absorbs more energy, while non-water
263 generally has a higher reflectance. Thus, a high contrast in reflectance values between water and non-
264 water landcover classes is particularly dominant in the NIR spectral band compared to the visible R, G
265 and B bands. To this regard, we test the influence of the NIR spectral band on the water segmentation
266 performance. Additionally, we also consider slope information derived from the freely and globally
267 available Copernicus DEM. For this purpose, we compute slope in percent at the 30 m native spatial
268 resolution of the DEM and resample it to match the respective sample image resolution.

269 *Data augmentation (DA)*

270 Data augmentation provides a way to learn invariance to changes in the augmented domains beyond what
271 is present in the raw training images. Remote sensing image properties are affected by changes in
272 atmospheric conditions, land-use / land-cover, seasonality and other scene and image properties such as
273 sun elevation or radiometric resolution. Therefore, even large remote sensing training datasets may not
274 cover all eventualities that may occur in real-world applications. To this regard, we test the influence of
275 different data augmentation techniques on the segmentation performance. In particular, we apply
276 training- and test-time augmentations. For this experiment the training dataset is augmented with random
277 contrast, brightness, scale and rotation. Factors are randomly applied within predefined ranges to contrast
278 $[-0.1, 0.1]$, brightness $[-0.1, 0.1]$ and scale $[0.9, 1.1]$. Rotation is performed in steps of 90 degrees. All
279 augmentations are applied with equal probability. During test-time we predict on the original input image
280 and five randomly augmented versions of it. The final prediction is based on the averaged class
281 probabilities.

282 *Pre-training (PT)*

283 Starting from pre-trained weights can improve performance, even if these have been trained on different
284 image types, because low-level features that are being learned in early network layers are similar across
285 image domains (Kaiser et al., 2017). In this experiment we compare the influence of different pre-trained
286 weights on the segmentation performance. We initialize weights either randomly as described in He et
287 al. (2015b) or we use pre-trained weights that were previously learned on ImageNet (Deng et al., 2009)
288 or on a global reference dataset for water segmentation in Sentinel-2 satellite images (Wieland et al., in
289 preparation). Since ImageNet weights only cover R-G-B bands, we initialize the weights of additional
290 image bands (N and slope) randomly. In case of Sentinel-2 weights, we pre-train the model specifically
291 on the same spectral bands and slope information. The assumption here is that while in general pre-
292 training has been reported to boost performance, application and domain specific pre-training on
293 Sentinel-2 satellite images may be superior compared to pre-training on ImageNet (despite the lower 10
294 m spatial resolution of Sentinel-2).

295 *Transfer model – satellite to aerial (TM)*

296 In this experiment we investigate if a model trained on satellite images can be transferred to aerial images
297 with reduced radiometric resolution and enhanced contrast. We also aim at answering the question,
298 whether noisy training data can improve the generalization ability across images acquired by different
299 sensors and acquisition platforms. Specifically, we test the influence of different training approaches and

300 compare the performance of models trained solely on the high-quality IKONOS dataset or the noisy
301 Mapbox dataset. We further evaluate the influence of contrast enhancement on the performance. In this
302 case, we apply a percentage linear contrast stretch between the 2nd and 98th percentiles separately for
303 each spectral band. Slope information is not stretched.

304 *General considerations*

305 For all experiments we evaluate the performance of the trained models under consideration of varying
306 data availability across three test scenarios.

- 307 • **Test scenario I:** Segment normal water in images of the same sensors. This scenario applies the
308 trained models on images of the same sensors and for the same task that they have been trained for.
309 Test data are the test split of our global reference dataset.
- 310 • **Test scenario II:** Segment normal water in images of other sensors. This scenario applies the trained
311 models on images of other sensors but for the same task that they have been trained for. Test data are
312 independent images taken during normal water conditions over locations and by sensors that were
313 not used for training.
- 314 • **Test scenario III:** Segment flood water in images of other sensors. This scenario applies the trained
315 models on images of other sensors and for another task that they have been trained for. Test data are
316 independent co-flood images. The target class is flood water, which due to its largely differing
317 spectral characteristics and appearance in the images can be considered as a different task than normal
318 water segmentation. Images are acquired over the same locations as in test scenario II to minimize
319 the effect of the background class on the prediction results and hence isolate the influence of different
320 sensors and water conditions.

321 We use reference data acquired by satellites with known preprocessing, radiometry and additional NIR
322 spectral band as basis for all experiments. Only the transfer experiment (TM) makes use of the Mapbox
323 and aerial image reference dataset, with reduced radiometric resolution, enhanced contrast and for which
324 a NIR spectral band is not available. The input image feature space in all experiments is standardized to
325 zero mean and unit variance with mean and standard deviation being computed on the training dataset
326 and applied to the validation and testing datasets. The training set is shuffled once between every epoch.

327 Initial learning rate and weight decay hyperparameters are optimized for each network setup specifically
328 using a tree-structured Parzen estimator with pruning (Akiba et al., 2019) over 15 trials (for five epochs
329 each) while maximizing the *IoU* score on the validation dataset. As search space we define log uniform

330 distributions (from $1e^{-5}$ to $1e^{-1}$) for the initial learning rate and weight decay. A weighted combination of
331 cross-entropy (L_{CE}) and Lovász loss ($L_{Lovász}$) is used for a direct optimization of the *IoU* score (Rakhlin
332 et al., 2018) during backpropagation (Equation 1).

$$L = (1 - \alpha_{Loss}) * L_{CE} + \alpha_{Loss} * L_{Lovász} \quad \text{Equation 1}$$

333 With α_{Loss} being the weighting coefficient to weight the two loss functions. Following the results of
334 (Helleis et al., 2022) we set $\alpha_{Loss} = 0.5$. To further account for class imbalance, we apply weights on
335 the positive class which are computed on the class distribution of the training dataset. We step-wise
336 reduce the learning rate by a factor of 0.1 if no improvement is seen for three epochs and early stopping
337 is applied in case of no improvement for nine consecutive epochs. For the optimization algorithm we
338 select Adam (Kingma and Lei, 2015). For model evaluation we track cross-entropy Lovász loss, *IoU*
339 score, Precision (*Prec*) and Recall (*Rec*). In order to account for the focus on applications in emergency
340 response, where computation time is a critical performance criterion, we also report model *throughput*
341 measured in megapixel per second (mp/s). Measurements per experiment are averaged across five
342 prediction runs on 5,000 tiles with shape (256, 256, 3). We use Pytorch as deep learning framework and
343 train the networks in batches of 24 until convergence on two NVIDIA Tesla K80 GPUs.

344 Since all CNNs make predictions on local windows, higher prediction errors towards the image borders
345 may be observed. Therefore, during inference we expand the input image with mirror-padding, split it
346 into overlapping tiles, run the predictions over batches of tiles, blend the prediction tiles to reconstruct
347 the expanded input image’s x-y-shape and un-pad the resulting prediction image. We use a tapered cosine
348 window function to weight pixels when blending overlapping tiles together (Wieland and Martinis,
349 2019).

350 To assess the segmentation performance, standard accuracy metrics Intersection over Union score (*IoU*),
351 Precision (*Prec*) and Recall (*Rec*) are reported for a threshold of 0.5 on the prediction probabilities. We
352 do not consider cross-validation due to the computational overhead of running each experiment multiple
353 times. However, to ensure that results are reproducible and comparable on the same system we fix
354 random seeds for all components of the experimental setup that involve randomness (e.g., shuffling of
355 datasets, noisy hidden layers, weight initialization, etc.). We also enforce deterministic GPU floating
356 point calculations.

357

358

359 **4 Results**

360 *Baseline model (BM)*

361 Table 2 shows the results of different decoder-encoder combinations on the three test scenarios. Models
 362 have been trained on IKONOS R-G-B satellite images. The results on test scenario I, which compares
 363 against a test split of the reference data (same sensor, same task), show that models based on U-Net
 364 decoder perform better than Deeplab-V3+ models. This accounts for their accuracy (0.73 mean *IoU*
 365 across all U-Net models compared to 0.72 mean *IoU* for Deeplab-V3+) as well as inference speed (18.80
 366 mean *throughput* for U-Net compared to 5.03 mean *throughput* for Deeplab-V3+). Similarly, Mobilenet-
 367 V3 encoder models provide better results than ResNet-50 and EfficientNet-B4 (0.74 mean *IoU* and 20.5
 368 mean *throughput* for Mobilenet-V3 compared to 0.71 mean *IoU* and 6.91 mean *throughput* for
 369 EfficientNet-B4, 0.73 mean *IoU* and 8.30 mean *throughput* for ResNet-50).

370 Test scenario II compares model results against an independent test dataset of water scenes from satellite
 371 sensors that have not been used during training (different sensors, same task). All models show a clear
 372 performance decrease when transferred to other sensors (0.31 mean drop in *IoU* across all models). In
 373 this scenario models based on U-Net decoder perform better than Deeplab-V3+ models (0.41 mean *IoU*
 374 for U-Net compared to 0.40 mean *IoU* for Deeplab-V3+). Mobilenet-V3 encoder models provide better
 375 results than ResNet-50 and EfficientNet-B4 models (0.45 mean *IoU* for Mobilenet-V3 compared to 0.41
 376 mean *IoU* for EfficientNet-B4 and 0.37 mean *IoU* for ResNet-50).

377 *Table 2: Results for different baseline models trained on IKONOS R-G-B satellite images. Test scenario I*
 378 *compares against a test split of the reference data (same sensors, same task); Test scenario II against an*
 379 *independent test dataset of water scenes from different sensors (different sensors, same task); Test scenario III*
 380 *against an independent test dataset of flood scenes from different sensors (different sensors, different task).*

ID	Decoder	Encoder backbone	Through put (mp/s)	Test scenario I			Test scenario II			Test scenario III		
				<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
BM-0	U-Net	Mobilenet-V3	32.69	0.74	0.85	0.82	0.45	0.47	0.85	0.45	0.48	0.78
BM-1	U-Net	ResNet-50	13.20	0.73	0.85	0.81	0.38	0.39	0.76	0.41	0.43	0.78
BM-2	U-Net	EfficientNet-B4	10.58	0.71	0.89	0.77	0.41	0.41	0.80	0.42	0.45	0.78
BM-3	DL-V3+	Mobilenet-V3	8.40	0.73	0.85	0.82	0.45	0.46	0.86	0.42	0.44	0.80
BM-4	DL-V3+	ResNet-50	3.36	0.72	0.84	0.81	0.35	0.40	0.73	0.36	0.37	0.77
BM-5	DL-V3+	EfficientNet-B4	3.23	0.70	0.81	0.79	0.40	0.45	0.72	0.42	0.46	0.77

381

382 Test scenario III compares model results against an independent test dataset of flood scenes from sensors
 383 that have not been used during training (different sensors, different task). A general performance decrease
 384 can be observed compared to test scenario I. Similar to the previous tests, also in this scenario models
 385 based on U-Net perform better than Deeplab-V3+ models (0.43 mean *IoU* for U-Net compared to 0.40
 386 mean *IoU* for Deeplab-V3+). Mobilenet-V3 encoder models outperform ResNet-50 and EfficientNet-B4
 387 encoder models (0.44 mean *IoU* for Mobilenet-V3 compared to 0.42 mean *IoU* for EfficientNet-B4 and
 388 0.39 mean *IoU* for ResNet-50).

389 In summary, the U-Net Mobilenet-V3 model provides best results across all test scenarios while being
 390 small and fast during inference compared to the other models.

391 *Input bands (IB)*

392 Based on the results of the initial experiment, we use a U-Net model with Mobilenet-V3 encoder and
 393 train it on IKONOS satellite images with varying input bands. From Table 3 it can be seen that compared
 394 to the baseline (IB-0) adding the NIR spectral band (IB-1) clearly improves the segmentation results
 395 across all test scenarios by 0.07 *IoU* for test scenario I, 0.17 *IoU* for scenario II and 0.32 *IoU* for scenario
 396 III. Also adding slope information (IB-2) improves the baseline R-G-B model by 0.02 *IoU* for scenario
 397 I, 0.03 *IoU* for scenario II and 0.02 *IoU* for scenario III. The combined R-G-B with NIR and Slope model
 398 (IB-3) shows best performance and exceeds the baseline by 0.07 *IoU* for scenario I, 0.20 *IoU* for scenario
 399 II and 0.35 *IoU* for scenario III.

400 *Table 3: Results of a U-Net model with Mobilenet-V3 encoder trained on IKONOS satellite images with varying*
 401 *input bands.*

ID	Input bands	Throughput (mp/s)	Test scenario I			Test scenario II			Test scenario III		
			<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
IB-0	R-G-B	32.69	0.74	0.85	0.82	0.45	0.47	0.85	0.45	0.48	0.78
IB-1	R-G-B-NIR	32.74	0.81	0.92	0.87	0.62	0.79	0.78	0.77	0.83	0.93
IB-2	R-G-B-Slope	32.98	0.76	0.89	0.82	0.48	0.55	0.78	0.47	0.49	0.80
IB-3	R-G-B-NIR-Slope	33.00	0.81	0.92	0.87	0.65	0.84	0.73	0.80	0.87	0.91

402

403 *Data augmentation (DA)*

404 Table 4 depicts the influence of different data augmentation techniques on the segmentation results for
 405 the three test scenarios. Based on the results of previous experiments, we use a U-Net model with

406 Mobilenet-V3 encoder trained on IKONOS R-G-B-NIR-Slope satellite images. DA-0 uses the training
 407 dataset as is without augmentation. The results indicate that when augmentation (DA-1) is applied to the
 408 training dataset, better performance can be achieved on all test scenarios compared to using the
 409 unmodified training data. The effects of augmentation are, however, minor. DA-2 applies test-time
 410 augmentation. For test scenarios I and III, the results are comparable to DA-1. On test scenario II test-
 411 time augmentation shows negative effects on the accuracy. Since every tile needs to be augmented and
 412 predicted multiple times before the final prediction can be assigned, inference times of this approach are
 413 longer compared to DA-2.

414 *Table 4: Results of a U-Net model with Mobilenet-V3 encoder trained on IKONOS R-G-B-NIR-Slope satellite*
 415 *images for different data augmentation techniques.*

ID	AUG Train-time	AUG Test-time	Test scenario I			Test scenario II			Test scenario III		
			<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
DA-0	False	False	0.81	0.92	0.87	0.65	0.84	0.73	0.80	0.87	0.91
DA-1	True	False	0.82	0.92	0.87	0.66	0.92	0.75	0.82	0.90	0.91
DA-2	True	True	0.82	0.92	0.87	0.62	0.93	0.66	0.82	0.90	0.91

416

417 *Pre-training (PT)*

418 Table 5 shows the results of a U-Net model with Mobilenet-V3 encoder trained on IKONOS R-G-B-
 419 NIR-Slope satellite images with train-time augmentation under consideration of different pre-trained
 420 weights. When training and testing on data of the same sensor and task (test scenario I) no improvement
 421 can be observed by using pre-trained weights. However, the generalization ability of the model improves
 422 with higher *IoU* being measured on test scenarios II and III when ImageNet weights are used to initialize
 423 R-G-B bands. Initializing weights for all bands (R-G-B-NIR-Slope) by pre-training on Sentinel-2 images
 424 seems to have a negative impact on model transferability to scenarios II and III.

425 *Table 5: Results of a U-Net model with Mobilenet-V3 encoder trained on IKONOS R-G-B-NIR-Slope satellite*
 426 *images with train-time augmentation under consideration of different pre-trained weights.*

ID	Weights	Test scenario I			Test scenario II			Test scenario III		
		<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
PT-0	Random	0.82	0.92	0.87	0.66	0.92	0.75	0.82	0.90	0.91
PT-1	ImageNet (R-G-B) Random (NIR-Slope)	0.82	0.90	0.89	0.68	0.92	0.74	0.84	0.89	0.93

PT-2	Sentinel-2	0.82	0.91	0.88	0.63	0.75	0.80	0.76	0.82	0.92
-------------	------------	------	------	------	------	------	------	------	------	------

427

428 *Transfer model – satellite to aerial (TM)*

429 Table 6 tests the transferability of a U-Net model with Mobilenet-V3 encoder between satellite and aerial
430 images for different training setups. Since the aerial images used in this study do not have a NIR spectral
431 band available, we train and test with R-G-B-Slope images. Training and testing on satellite images (TM-
432 0) performs significantly better than testing the same model on aerial images for all three scenarios.
433 Training on Mapbox images (TM-1 and TM-2) performs better on aerial images than on satellite images.
434 The performance difference between aerial and satellite images is, however, much smaller compared to
435 TM-0. Training Mapbox with pre-trained weights from satellite images (TM-2) shows better
436 performance than using ImageNet weights. The best overall performance, however, can be reached when
437 a model is trained on a joined training dataset of contrast enhanced satellite images and Mapbox images
438 (TM-3). This model shows the best performance on both satellite and aerial images across all test
439 scenarios.

440 *Table 6: Transferability of a U-Net model with Mobilenet-V3 encoder between satellite and aerial images for*
441 *different training setups.*

ID	Training	Testing	Test scenario I			Test scenario II			Test scenario III		
			IoU	Prec	Rec	IoU	Prec	Rec	IoU	Prec	Rec
TM-0	Satellite R-G-B-Slope, pretrained on ImageNet, train-time augment.	Aerial	0.28	0.31	0.81	0.10	0.12	0.88	0.15	0.25	0.31
		Satellite	0.77	0.90	0.82	0.63	0.76	0.79	0.67	0.78	0.81
TM-1	Mapbox R-G-B-Slope, pretrained on ImageNet, train-time augment.	Aerial	0.70	0.83	0.79	0.54	0.82	0.62	0.64	0.88	0.69
		Satellite	0.67	0.76	0.84	0.54	0.60	0.79	0.62	0.72	0.84
TM-2	Mapbox R-G-B-Slope, pretrained on TM-0 train-time augment.	Aerial	0.72	0.84	0.80	0.56	0.89	0.57	0.62	0.74	0.83
		Satellite	0.69	0.84	0.78	0.55	0.54	0.91	0.64	0.69	0.89
TM-3	TM-0 and TM-1 trained together with contrast enhancement	Aerial	0.72	0.84	0.80	0.78	0.82	0.82	0.70	0.86	0.74
		Satellite	0.79	0.89	0.87	0.70	0.78	0.89	0.82	0.89	0.90

442

443 Since contrast enhancement of the satellite images seems to have a positive effect on the segmentation
444 performance, we also applied this additional pre-processing step to the NIR spectral band and trained a
445 model on contrast enhanced R-G-B-NIR bands and slope information. Table 7 indicates that when
446 training and testing on data of the same sensor, the effect of contrast enhancement is not seen and

447 performances on test scenario I are equal between PT-1 and PT-1 Contrast. However, the model trained
 448 on contrast enhanced data (PT-1 Contrast) transfers better to test scenarios II and III than the same model
 449 trained on the original top of canopy reflectance values (PT-1).

450 *Table 7: Influence of contrast enhancement on best performing U-Net model with Mobilenet-V3 encoder with pre-*
 451 *trained ImageNet weights and train-time augmentation for R-G-B-NIR-Slope satellite images (PT-1).*

ID	Weights	Test scenario I			Test scenario II			Test scenario III		
		<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
PT-1	ImageNet (R-G-B) Random (NIR-Slope)	0.82	0.90	0.89	0.68	0.92	0.74	0.84	0.89	0.93
PT-1 Contrast	ImageNet (R-G-B) Random (NIR-Slope) Contrast enhancement	0.82	0.90	0.89	0.72	0.82	0.86	0.85	0.92	0.92

452

453 To further evaluate the usefulness of our results for flood mapping applications we also tested the
 454 influence of adding samples of the target (flood water) domain to the training data. Due to the limited
 455 availability of flood water samples, we divided the six available flood images from test scenario III into
 456 training and test splits. The best-performing model for aerial and satellite images (TM-3) is then retrained
 457 including data from the target (flood water) domain. Specifically, we add two of the flood images to the
 458 training data and use the remaining four flood images as independent test data. Table 8 provides an
 459 overview of the data scenarios and compares the results for simply applying the TM-3 model (trained
 460 without images of the target domain) with TM-3 Flood (trained including images of the target domain).
 461 It can be seen that adding few flood water samples during training could slightly improve the results on
 462 test images that show flood situations. An increase in *Rec* of 0.02 indicates that the improvement is
 463 mainly related to more flood water pixels being correctly detected.

464 *Table 8: Influence of retraining TM-3 with limited training data of the target (flood water) domain.*

ID	Training	Testing	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
TM-3	TM-3 without flood water images	Landes, France (GeoEye-1), Brück, Germany (WorldView-3), Altenahr, Germany (DLR MACS), Dresden, Germany (DLR 3K)	0.73	0.89	0.80
TM-3 Flood	TM-3 with flood water images from Bihar, India (GeoEye-1) and Minden, Germany (DLR 4K)	Landes, France (GeoEye-1), Brück, Germany (WorldView-3), Altenahr, Germany (DLR MACS), Dresden, Germany (DLR 3K)	0.74	0.89	0.82

465

466 Finally, we compared the influence of image resolution on the segmentation performance. All previous
 467 experiments used the original resolution of the test images. As can be seen from Table 9, resampling of
 468 the test images leads to worse performance on all test scenarios. Predicting on the original image
 469 resolution, despite varying spatial resolution across the test images seems to be the best choice. This is
 470 further underlined by the observation that the closer the resampled resolution is to the average original
 471 image resolution (0.8 m for scenario I and 0.5 m for scenarios II and III) the better the results.

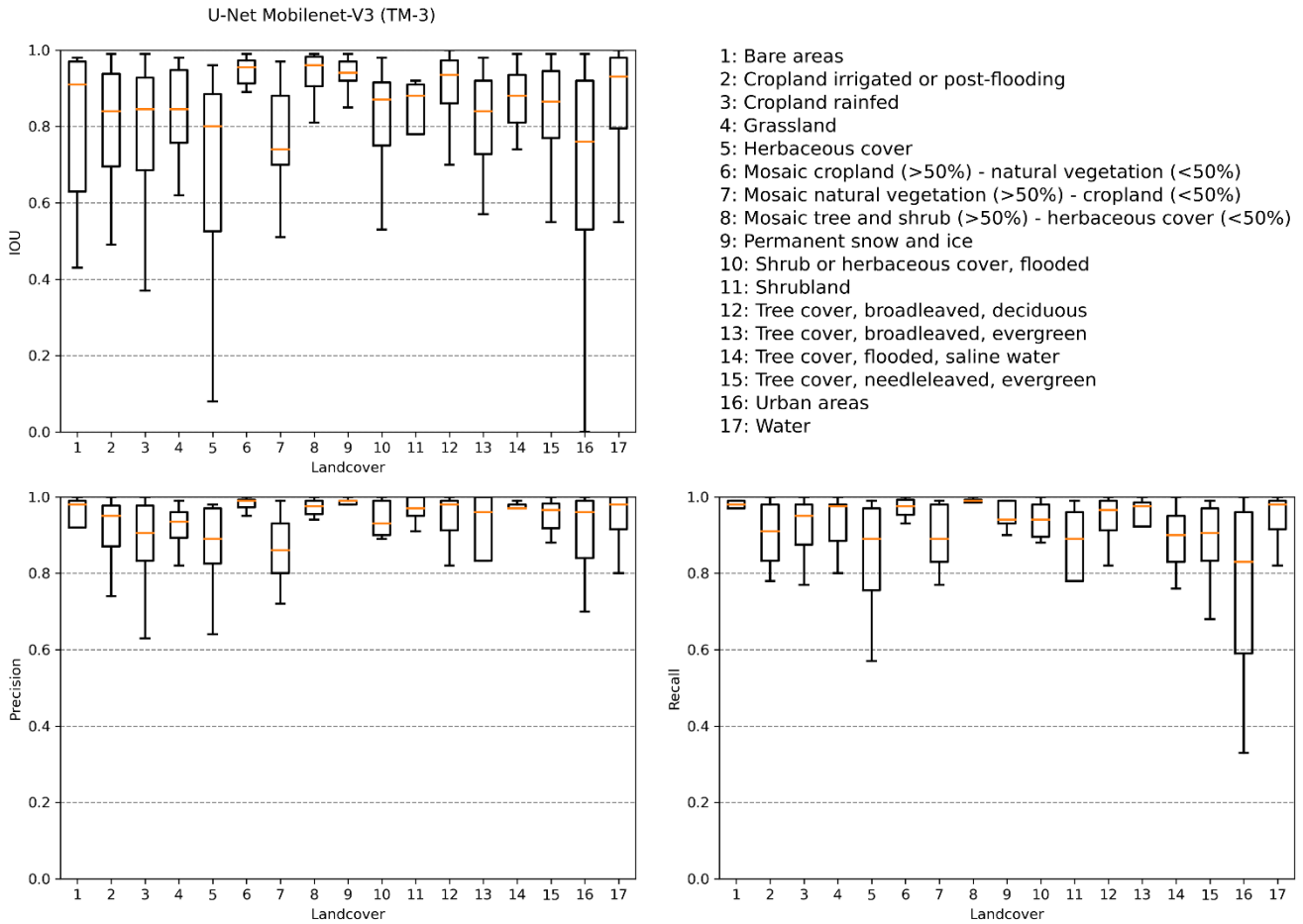
472 *Table 9: Influence of image spatial resolution on best performing U-Net model with Mobilenet-V3 encoder, pre-*
 473 *trained ImageNet weights and train-time augmentation for contrast enhanced R-G-B-NIR-Slope satellite images*
 474 *(PT-1 Contrast).*

ID	Resolution (m)	Test scenario I			Test scenario II			Test scenario III		
		<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>	<i>IoU</i>	<i>Prec</i>	<i>Rec</i>
RES-0	Varying original	0.82	0.90	0.89	0.72	0.82	0.86	0.85	0.92	0.92
RES-1	0.25	0.76	0.81	0.91	0.68	0.79	0.83	0.80	0.85	0.90
RES-2	0.50	0.80	0.87	0.89	0.72	0.84	0.83	0.84	0.90	0.92
RES-3	0.75	0.82	0.89	0.90	0.62	0.73	0.76	0.79	0.91	0.82
RES-4	1.00	0.80	0.89	0.87	0.56	0.70	0.70	0.76	0.90	0.80

475

476 *Summary of results for best performing model*

477 Figure 3 shows the results of the best-performing U-Net Mobilenet-V3 model with pre-trained ImageNet
 478 weights trained on contrast enhanced and augmented R-G-B-Slope IKONOS and Mapbox data (TM-3)
 479 on test scenario I grouped by predominant landcover of the test samples (aerial and satellite test datasets
 480 combined). Overall good performance across all landcover types can be observed with median values for
 481 all metrics above 0.75. Performance over the majority of landcover types, moreover, shows relatively
 482 small variation. Few landcover types, however, highlight a larger spread of the performance with lower
 483 quartiles for *IoU* being below 0.70. This becomes evident in particular for bare areas (class 1), herbaceous
 484 cover (class 5) and urban areas (class 16). A closer look at *Prec* and *Rec* reveals that largely *Prec* is
 485 affected for water bodies located in bare areas, which indicates a tendency to overestimate water in these
 486 areas. Herbaceous cover and urban areas show a larger spread of *Rec* values, which indicates a tendency
 487 to underestimate water.

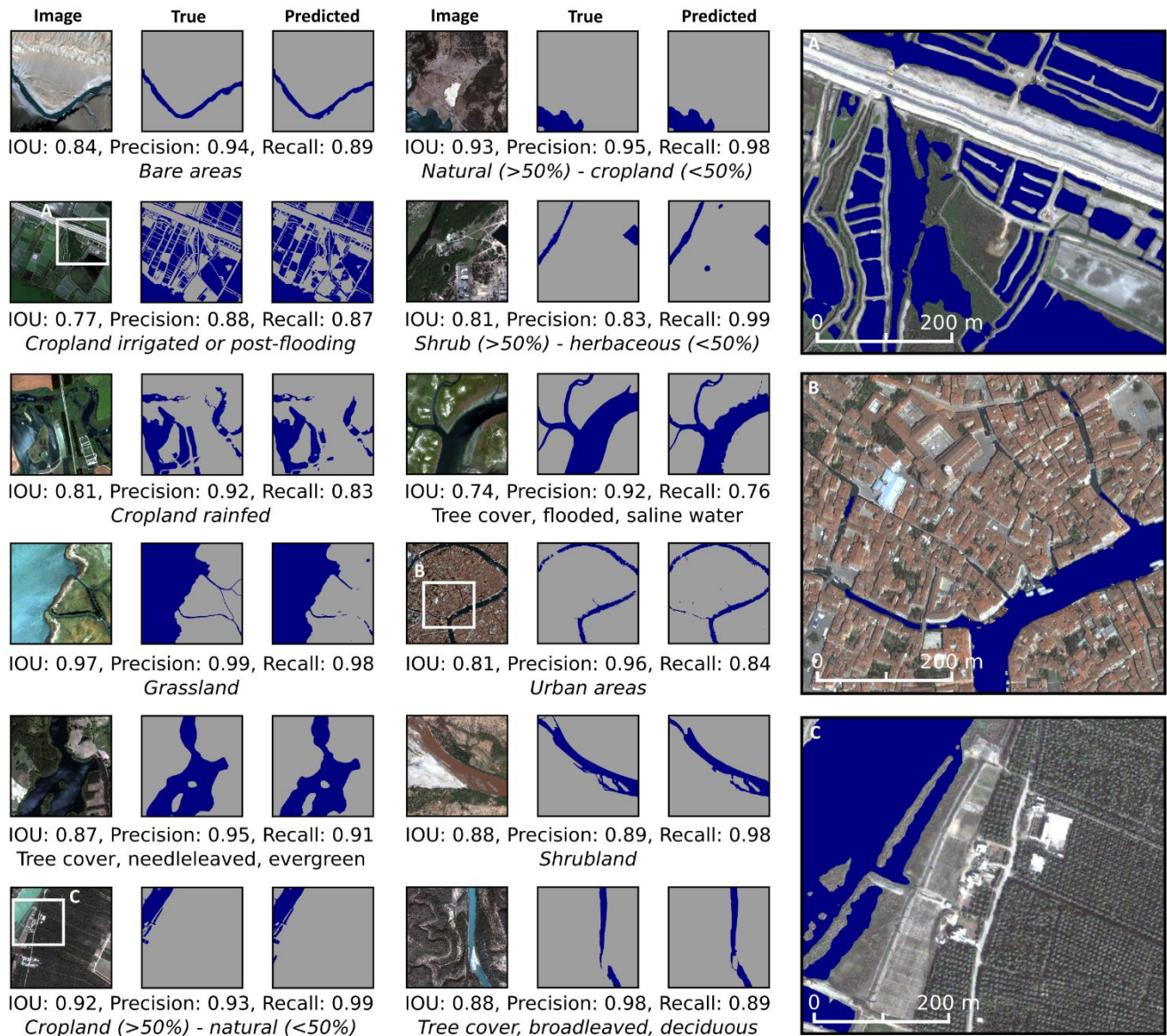


488

489 *Figure 3: Results of best-performing U-Net Mobilenet-V3 model with pre-trained ImageNet weights trained on*
 490 *contrast enhanced and augmented R-G-B-Slope IKONOS and Mapbox images (TM-3) on test scenario I (aerial*
 491 *and satellite test datasets combined) grouped by predominant landcover.*

492 Figure 4 shows results of the best-performing overall model (TM-3) for selected scenario I test images
 493 per landcover. All water predictions show good performance and match well the reference masks. Some
 494 minor issues with false positives and negatives can be observed for example in water look-alike areas
 495 (e.g., shadows or moist soil) or where a clear border between water and land is difficult to define even
 496 by visual image interpretation. The “cropland irrigated or post-flooding” sample (zoom A) depicts a
 497 complex scene with a mixture of different shaped and coloured water basins and canals. Despite the large
 498 variety of water appearance in the scene, the overall accuracy of the predicted water mask is good with
 499 an *IoU* of 0.77, *Prec* and *Rec* above 0.85. In the “urban areas” sample (zoom B) it can be seen that narrow
 500 water canals of central Venice in between dense built-up areas are not sufficiently segmented thus leading
 501 to false negatives and a lower *Rec*. Similar effects can be observed in other urban samples (see also
 502 Figure 3). Despite the very good model performance on the “cropland (>50%) – natural (<50%)” sample
 503 (zoom C), it is a typical example for a land-water border that even by visual image interpretation is not

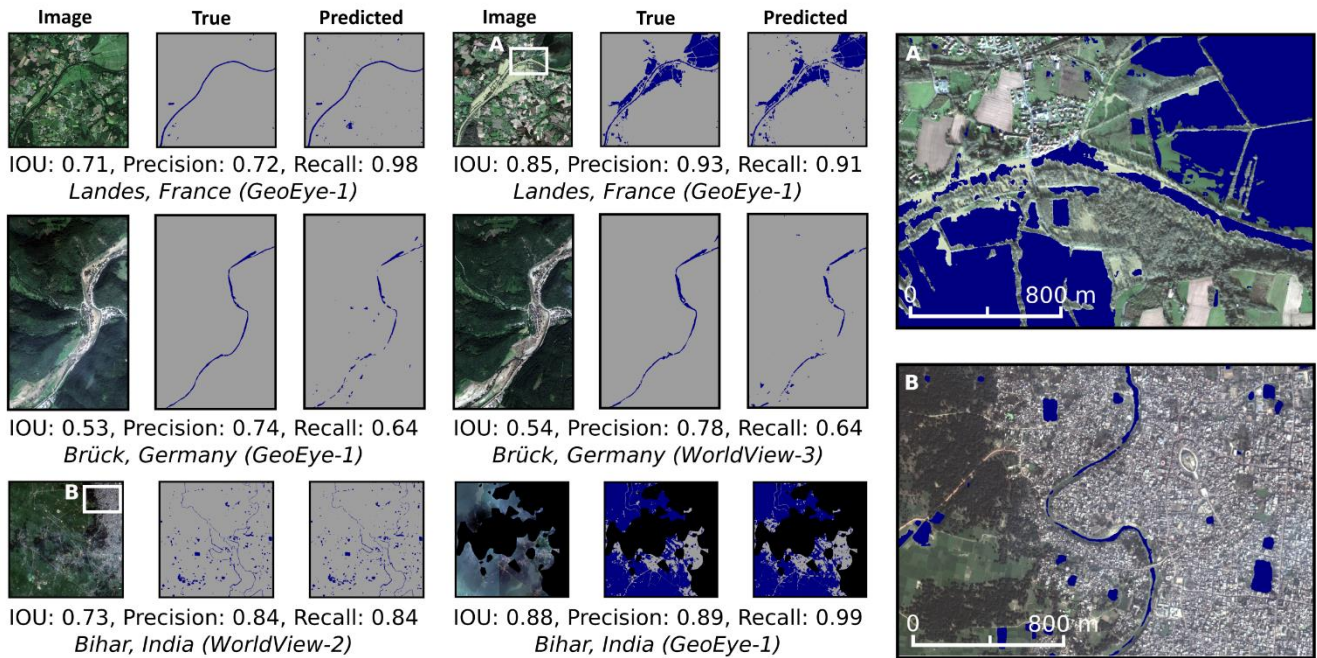
504 clearly identifiable. Especially shallow water areas and muddy river banks in the lower left corner make
 505 it difficult to identify a clear land-water border.



506
 507 *Figure 4: Results of best-performing model (TM-3) for randomly selected scenario I test images per landcover.*
 508 *Water is outlined in blue, background in grey. Each sample covers an area of approximately 1.3 km².*

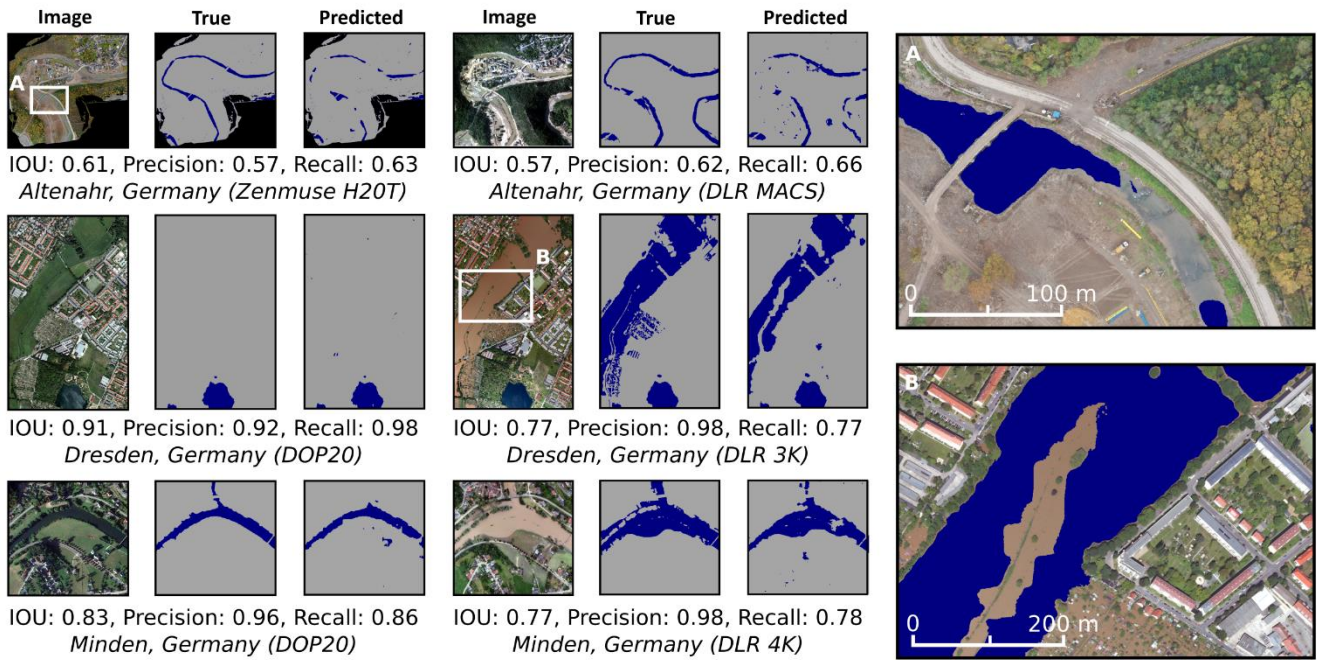
509 Figures 5 shows results of the U-Net Mobilenet-V3 model with pre-trained ImageNet weights trained on
 510 contrast enhanced and augmented IKONOS and Mapbox data (TM-3) for scenario II and III satellite test
 511 images. The flood water image of Landes acquired by the GeoEye-1 satellite highlights the overall good
 512 performance of the model across water bodies with heterogeneous appearance. False negatives are
 513 mainly observed in water with high sediment content and along fuzzy land-water borders with partially
 514 overlapping vegetation. The normal water WorldView-2 scene of Bihar (zoom B) indicates good

515 performance with only few false positives, despite some visible building shadows. Water bodies inside
 516 the city are segmented with only minor misclassifications even in muddy areas along the main river.



517
 518 *Figure 5: Results of best-performing model (TM-3) for scenario II (left) and III (right) satellite test images. Water*
 519 *is outlined in blue, background in grey and no-data areas in black. The area covered by each image varies and is*
 520 *depicted in Table 1.*

521 Figure 6 shows results of the U-Net Mobilenet-V3 model with pre-trained ImageNet weights trained on
 522 contrast enhanced and augmented IKONOS and Mapbox data (TM-3) for scenario II and III aerial test
 523 images. The normal water Zenmuse H20T image acquired by UAV of Altenahr (zoom A) shows
 524 problems with false negatives related to shallow transparent water that may cause confusion with the
 525 sandy river bottom. This is reflected in the low *Rec* values for this scene. The DLR 3K aerial flood water
 526 scene of Dresden (zoom B) shows low *Rec* values and has problems with false negatives that are likely
 527 related to the large amount of sediment in the water in combination with partially inundated vegetation
 528 and allotment gardens.



529

530 *Figure 6: Results of best-performing model (TM-3) for scenario II (left) and III (right) aerial test images. Water*
531 *is outlined in blue, background in grey and no-data areas in black. The area covered by each image varies and is*
532 *depicted in Table 1.*

533 **5 Discussion**

534 The objective of this study is to train a model for semantic segmentation of normal water and flood water
535 bodies in very high-resolution satellite and aerial images. Performance evaluation of CNN architectures
536 and encoders indicate superior performance of a U-Net Mobilenet-V3 model when considering *IoU*,
537 *Prec*, *Rec* and *throughput* (Table 2). Compared to other studies, we focus specifically on multi-sensor
538 generalization ability, global applicability, simplicity and processing speed to produce timely, accurate
539 and relevant information under varying data availability scenarios. Our work is thus targeted towards
540 operational use in rapid mapping activities to support situational awareness in flood emergency response.
541 Similar to Yang et al. (2020) we can report good performance for all tested models with a ResNet-50
542 encoder. Following their conclusions, we have not tested the more complex ResNet-100, which
543 performed worse in their experiments. Instead we focussed on the more recently developed Mobilenet-
544 V3 and EfficientNet-B4, which showed comparable results than ResNet-50 at higher *throughput*.
545 Contrary to the findings of Duan and Hu (2020), who compare a multi-scale refinement network with U-
546 Net, SegNet and DeepLab-V3+ networks for water segmentation in WorldView-3 and GaoFen-2 satellite
547 images, our experiments indicate that U-Net-based models outperform models based on the more
548 complex DeepLab-V3+ architecture both in terms of accuracy metrics and *throughput*. While the

549 difference in *IoU* between the two architectures is minimal in case of test scenario I (same sensor, same
550 task), it becomes more evident in test scenarios II (different sensor, same task) and III (different sensor,
551 different task). Further performance improvements for DeepLab-V3+ could be achieved by longer
552 training without early stopping. More computational resources would be required for a more thorough
553 comparison of the architectures, since deeper architectures like DeepLab-V3+ commonly take longer to
554 converge and hence their training in our experiments may be sub-optimal given the relatively short
555 patience before learning rate reduction and early stopping. Independent on training times, the main
556 decision criterion in this study to select a U-Net-based model has been the significantly higher throughput
557 during inference.

558 Duan and Hu (2020) use the Deepglobe dataset (Demir et al., 2018) and expand it with a custom dataset
559 of GaoFen-2 images over China. We found the Deepglobe dataset not to be suitable for our study, due to
560 limited geographical coverage and relatively poor annotation quality for the water class. Therefore, we
561 acquired a comprehensive multi-sensor reference dataset that we used to train, validate and test the water
562 segmentation. The reference dataset has been sampled globally to account for a large variety of climatic,
563 atmospheric and land-cover conditions. To overcome the limited geographical coverage of the IKONOS
564 data archive that we had access to, we complement the reference dataset with freely available optical
565 images from a Web Map Service and water masks from OpenStreetMap. Nevertheless, even the
566 combined reference dataset shows a geographical bias towards Europe, Middle East and India (Figure
567 1). Also, the acquired independent test datasets for other sensors (test scenario II) and tasks (test scenario
568 III) show a geographical bias towards Europe and India. Similar to the reference dataset (test scenario I),
569 this is caused by limited availability of suitable test images from other locations. The aim in this study
570 was to compile test data from different sensors covering the same location during normal and flood water
571 situations. This constraint strongly reduced the number of available and accessible samples. The dataset
572 should be further extended and diversified in future to improve the significance of the results regarding
573 model transferability.

574 To allow for a high degree of transferability across sensors, we focussed on the use of R-G-B image
575 bands, which are acquired by almost all available optical sensors independent on the acquisition platform.
576 For sensor-specific applications, however, including additional spectral bands like the NIR band into the
577 input feature space has clearly a positive effect on the performance (Table 3 and Table 7). Additionally,
578 slope information derived from a DEM can further improve results. Since high-quality DEMs like the
579 Copernicus DEM are freely available at global coverage, slope information can be added independently
580 for all sensors.

581 Sample water masks are derived by support of manual interpretation and digitization of image tiles,
582 which means that the results are to be regarded as relative to the performance of a human analyst and not
583 to an absolute ground-truth. Especially the fractal geometry of the border between land and water may
584 lead to label inconsistencies when trying to define a hard class boundary. Their manual delineation is at
585 least partially a matter of subjective interpretation and may introduce a bias in the reference dataset and
586 hence into the performance results. In this context, it also needs to be considered which level of detail
587 and resolution is actually required by an application. Wetland monitoring for example may have different
588 requirements compared to flood emergency response. In this study we focus on emergency response and
589 have followed as much as possible best-practises of rapid mapping procedures as being conducted at the
590 Centre for Satellite-based Crisis Information (ZKI) of the German Aerospace Center (DLR) (Lechner
591 and Gähler, 2017).

592 Results show that more efforts should be spend on preparing high-quality training data and dedicate more
593 research towards reference data preparation and its incremental improvement. The target domain is
594 affected by varying scene and image properties that are dependent on sensor characteristics (e.g.,
595 radiometry, spatial resolution), atmospheric conditions, land-use / land-cover of the background class
596 and appearance of the water class. Despite our aim to cover such variations as part of the reference
597 dataset, we decided to apply image augmentation to artificially increase the training sample size and to
598 cover a larger range of conditions that may occur during inference in real-world applications. The need
599 for augmentation has been reinforced by the outcomes of an experiment to test the influence of different
600 levels and combinations of data augmentation on the segmentation results (Table 4). Although we are
601 not specifically isolating the effects of single augmentations (contrast, brightness, scale and rotation) on
602 the performance, it can be noted that scaling images during augmentation seems to improve the
603 performance on images of a fraction of the original resolution. Test-time augmentation as well as
604 resampling images before prediction (Table 9) did not improve results. Predicting on the original image
605 and resolution seems to be the best choice, despite varying spatial resolutions across the test images.

606 Similar to what is reported by Gebrehiwot et al. (2019) for the task of extracting flooded areas from UAV
607 images, our results confirm that adapting pre-trained models is highly beneficial and can lead to superior
608 performance compared to training with randomly initialized weights. While He et al. (2018) challenge
609 the commonly used technique of pre-training on ImageNet and fine-tuning on custom datasets, we can
610 report that for our use case pre-training on not remote sensing specific image databases like ImageNet
611 helped to improve the results. Pre-training on a domain-specific water dataset of lower resolution
612 Sentinel-2 images did not improve performance on tests scenario I and even produced lower accuracies

613 when transferred on test scenario II and III compared to a model with random weight initialization (Table
614 5). A possible explanation for this behaviour could be that pre-training on Sentinel-2 emphasizes the
615 importance of spectral intensity values for the final model, which could hamper the transferability to
616 other sensors and tasks. In contrast, models pre-trained on ImageNet are known to be biased towards
617 texture, which may increase the robustness towards a range of image distortions and thus support model
618 transferability (Geirhos et al., 2019). This assumption is further supported by the observation that
619 applying percentage linear contrast stretch improved model accuracies on test scenarios II and III (Table
620 6 and Table 7). Similar to the findings of Kaiser et al. (2017)), our results indicate that by adding noisy
621 large-scale training data (in our case from Mapbox images with OpenStreetMap annotations)
622 performance improvements can be achieved with relatively little manual annotation effort.

623 Moreover, we could show that adding a limited number of training samples from the target (flood water)
624 domain, could further improve results on test scenario III (Table 8). This is particularly relevant for
625 applications targeted towards flood disaster response. More flood samples are, however, required to
626 further test and verify the effectiveness of this approach. Further approaches to improve model accuracy,
627 especially on test scenarios II and III, include fine-tuning on limited amounts of domain-specific training
628 data (e.g., from flood events and / or other sensors), multi-source domain adaptation by data
629 standardization (Tasar et al., 2020) or few-shot learning (Rußwurm et al., 2020).

630 **6 Conclusions**

631 In this study, we combined different architectures (U-Net and DeepLab-V3+) with various encoder
632 backbones (Mobilenet-V3, ResNet-50 and EfficientNet-B4) and tested their ability to delineate normal
633 and flood water under varying environmental conditions and data availability scenarios. We introduced
634 a reference dataset of 1,120 globally sampled images with quality checked binary water masks to train,
635 validate and test the CNN models. Compared to previous studies, we considered a large variety of
636 available satellite and aerial sensors as input, namely IKONOS, GeoEye-1, WorldView-2, WorldView-
637 3 and four different airborne cameras. Across several experiments we identified the superior performance
638 of a U-Net Mobilenet-V3 model with initial weights being pre-trained on ImageNet. The best-performing
639 model indicates good generalization ability across sensors and varying environmental conditions. In this
640 context, not only the combination of architecture, encoder and weight initialization is relevant, but also
641 the way training data are augmented. While a model trained with R-G-B-NIR-Slope input feature space
642 produces the highest accuracies on all test scenarios (PT-1 Contrast), models without the NIR spectral
643 band can achieve comparable results (TM-3). The benefit of a model that requires only R-G-B-Slope

644 input features is that it is applicable across a larger spectrum of sensors (e.g., none of the aerial images
645 used in this study had an additional NIR spectral band available). Contrast enhancement and noisy
646 training data from Mapbox and OpenStreetMap could improve performance on all test scenarios. By
647 successfully applying the model to four exemplary flood events, we could highlight the usefulness of this
648 work for rapid mapping activities to support situational awareness in emergency response. In particular,
649 we could show that it is possible to train a model that is able to cope with highly diverse data availability
650 scenarios in disaster situations. Our experiments further indicate that in order to improve flood mapping
651 results it is beneficial to add even limited amounts of flood water samples during model training. More
652 work is, however, required to further adapt to specific scene properties that are common during floods
653 and that have caused misclassifications in the test images. This includes high degrees of sediment in the
654 water, strong sun-light reflections on the water surface and presence of cloud and terrain shadows.
655 Ongoing and future works focus on acquiring more samples for fine-tuning models with domain-specific
656 training data from flood events and applying the method on large scale datasets of past floods (e.g., flood
657 in Germany 2013 and 2021). Finally, the proposed method will support complementation of existing
658 systematic water and flood monitoring services based on Sentinel-1 (Helleis et al., 2022; Twele et al.,
659 2016) and Sentinel-2 (Wieland and Martinis, 2019) with an ad-hoc component for on-demand very high-
660 resolution optical satellite and aerial imagery.

661 **Funding:** This work was supported by the German Federal Ministry of Education and Research (BMBF) as part
662 of the AIFER project [Grant No. 13N15525] and DLR-internal funding.

663 **Conflicts of interest:** The authors declare no conflict of interest. The founding sponsors had no role in the design
664 of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and in the
665 decision to publish the results.

666 **References**

- 667 Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M., 2019. Optuna: A Next-generation
668 Hyperparameter Optimization Framework. arXiv:1907.10902 [cs, stat].
- 669 Azimi, S.M., Kiefl, R., Gstaiger, V., Bahmanyar, R., Merkle, N., Henry, C., Rosenbaum, D., Kurz, F.,
670 2021. Automatic object segmentation to support crisis management of large-scale events. *Int.*
671 *Arch. Photogramm. Remote Sens. Spatial Inf. Sci.* XLIII-B2-2021, 433–440.
672 <https://doi.org/10.5194/isprs-archives-XLIII-B2-2021-433-2021>
- 673 Ball, J.E., Anderson, D.T., Sr, C.S.C., 2017. Comprehensive survey of deep learning in remote sensing:
674 theories, tools, and challenges for the community. *JARS* 11, 042609.
675 <https://doi.org/10.1117/1.JRS.11.042609>
- 676 Bonafilia, D., Tellman, B., Anderson, T., Issenberg, E., 2020. Sen1Floods11: a georeferenced dataset to
677 train and test deep learning flood algorithms for Sentinel-1, in: 2020 IEEE/CVF Conference on
678 Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 835–845.

679 <https://doi.org/10.1109/CVPRW50498.2020.00113>

680 Buchhorn, M., Lesiv, M., Tsendbazar, N.-E., Herold, M., Bertels, L., Smets, B., 2020. Copernicus
681 Global Land Cover Layers—Collection 2. *Remote Sensing* 12, 1044.
682 <https://doi.org/10.3390/rs12061044>

683 Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M., 2021. Swin-Unet: Unet-like
684 Pure Transformer for Medical Image Segmentation. *arXiv:2105.05537v1 [eess.IV]*.

685 Castillo-Navarro, J., Le Saux, B., Boulch, A., Audebert, N., Lefèvre, S., 2021. Semi-supervised
686 semantic segmentation in Earth Observation: the MiniFrance suite, dataset analysis and multi-
687 task network study. *Mach Learn.* <https://doi.org/10.1007/s10994-020-05943-y>

688 Chen, L.-C., Papandreou, G., Schroff, F., Adam, H., 2017. Rethinking Atrous Convolution for
689 Semantic Image Segmentation. *arXiv:1706.05587 [cs]*.

690 Chen, Y., Fan, R., Yang, X., Wang, J., Latif, A., 2018. Extraction of Urban Water Bodies from High-
691 Resolution Remote-Sensing Imagery Using Deep Learning. *Water* 10, 585.
692 <https://doi.org/10.3390/w10050585>

693 Demir, I., Koperski, K., Lindenbaum, D., Pang, G., Huang, J., Basu, S., Hughes, F., Tuia, D., Raskar,
694 R., 2018. DeepGlobe 2018: A Challenge to Parse the Earth through Satellite Images. 2018
695 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)
696 172–17209. <https://doi.org/10.1109/CVPRW.2018.00031>

697 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. ImageNet: A large-scale hierarchical
698 image database, in: 2009 IEEE Conference on Computer Vision and Pattern Recognition. pp.
699 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>

700 Ding, L., Lin, D., Lin, S., Zhang, J., Cui, X., Wang, Y., Tang, H., Bruzzone, L., 2022. Looking Outside
701 the Window: Wide-Context Transformer for the Semantic Segmentation of High-Resolution
702 Remote Sensing Images. *IEEE Transactions on Geoscience and Remote Sensing* 60, 1–13.
703 <https://doi.org/10.1109/TGRS.2022.3168697>

704 Duan, L., Hu, X., 2020. Multiscale Refinement Network for Water-Body Segmentation in High-
705 Resolution Satellite Imagery. *IEEE Geosci. Remote Sensing Lett.* 17, 686–690.
706 <https://doi.org/10.1109/LGRS.2019.2926412>

707 Fahrland, E., Jacob, P., Schrader, H., Kahabka, H., 2020. Copernicus Digital Elevation Model Product
708 Handbook (No. AO/1-9422/18/I-LG), Product handbook. Airbus.

709 Feng, W., Sui, H., Huang, W., Xu, C., An, K., 2019. Water Body Extraction From Very High-
710 Resolution Remote Sensing Imagery Using Deep U-Net and a Superpixel-Based Conditional
711 Random Field Model. *IEEE Geosci. Remote Sensing Lett.* 16, 618–622.
712 <https://doi.org/10.1109/LGRS.2018.2879492>

713 Gebrehiwot, A., Hashemi-Beni, L., Thompson, G., Kordjamshidi, P., Langan, T., 2019. Deep
714 Convolutional Neural Network for Flood Extent Mapping Using Unmanned Aerial Vehicles
715 Data. *Sensors* 19, 1486. <https://doi.org/10.3390/s19071486>

716 Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W., 2019. ImageNet-
717 trained CNNs are biased towards texture; increasing shape bias improves accuracy and
718 robustness. *arxiv.org/abs/1811.12231 [cs, q-bio, stat]*.

719 Gomes, R., Rozario, P., Adhikari, N., 2021. Deep Learning optimization in remote sensing image
720 segmentation using dilated convolutions and ShuffleNet, in: 2021 IEEE International
721 Conference on Electro Information Technology (EIT). pp. 244–249.
722 <https://doi.org/10.1109/EIT51626.2021.9491910>

723 Gu, X., Li, S., Ren, S., Zheng, H., Fan, C., Xu, H., 2022. Adaptive enhanced swin transformer with U-
724 net for remote sensing image segmentation. *Computers and Electrical Engineering* 102, 108223.
725 <https://doi.org/10.1016/j.compeleceng.2022.108223>

726 Guo, H., He, G., Jiang, W., Yin, R., Yan, L., Leng, W., 2020. A Multi-Scale Water Extraction

727 Convolutional Neural Network (MWEN) Method for GaoFen-1 Remote Sensing Images. *IJGI*
728 9, 189. <https://doi.org/10.3390/ijgi9040189>

729 Gupta, R., Goodman, B., Patel, N., Hosfelt, R., Sajeev, S., Heim, E., Doshi, J., Lucas, K., Choset, H.,
730 Gaston, M., 2019. Creating xBD: A Dataset for Assessing Building Damage from Satellite
731 Imagery. arXiv:1911.09296 [cs.CV].

732 Haklay, M., 2010. How Good is Volunteered Geographical Information? A Comparative Study of
733 OpenStreetMap and Ordnance Survey Datasets. *Environment and Planning B: Planning and*
734 *Design* 37, 682–703. <https://doi.org/10.1068/b35097>

735 Hänsch, R., Arndt, J., Lunga, D., Gibb, M., Pedelose, T., Boedihardjo, A., Petrie, D., Bacastow, T.M.,
736 2022. SpaceNet 8 - The Detection of Flooded Roads and Buildings. 2022 IEEE/CVF
737 Conference on Computer Vision and Pattern Recognition Workshops (CVPRW).

738 He, K., Girshick, R., Dollár, P., 2018. Rethinking ImageNet Pre-training. arXiv:1811.08883v1 [cs.CV].

739 He, K., Zhang, X., Ren, S., Sun, J., 2015a. Deep Residual Learning for Image Recognition.
740 arXiv:1512.03385 [cs].

741 He, K., Zhang, X., Ren, S., Sun, J., 2015b. Delving Deep into Rectifiers: Surpassing Human-Level
742 Performance on ImageNet Classification. arXiv:1502.01852 [cs].

743 Helleis, M., Wieland, M., Krullikowski, C., Martinis, S., Plank, S., 2022. Sentinel-1-based water and
744 flood mapping: benchmarking convolutional neural networks against an operational rule-based
745 processing chain. *Journal of Selected Topics in Applied Earth Observations and Remote*
746 *Sensing* 15, 2023–2036.

747 Howard, A., Sandler, M., Chen, B., Wang, W., Chen, L.-C., Tan, M., Chu, G., Vasudevan, V., Zhu, Y.,
748 Pang, R., Adam, H., Le, Q., 2019a. Searching for MobileNetV3, in: 2019 IEEE/CVF
749 International Conference on Computer Vision (ICCV). IEEE, Seoul, Korea (South), pp. 1314–
750 1324. <https://doi.org/10.1109/ICCV.2019.00140>

751 Hu, J., Shen, L., Albanie, S., Sun, G., Wu, E., 2020. Squeeze-and-Excitation Networks. *IEEE*
752 *Transactions on Pattern Analysis and Machine Intelligence* 42, 2011–2023.
753 <https://doi.org/10.1109/TPAMI.2019.2913372>

754 Huang, C., Chen, Y., Zhang, S., Wu, J., 2018. Detecting, Extracting, and Monitoring Surface Water
755 From Space Using Optical Sensors: A Review. *Rev. Geophys.* 56, 333–360.
756 <https://doi.org/10.1029/2018RG000598>

757 Iglovikov, V., Mushinskiy, S., Osin, V., 2017. Satellite Imagery Feature Detection using Deep
758 Convolutional Neural Network: A Kaggle Competition. arXiv:1706.06169 [cs].

759 Ireland, G., Volpi, M., Petropoulos, G., 2015. Examining the Capability of Supervised Machine
760 Learning Classifiers in Extracting Flooded Areas from Landsat TM Imagery: A Case Study
761 from a Mediterranean Flood. *Remote Sensing* 7, 3372–3399.
762 <https://doi.org/10.3390/rs70303372>

763 Kaiser, P., Wegner, J.D., Lucchi, A., Jaggi, M., Hofmann, T., Schindler, K., 2017. Learning Aerial
764 Image Segmentation From Online Maps. *IEEE Transactions on Geoscience and Remote*
765 *Sensing* 55, 6054–6068. <https://doi.org/10.1109/TGRS.2017.2719738>

766 Kingma, D.P., Lei, J., 2015. Adam: A Method for Stochastic Optimization. ArXiv:1412.6980v9 15.

767 Lechner, K., Gähler, M., 2017. Earth observation based crisis information — Emergency mapping
768 services and recent operational developments, in: 2017 4th International Conference on
769 Information and Communication Technologies for Disaster Management (ICT-DM). pp. 1–7.
770 <https://doi.org/10.1109/ICT-DM.2017.8275682>

771 Li, Y., Martinis, S., Wieland, M., 2019. Urban flood mapping with an active self-learning
772 convolutional neural network based on TerraSAR-X intensity and interferometric coherence.
773 *ISPRS Journal of Photogrammetry and Remote Sensing* 152, 178–191.
774 <https://doi.org/10.1016/j.isprsjprs.2019.04.014>

775 Long, Y., Xia, G.-S., Li, S., Yang, W., Yang, M.Y., Zhu, X.X., Zhang, L., Li, D., 2021. On Creating
776 Benchmark Dataset for Aerial Image Interpretation: Reviews, Guidances, and Million-AID.
777 IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 14, 4205–
778 4230. <https://doi.org/10.1109/JSTARS.2021.3070368>

779 Mapbox Satellite: global base map & satellite imagery [WWW Document], 2021. URL
780 <https://www.mapbox.com/maps/satellite> (accessed 11.23.21).

781 Martinis, S., Twele, A., Plank, S., Zwenzner, H., Danzeglocke, J., Strunz, G., Lüttenberg, H.-P., Dech,
782 S., 2017. The International Charter ‘Space and Major Disasters’: DLR’s Contributions to
783 Emergency Response Worldwide. PFG 85, 317–325. [https://doi.org/10.1007/s41064-017-0032-](https://doi.org/10.1007/s41064-017-0032-1)
784 1

785 Mattyus, G., Wang, S., Fidler, S., Urtasun, R., 2015. Enhancing Road Maps by Parsing Aerial Images
786 Around the World, in: 2015 IEEE International Conference on Computer Vision (ICCV). IEEE,
787 Santiago, pp. 1689–1697. <https://doi.org/10.1109/ICCV.2015.197>

788 McFeeters, S.K., 1996. The use of the Normalized Difference Water Index (NDWI) in the delineation
789 of open water features. International Journal of Remote Sensing 17, 1425–1432.
790 <https://doi.org/10.1080/01431169608948714>

791 Mnih, V., Hinton, G., 2012. Learning to Label Aerial Images from Noisy Data, in: Proceedings of the
792 29 Th International Conference on Machine Learning. Edinburgh, p. 8.

793 Neupane, B., Horanont, T., Aryal, J., 2021. Deep Learning-Based Semantic Segmentation of Urban
794 Features in Satellite Images: A Review and Meta-Analysis. Remote Sensing 13, 808.
795 <https://doi.org/10.3390/rs13040808>

796 Olson, D.M., Dinerstein, E., Wikramanayake, E.D., Burgess, N.D., Powell, G.V.N., Underwood, E.C.,
797 D’amico, J.A., Itoua, I., Strand, H.E., Morrison, J.C., Loucks, C.J., Allnutt, T.F., Ricketts, T.H.,
798 Kura, Y., Lamoreux, J.F., Wettengel, W.W., Hedao, P., Kassem, K.R., 2001. Terrestrial
799 Ecoregions of the World: A New Map of Life on Earth: A new global map of terrestrial
800 ecoregions provides an innovative tool for conserving biodiversity. BioScience 51, 933–938.
801 [https://doi.org/10.1641/0006-3568\(2001\)051\[0933:TEOTWA\]2.0.CO;2](https://doi.org/10.1641/0006-3568(2001)051[0933:TEOTWA]2.0.CO;2)

802 OpenStreetMap [WWW Document], 2021. . OpenStreetMap. URL <https://www.openstreetmap.org/>
803 (accessed 11.23.21).

804 Pan, S.J., Yang, Q., 2010. A Survey on Transfer Learning. IEEE Transactions on Knowledge and Data
805 Engineering 22, 1345–1359. <https://doi.org/10.1109/TKDE.2009.191>

806 Pekel, J.-F., Cottam, A., Gorelick, N., Belward, A.S., 2016. High-resolution mapping of global surface
807 water and its long-term changes. Nature 540, 418–422. <https://doi.org/10.1038/nature20584>

808 Rahnemoonfar, M., Chowdhury, T., Sarkar, A., Varshney, D., Yari, M., Murphy, R.R., 2021.
809 FloodNet: A High Resolution Aerial Imagery Dataset for Post Flood Scene Understanding.
810 IEEE Access 9, 89644–89654. <https://doi.org/10.1109/ACCESS.2021.3090981>

811 Rakhlin, A., Davydow, A., Nikolenko, S., 2018. Land Cover Classification from Satellite Imagery with
812 U-Net and Lovász-Softmax Loss, in: 2018 IEEE/CVF Conference on Computer Vision and
813 Pattern Recognition Workshops (CVPRW). IEEE, Salt Lake City, UT, USA, pp. 257–2574.
814 <https://doi.org/10.1109/CVPRW.2018.00048>

815 Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional Networks for Biomedical Image
816 Segmentation, in: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (Eds.), Medical Image
817 Computing and Computer-Assisted Intervention – MICCAI 2015. Springer International
818 Publishing, Cham, pp. 234–241. https://doi.org/10.1007/978-3-319-24574-4_28

819 Rottensteiner, F., Sohn, G., Gerke, M., Wegner, J.D., 2013. ISPRS Test Project on Urban Classification
820 and 3D Building Reconstruction (No. 12).

821 Rußwurm, M., Wang, S., Korner, M., Lobell, D., 2020. Meta-Learning for Few-Shot Land Cover
822 Classification, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition

823 Workshops (CVPRW). IEEE, Seattle, WA, USA, pp. 788–796.
824 <https://doi.org/10.1109/CVPRW50498.2020.00108>
825 Schreier, G., Dech, S., Diedrich, E., Maass, H., Mikusch, E., 2008. Earth observation data payload
826 ground segments at DLR for GMES. *Acta Astronautica, Touching Humanity - Space for*
827 *Improving Quality of Life. Selected Proceedings of the 58th International Astronautical*
828 *Federation Congress, Hyderabad, India, 24-28 September 2007* 63, 146–155.
829 <https://doi.org/10.1016/j.actaastro.2007.12.010>
830 Simonyan, K., Zisserman, A., 2015. Very Deep Convolutional Networks for Large-Scale Image
831 Recognition. arXiv:1409.1556 [cs].
832 Tan, M., Le, Q.V., 2019. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks.
833 arXiv:1905.11946 [cs, stat].
834 Tasar, O., Tarabalka, Y., Giros, A., Alliez, P., Clerc, S., 2020. StandardGAN: Multi-source Domain
835 Adaptation for Semantic Segmentation of Very High Resolution Satellite Images by Data
836 Standardization, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition
837 Workshops (CVPRW). IEEE, Seattle, WA, USA, pp. 747–756.
838 <https://doi.org/10.1109/CVPRW50498.2020.00104>
839 Twele, A., Cao, W., Plank, S., Martinis, S., 2016. Sentinel-1-based flood mapping: a fully automated
840 processing chain. *International Journal of Remote Sensing* 37, 2990–3004.
841 <https://doi.org/10.1080/01431161.2016.1192304>
842 Voigt, S., Giulio-Tonolo, F., Lyons, J., Kučera, J., Jones, B., Schneiderhan, T., Platzeck, G., Kaku, K.,
843 Hazarika, M.K., Czarán, L., Li, S., Pedersen, W., James, G.K., Proy, C., Muthike, D.M.,
844 Bequignon, J., Guha-Sapir, D., 2016. Global trends in satellite-based emergency mapping.
845 *Science* 353, 247–252. <https://doi.org/10.1126/science.aad8728>
846 Wieland, M., Helleis, M., Fichtner, F., Krullikowski, C., Martinis, S., Plank, S., Motagh, M., in
847 preparation. S1S2-Water: A global dataset for semantic segmentation of water bodies from
848 Sentinel-1 and Sentinel-2 data.
849 Wieland, M., Martinis, S., 2019. A modular processing chain for automated flood monitoring from
850 multi-spectral satellite data. *Remote Sensing* 11, 2330.
851 Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K., 2017. Aggregated Residual Transformations for Deep
852 Neural Networks. <https://doi.org/10.48550/arXiv.1611.05431>
853 Xu, Z., Zhang, W., Zhang, T., Yang, Z., Li, J., 2021. Efficient Transformer for Remote Sensing Image
854 Segmentation. *Remote Sensing* 13, 3585. <https://doi.org/10.3390/rs13183585>
855 Yang, F., Feng, T., Xu, G., Chen, Y., 2020. Applied method for water-body segmentation based on
856 mask R-CNN. *J. Appl. Rem. Sens.* 14, 1. <https://doi.org/10.1117/1.JRS.14.014502>
857 Yuan, W., Wang, J., Xu, W., 2022. Shift Pooling PSPNet: Rethinking PSPNet for Building Extraction
858 in Remote Sensing Images from Entire Local Feature Pooling. *Remote Sensing* 14, 4889.
859 <https://doi.org/10.3390/rs14194889>
860 Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J., 2017. Pyramid Scene Parsing Network. arXiv:1612.01105
861 [cs].
862