

Convolutional Neural Networks for Water segmentation using Sentinel-2 Red, Green, Blue (RGB) composites and derived Spectral Indices

Thomas James¹, Calogero Schillaci^{2*}, Aldo Lipani¹

*corresponding author

¹ Department of Civil, Environmental and Geomatic Engineering, Gower Street, London WC1E 6BT, University College London

² Department of Agricultural and Environmental Sciences, Milan Via Celoria 2, University of Milan

Abstract

Near-real time water segmentation with medium resolution satellite imagery plays a critical role in water management. Automated water segmentation of satellite imagery has traditionally been achieved using spectral indices. Spectral water segmentation is limited by environmental factors and requires human expertise to be applied effectively. In recent years the use of convolutional neural networks (CNN's) for water segmentation has been successful when used on high-resolution satellite imagery, but to a lesser extent for medium resolution imagery. Existing studies have been limited to geographically localised datasets and reported metrics have been benchmarked against a limited range of spectral indices. This study seeks to determine if a single CNN based on Red, Green, Blue (RGB) image classification can effectively segment water on a global scale and outperform traditional spectral methods. Additionally, this study evaluates the extent to which smaller datasets (of very complex pattern e.g harbor megacities) can be used to improve a globally applicable CNN's within a specific region. Multispectral imagery from the European Space Agency, Sentinel-2 satellite (10 m spatial resolution) was sourced. Test sites were selected in Florida, New York, and Shanghai to represent a globally diverse range of water body typologies. Region-specific spectral water segmentation algorithms were developed on each test site, to represent benchmarks of spectral index performance. DeepLabV3-ResNet101 was trained on 33,311 semantically labelled true-colour samples. The resulting model was retrained on three smaller subsets of the data, specific to New York, Shanghai and Florida. CNN predictions reached a maximum mean intersection over union result of 0.986 and F1-Score of 0.983. At the Shanghai test site, the CNN's predictions outperformed the spectral benchmark, primarily due to the CNN's ability to process contextual features at multiple scales. In all test cases, retraining the networks to localised subsets of the dataset improved the localised region's segmentation predictions. The CNN's presented are suitable for cloud-based deployment and could contribute to the wider use of satellite imagery for water management.

1. Introduction

Near-real time mapping of water bodies from satellite imagery plays a critical role in water management. The continuous monitoring of environmental change over time, such as estimation of water availability, prediction of floods, and droughts, is essential to human activities such as agriculture, hydrology, and management (Molden, 2007; Schanze, et al., 2006, Ferral, et al., 2019). As a result, there has been significant interest in identifying methods of accurately automating the water segmentation of satellite imagery.

40 A large body of research has been devoted to the development of Spectral Indices (SIs) to automate water
41 mapping tasks (McFeeters, 1996; Feyisa, et al., 2014; Xu, 2006; Jain, et al., 2020, Zhou et al 2017). SIs are
42 the most prominent tool for automated water mapping and are widely integrated with geospatial software
43 platforms and application programming interfaces (API's). SIs classify each image pixel independently
44 without processing the contextual features of an image. Subsequently, the performance of SIs is hindered
45 by features such as shadow or bright objects such as roofs and solar panels. Additionally, the process of
46 selecting and optimally thresholding a SI is a complex arduous task that must be performed by an
47 experienced professional.

48 In recent years, the expansive growth in the availability and capabilities of graphics processing units
49 (GPU's) has driven the development of sophisticated deep learning (DL) architectures, and more
50 specifically, convolutional neural networks (CNN's). Innovations in CNN architecture has enabled
51 multiscale contextual detection of features within a scene (Chen, et al., 2017). This has led to a surge of
52 interest in state-of-the-art CNN applications to classify land with semantic segmentation (Hoeser and
53 Kuenzer, 2020; Tsagkatakis, et al., 2019). CNN's been hugely successful when used on very high-
54 resolution imagery ($< 1 \text{ m} \times \text{pixel}$), with reported overall accuracy scores that exceed 99% (Talal, et al.,
55 2018; Chen, et al., 2018). CNN's have been less successful on medium resolution imagery, achieving
56 segmentation results ranging from 84% to 97% overall accuracy (Isikdogan, et al., 2017; Wang, et al., 2020;
57 Wieland and Martinis, 2020). Medium resolution imagery contributes to the majority of land mapping
58 activities due to their typically higher spatial and temporal resolution, highlighting a need for further
59 development within this field (Belward, A. and Skoien, J., 2015).

60 Studies tend to be localized to specific geographic regions and have benchmarked CNN predictions against
61 a small group of spectral water segmentation indices, most often Normalized difference water index
62 (NDWI) and Modified Normalized Difference Water Index mNDWI (Isikdogan, et al., 2017; Wang, et al.,
63 2020; Guo, et al., 2020). This study seeks to determine how effective CNN's are in on a global scale, and
64 if CNN's are able to outperform a wide range of spectral methods. Regarding the use of machine learning
65 (ML) for water mapping, Land Remote-Sensing Satellite (Landsat) imagery was used in a ML framework
66 in Nepal (Acharya et al., 2018, 2019) and China (Jiang et al., 2018), where the latter assessed also the
67 performance of the surface water extraction for the entire scene. While the subpixel surface water coverage
68 in urban environments was object of investigation in Sun et al. (2017), and a focus on detection of subpixel-
69 scale inundation was proposed by Jones (2019).

70 Regarding CNN for Remote sensing classification, a recent increase in the output of literature could be seen
71 by a systematic search carried out in Scopus; the query included title abstract and keywords, ("water AND
72 segmentation AND with AND convolutional AND neural AND networks") and it was limited for document
73 type (articles and reviews) and subject area earth and environmental sciences, resulted in 66 research paper
74 between the 2017 and 2020. Out of this papers (Wieland and Martinis, 2020) used CNN and Sentinel-2
75 multispectral imagery to describe a methodology to map large-scale surface water change after drought in
76 Germany. Hughes et al., (2020) used the CNN for classify Synthetic Aperture Radar (SAR) imagery.

77 The first water separation index developed for a multispectral sensor was the (NDWI (McFeeters, 1996).
78 The index was built initially for a Landsat Thematic Mapper (TM) and uses the Near Infrared (NIR) band
79 and Green band to delineate open water features, excluding soil and terrestrial vegetation. There are
80 significant challenges associated with mapping shallow water due to shadow from large physical structures
81 from built-up areas. Xu, (2006) modified the NDWI with mNDWI, replacing the NIR band with short-
82 wavelength infrared (SWIR) band to better partition built-up areas. The resolution performance of mNDWI
83 is limited by the typically lower resolution of the SWIR band. The mNDWI also produces a higher
84 occurrence of false positives in shadow areas, such as cloud shadow, or on dark surfaces such as roads.

85 Feyisa et al. (2014) addressed this shadowing problem with two automated water extraction indices:
 86 $AWEI_{nsh}$ and $AWEI_{sh}$, optimized for environments with no shadow and shadow. The $AWEI_{sh}$ removes
 87 shadow pixels, while $AWEI_{nsh}$ has been designed specifically for urban areas. An alternative method was
 88 proposed by Mishra and Prasad (2015) to improve detail the detection of shallow water. This was achieved
 89 simply through the addition of an index using blue and NIR band. Jain, et al. (2020) built upon I, with PI,
 90 demonstrating a reduction in noise with the SWIR band instead of the NIR band. Errors often occur from
 91 spectral diversity within the water. Turbid water has higher reflectance in the NIR and above bands due to
 92 high concentrations of suspended sediment. This can be corrected by integrating the normalized difference
 93 built-up (NDBI) index (Zha, et al., 2004). False negatives can occur from water bodies that contain high
 94 concentrations of phytoplankton (Chen, et al., 2015). This can be corrected using the normalized difference
 95 vegetation index (NDVI) (Tarpley, et al., 2015). Table 1 summarizes all SIs described in this literature
 96 review.

97 **Table 1:** A summary table of all Spectral Indices related to water segmentation.

Indices	Equation	Merit	Limitation	Reference
NDWI	$(\rho_{Green} - (NIR)) / (\rho_{Green} + (NIR))$	NIR channel has higher resolution capabilities that other sensors.	Less capable of delineating between built-up areas and water.	(McFeeters, 1996)
mNDWI	$(\rho_{Green} - (SWIR_1)) / (\rho_{Green} + (SWIR_1))$	Use of the SWIR band offers greater contrast between built-up areas and water bodies.	The SWIR bands are less capable at higher resolutions. Typically produces false positives on roads, shadows and dark surfaces.	(Xu, 2006)
$AWEI_{nsh}$	$4 (\rho_{Green} - (SWIR_1)) - 0.5 (NIR) + 2.75 (SWIR_2)$	Capable of delineating water and dark surfaces that occur from shadow in built up urban areas	Typically produces false positives on roads, shadows and dark surfaces.	(Feyisa, et al., 2014)
$AWEI_{sh}$	$(\rho_{Blue} + 2.3 \rho_{Green} - 1.5 ((NIR) + (SWIR_2))) / ((\rho_{Green} + (NIR) + (SWIR_1) + (SWIR_2)))$	Removed shadow pixels.	High albedo surfaces such as snow, white roofs and crop-coverings can produce false positives.	(Feyisa, et al., 2014)
I	$(\rho_{Green} - (NIR)) / (\rho_{Green} + (NIR)) + (\rho_{Blue} - (NIR)) / (\rho_{Blue} + (NIR))$	Improves the detail of shallow water detection.	Excess spectral noise.	(Mishra and Prasad, 2015)
PI	$(\rho_{Green} - (SWIR_1)) / (\rho_{Green} + (NIR)) + (\rho_{Blue} - (SWIR_1)) / (\rho_{Blue} + (NIR))$	Noise reduction resulting from SWIR use.	SWIR bands are less capable at higher resolutions.	(Jain, et al., 2020)
NDBI	$((SWIR_1) - (NIR)) / ((SWIR_1) + (NIR))$	Identifying built up areas. Capable of isolating narrow water bodies.	Only applicable in areas of dense vegetation. Noise occurs from any vegetation.	(Zha, et al., 2004)
NDVI	$((NIR) - \rho_{Red}) / ((NIR) + \rho_{Red})$	Can be used for calibrating against high water phytoplankton content.	Water bodies with low reflectance in both red and NIR can produce false positives.	(Tarpley, et al., 2015)

98

99

100 The rationale of this work is to develop a widely usable application Copernicus Sentinel-2 multispectral
 101 and true-color imagery, Red, Green and Blue (RGB) composite, for selected sites will be accessed and
 102 labelled manually, existing spectral algorithms will be fine-tuned to generate benchmarks that represent the
 103 optimal capabilities of spectral water segmentation methods, sites containing complex and diverse
 104 waterbodies were selected. This study also investigates the potential of geographically localizing CNN's
 105 with the use of smaller subsets of data through transfer learning. The results of this study hope to contribute

106 to the development of automated water segmentation tools to streamline access to earth observation
107 analytics.

108 The aim of this work is twofold: i) Determine if water segmentation using CNN's on Sentinel-2 row
109 imagery can outperform multispectral water segmentation indices and, ii) determine if transfer learning
110 with small geographically localized datasets can improve CNN water segmentation's performance in
111 specific regions.

112 2. Material and methods

113 2.1 Data Preparation

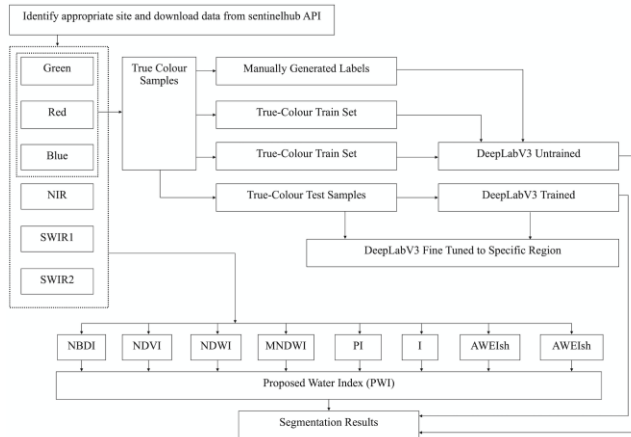
114 Suitable areas of interest were selected using google earth imagery. Satellite data from the sites was
115 downloaded from the Sentinel-hub API and labelled. Test sites from Florida, New York and Shanghai
116 were set aside for testing.

117 As part of the Copernicus programme of the European Commission (EC), the European Space Agency
118 (ESA) has launched the Sentinel-2 constellation (Drusch et al.,2012). The Copernicus programme aims to
119 enable atmospheric, land and marine environment monitoring, climate change research, emergency
120 management, and support security. The constellation consists of two satellites, 2A and 2B (Drusch, et al.,
121 2012). The purpose of the Sentinel-2 mission is to monitor global land surfaces and coastal waters
122 continuously. The Sentinel-2 constellation systematically acquires imagery between -56° to 84° latitude.
123 Sentinel-2 is sun-synchronous at 786 km altitude with $14 + 3/10$ revolutions per day. The Sentinel-2
124 satellites are equipped with filter-based push-broom imager multispectral (MSI) sensors. The bands at 10
125 m resolution are the blue (458 to 523 nm), green (543 to 578 nm), red (650 to 680 nm) and near-infrared
126 (NIR) (785 to 900 nm). There are 6 bands of 20 m spatial resolution, four of which are narrow bands (689
127 to 713 nm, 733 to 748 nm, 773 to 793 nm, 855 to 875 nm), primarily used for vegetation
128 characterisations, and two SWIR-1 (1565 to 1655 nm), SWIR-2 (2100 to 2280 nm) used for detecting
129 clouds, snow and ice and vegetation moisture measurements. There are 3 bands of 60 m spatial
130 resolutions: aerosols (433 to 453 nm), important for analysing the oceanic ecosystem and water vapour
131 (935 to 955 nm) Shortwave infrared for Cirrus detection (1360 to 1390 nm); these bands are used for
132 atmospheric corrections (Gascon, et al., 2017). The Sentinel-2 scenes were accessed using the Sentinel-
133 hub Web Coverage Service (WCS).

134 The Sentinel hub application programming interface (API) was used to source all data used within this
135 study. The API enabled programmatic processing and integration of satellite data into a Python
136 environment. The data is made available through two different levels: Level-1C (L1C) and Level-2A
137 (L2A). L1C corresponds to top-of-atmosphere (TOA) observations, while the L2A is an atmospherically
138 corrected bottom-of-atmosphere (BOA) product. Specific layer configurations were set up to generate the
139 data. Bands 1 to 12 were sourced from both L1C and L2A products and corresponding true-colour
140 composite were accessed for each selected scene.

141 2.1.1 Site Selection and Data Acquisition

142

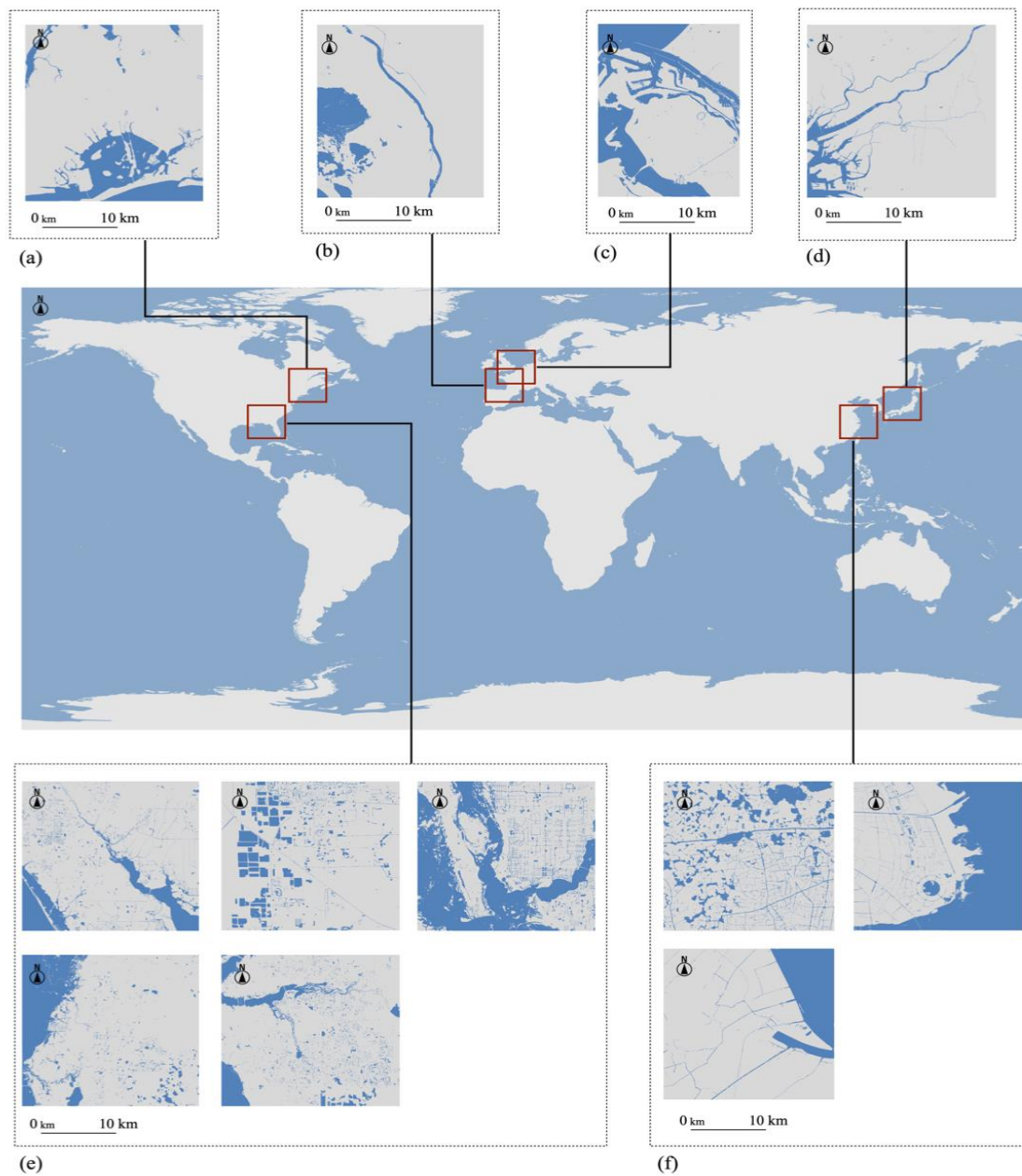


143

144 Figure 1. Flowchart to show the stages of model training and proposed water segmentation index development.

145 The workflow is summarised in a flowchart (Figure 1). Sites were identified with the aid of the google
 146 earth imagery platform. Areas of heavily built-up and complex sea to land interfaces or locations with
 147 densely packed diverse inland water bodies were selected. The training data was selected from the
 148 Netherlands, Osaka, New York, and Florida (Figure 2). The Sentinel-hub EO Browser was used to quick
 149 search for suitable acquisition dates. From selected sample sites and acquisition dates, the 10 m
 150 resolution, 12 multi-spectral layers data and a corresponding true colour reconstruction was requested
 151 from the Sentinel-hub API. Samples downloaded for the labelling purposes were both L1C and L2A
 152 products and were filtered at 0% cloud coverage to prevent any incorrect labelling. The samples intended
 153 for labelling were also selected from summer months, to prevent mislabelling due to periodic snow or ice
 154 within the scene. Once labelled, additional scenes at the same location were downloaded at both the L1C
 155 and the L2A processing levels, with a maximum cloud cover filter of 20%.

156



157
 158 Figure 2. Labels for all samples used in the training process. a) New York, b) Marseille, c) Rotterdam, d) Osaka, e)
 159 Florida region, d) Shanghai and surrounding region.

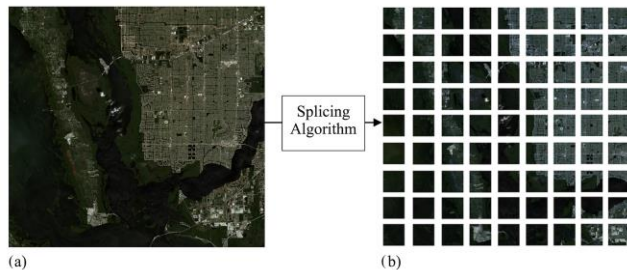
160
 161 Every sample used in the training process and for evaluation purposes was labelled by photointerpretation
 162 and validated using at least 3 images throughout the year. The labelling process was aided by photoshop
 163 tools, primarily the ‘magic wand’ tool. The magic wand tool accepts a colour value of the selected pixels
 164 and expands the selection area to all neighbouring pixels of a similar colour value to build the “Region of

165 Interest” ROI. Large water bodies with homogenous colour could be quickly labelled, however scenes
166 containing variable water texture required a more fragmented and attentive approach. Small localised
167 water bodies required meticulous examination to ensure they were not missed. To reduce the number of
168 falsely identified pixels, the scene was checked against 0.3 m resolution imagery obtained by
169 miscellaneous sources through google earth imagery at a scale of at least 1:1000.

170 Once the sites boundaries were defined and were fully labelled, imagery of the same region were
171 downloaded and matched with the original scene labels. The additional scenes were within the closest
172 possible time periods to the original image to reduce any changes that may have occurred over time.
173 These duplicates were chosen to represent the variability in both atmospheric and surface properties that
174 can be expected within a scene.

175 The CNN used required input channels with dimensions of $3 \times 244 \times 244$, and labels with dimensions $1 \times$
176 244×244 . To preserve the 10 m resolution of the original samples, the images and the labels were spliced
177 into sub-samples with dimensions of 224×244 (Figure 3). The number of spliced subsets equates to the
178 number of samples stated for model training.

179



180

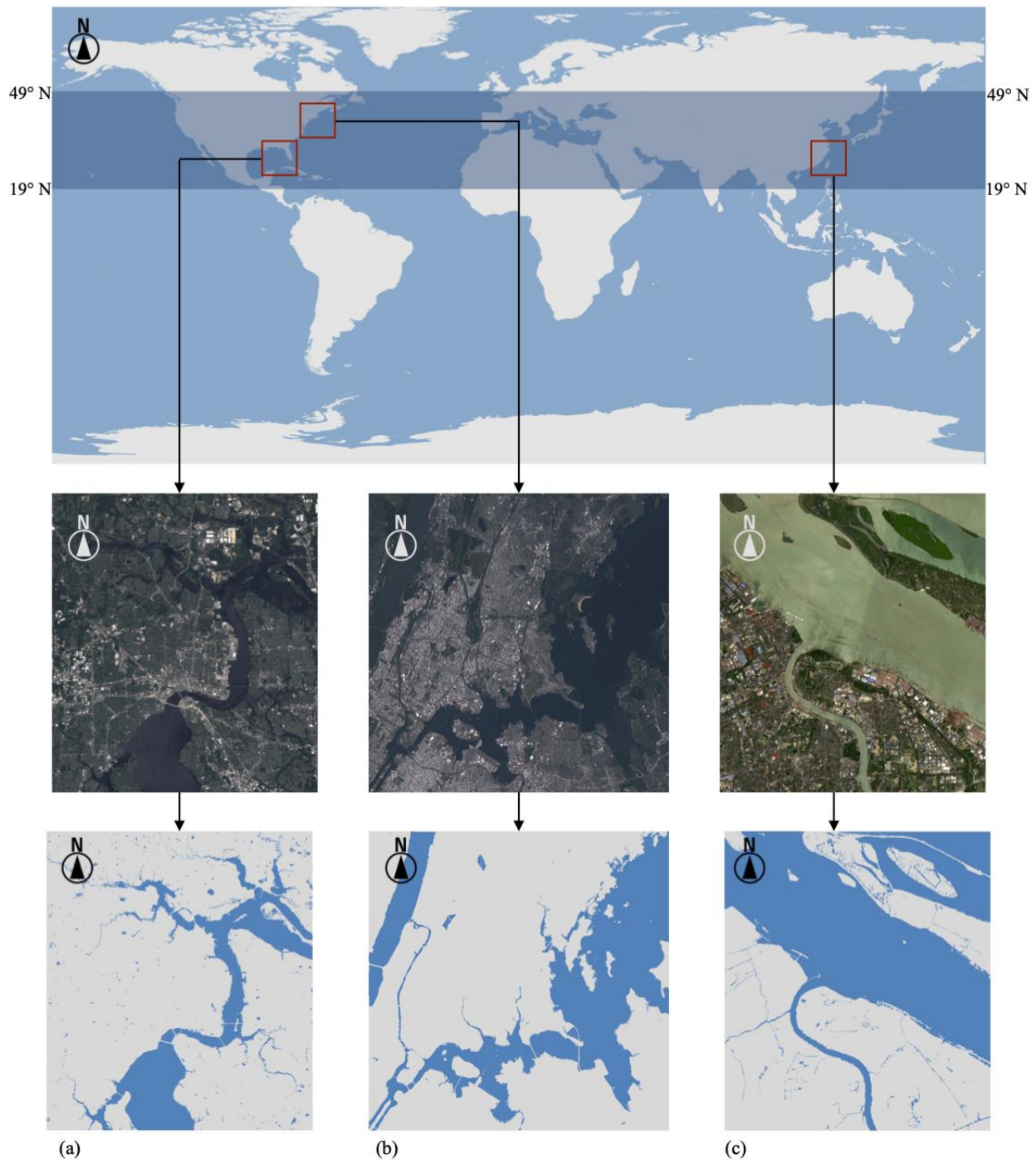
181

182 **Figure 3.** Input and output of splicing algorithm used to generate 244 x 244 pixel samples compatible with
183 the DeepLabV3 model. a) Original true-colour image of Fort Myers, Florida. b) Corresponding true colour
184 images spliced into 244 x 244 samples.

185

186 2.1.2 Test Sites

187 Three areas were delegated and preserved specifically as a benchmark for evaluating water mask
188 predictions. The areas were not exposed to the CNN at any stage of the training process. These three sites
189 have been displayed in Figure 4 and Table 2. All three test sites cover a mixture of heavily urbanised and
190 rural land use. The first evaluation area was a 21.96 km \times 19.52 km region covering Jacksonville Florida.
191 The area was chosen due to the extraordinary density of small lakes within the land, and the complex
192 meandering inland river network. The second area was a 19.52 km \times 19.52 km region covering New
193 York. The area was chosen primarily due to the densely packed tall buildings with extensive shadowed
194 regions. The third area chosen was a 21.96 km \times 19.52 km region of the northern section of Shanghai
195 City, this area enabled the model to be evaluated on a transient intertidal zone with high levels of
196 suspended sediment. Additionally, the area has a very high density of both large and small boats.



198

199 Figure 4. Test sites chosen to test the quality of the model for a) Jacksonville, Florida, b) New York, c) Shanghai.

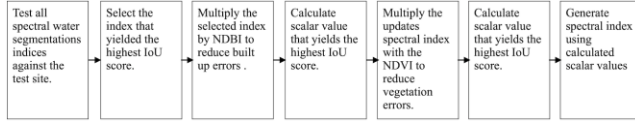
200 Each test site depicted in true colour form and corresponding binary classification label.

201 **2.2. Benchmarks**202 The mIoU values for a range of SIs were calculated. SIs were developed for each test site through
203 parameterisation with NBDI and NVDI indices.

204 2.2.1 Development of Spectral Benchmark

205 To test the performance of the model against benchmark SI, an optimised SI was generated for each
 206 evaluation sample to represent a theoretical best-case performance for what could be achieved through
 207 spectral methods. The process has been summarised in figure 5.

208



209

210 Figure 5. Flowchart to describe the stages of development for a spectral benchmark at each test site.

211 2.2.2 Algorithm tuning

212 A novel method was developed using regression curves was used to ‘fine-tune’ the algorithms. A
 213 preliminary assessment of the available SI showed that the index I, AWEIsh, and NDWI, demonstrated
 214 the highest-ranking performance for Florida, Shanghai and New York respectively. To ‘fine-tune’ this
 215 algorithm, it was multiplied by the NDBI and NDVI at different scalars. The scalars were plotted against
 216 mIoU results in a regression curve. The optimised scalar values were derived from the regression curve to
 217 produce an optimised SI for each test site. The finalised optimised spectral algorithms have been denoted
 218 in equations 1, 2 and 3.

219

$$PWI_{Florida} = -0.4 \frac{(SWIR_2) - (NIR)}{(SWIR_2) + (NIR)} + \frac{\rho_{Green} - (NIR)}{\rho_{Green} + (NIR)} + \frac{\rho_{Blue} - (NIR)}{\rho_{Blue} + (NIR)} + 0.2 \frac{(NIR) - \rho_{Red}}{(NIR) + \rho_{Red}} \quad (1)$$

220

221 **Equation 1:** Proposed spectral index for the Florida test scene: $PWI_{Florida}$

222

223

$$PWI_{Shanghai} = 0.5 \frac{(SWIR_2) - (NIR)}{(SWIR_2) + (NIR)} + 4(\rho_{Green} - (SWIR_1)) - \frac{0.25(NIR) + 2.75(SWIR_2)}{\rho_{Green} + (SWIR_1) + (SWIR_2) + (NIR)} - 3.4 \frac{(NIR) - \rho_{Red}}{(NIR) + \rho_{Red}} \quad (2)$$

224

225 **Equation 2:** Proposed spectral index for the Shanghai test scene: $PWI_{Shanghai}$

226

$$PWI_{New York} = -0.1 \frac{(SWIR_2) - (NIR)}{(SWIR_2) + (NIR)} + \frac{\rho_{Green} - (NIR)}{\rho_{Green} + (NIR)} + \frac{\rho_{Blue} - (NIR)}{\rho_{Blue} + (NIR)} + 0.3 \frac{(NIR) - \rho_{Red}}{(NIR) + \rho_{Red}} \quad (3)$$

227

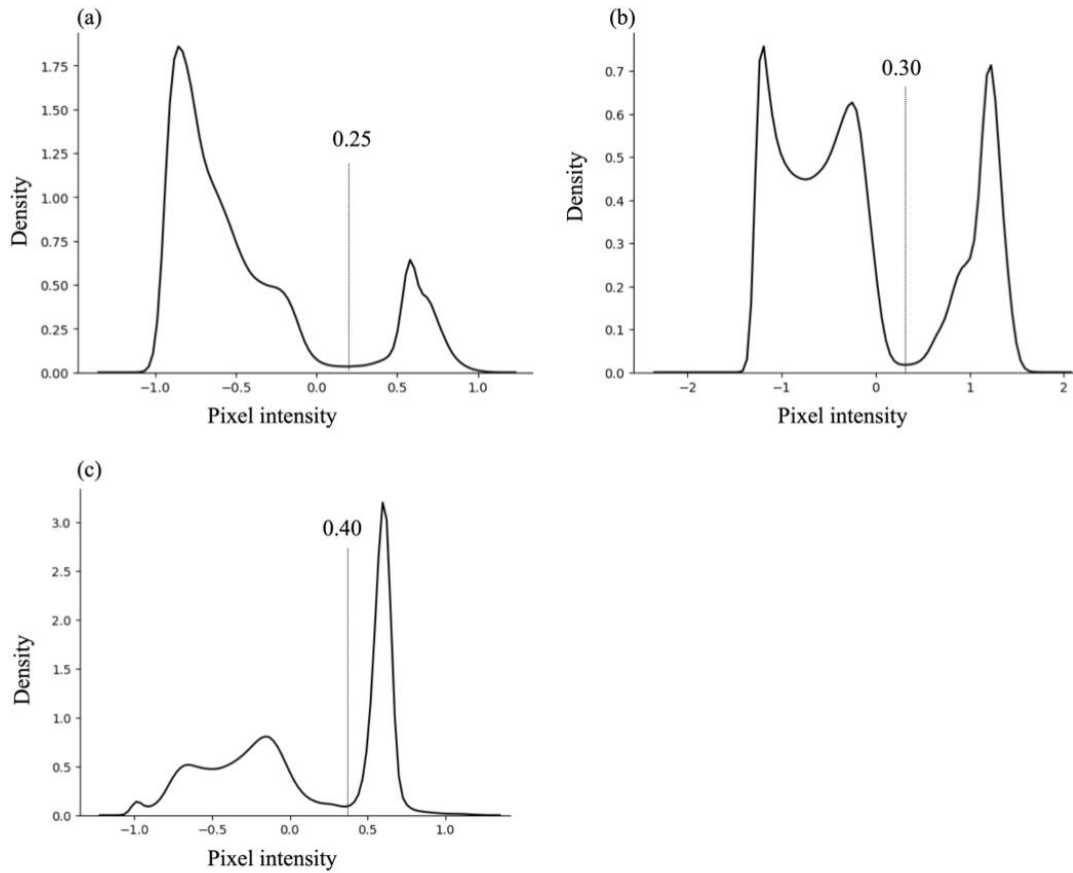
228 **Equation 3:** Proposed spectral index for the New York test scene: $PWI_{New York}$

229

230 2.2.3 Benchmark Algorithm Threshold Optimisation

231 The optimum threshold value for the benchmark water indices was determined by plotting the result of
 232 the probability density function of pixel intensity values of the proposed water index. This enabled visual
 233 inspection of the distribution of pixel intensity. The segmentation threshold was chosen by visually
 234 identifying the lowest point between the two intensity peaks. This was performed iteratively in unison
 235 with the algorithm tuning step (Figure 6).

236



238

239 Figure 6. Distribution plots to show the results of the probability density function of pixel intensity values
 240 of the proposed water index for a) Florida, b) New York, c) Shanghai. Includes indication of the values of
 241 optimal thresholds.

242 **Table 2:** Summary table for test site data

	Jacksonville Florida	New York	Shanghai
Coordinates (WGS84)	-81.761727, 30.241694, -81.507889, 30.454348	-73.957111, 40.717802, -73.703274, 40.930456	120.69626, 30.873884, 120.950097, 31.086538
Sample features	Small densely distributed lakes and urban water bodies. Narrow, convoluted, intertidal rivers. Variable water texture and reflectivity.	Densely packed tall buildings and with extensive shadowing.	Large intertidal zone with complex tributaries and patched of sediment rich water. Variable water texture and reflectivity.
Dimensions	2196 pixels × 1952 pixels 21.96 (km) × 19.52 (km)	1952 pixels × 1952 pixels 19.52 km × 19.52 km	2196 pixels × 1952 pixels 21.96 km × 19.52 km
Product	LIC (BOA)	LIC (BOA)	LIC (BOA)

243

244

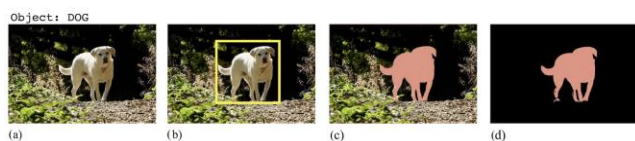
245 2.3. Image Segmentation and CNN

246 2.3.1 Semantic segmentation

247 The semantic segmentation refers to the process of making pixel-wise predictions for a given image
248 (Long, et al., 2015). The potential methods of scene classification have been depicted in figure 7.

249 Semantic segmentation differs from image recognition, object detection and instance segmentation in that
250 every pixel in the image is given a classification, in this case, black pixels represent planet earth, while
251 pink pixels in. For EO classification tasks, semantic segmentation is the classification method of choice,
252 due to their applicability to land surface classification and change detection tasks (Jain, et al., 2020;
253 Hoenser and Kuenzer, 2020).

254



255

256 Figure 7. Depiction of various computer vision classification tasks. a) Object Detection, b) Object Localisation, c)
257 Instance Segmentation, d) Semantic Segmentation.

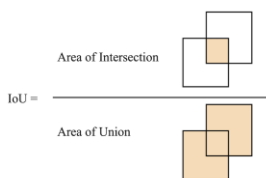
258

259 2.3.2 Evaluating Segmentation

260 The similarity of the segmentation prediction and ‘ground truth’ indicates the quality of the prediction.
261 Many evaluation criteria have been proposed to evaluate the quality of the performance of a given
262 segmentation method. Intersection over Union (IoU) and F1-scores are the two most frequently used
263 metrics of evaluation for computer vision semantic segmentation tasks. All the metrics used within this
264 study have been outlined in table 3.

265 IoU computes a ratio between the intersection and the union of the prediction and the ground truth. This
266 returns a value between 0 and 1. A value of 1 indicates a segmentation result that perfectly matches the
267 ground truth (Figure 8).

268



269

270 Figure 8. Visualised formula for the computation of the intersection over union (IoU) metric.

271 Where more than one class exists, the mean Intersection over Union (mIoU) can be calculated by taking
272 the mean of the IoU values across all the classes (Garcia-Garcia, et al., 2018).

273 The F1 score is a statistical metric for evaluating classification that represents the harmonic mean
274 between the precision and recall. The metric returns a result between 0 and 1, where 1 indicates both the
275 precision and the recall was perfect (Sasaki, 2007).

276 It is important to note that there is no metric that perfectly represents the quality of a semantic prediction.
 277 The IoU penalizes instances of incorrect classification more than the F1-score. Therefore, a lower IoU
 278 score can be expected.

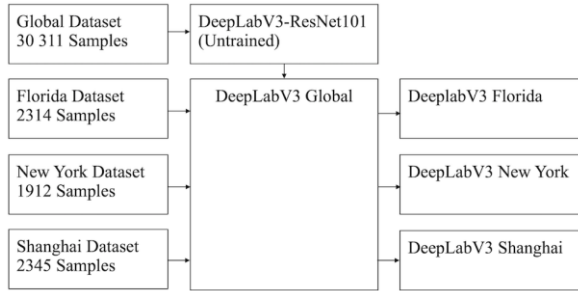
279 **Table 1:** Table to describe metrics used within this thesis.

Metric	Description	Formula
Total number of classes	Number of classes defined, (2 for a binary task).	K
True Positive	Sum of correctly identified water pixels.	TP
True Negative	Sum of correctly identified non-water pixels.	TN
False Positive	Sum of pixels incorrectly identified as water	FP
False Negative	Sum of pixels incorrectly identified was non-water	FN
Precision (P)	The proportion of water detected water pixels	$\frac{TP}{TP + FP}$
Recall (R)	The proportion of ground truth water pixels detected	$\frac{TP}{TP + FN}$
F_1 Score	The harmonic mean between the precision and recall	$2 \times \frac{R \times P}{R + P}$
Intersection over union (IoU)	The ratio between the intersection and the union of the predicted segmentation and the ground truth.	$\frac{TP}{TP + FP + FN}$
Mean Intersection over Union (mIoU)	The mean of the IoU values across all the classes	$\frac{1}{K} \sum_{k=1}^K \text{IoU}_k$

280

281 2.3.3 Choice of Model: DeepLabV3

282 A CNN that is tasked with the segmentation of these water bodies must have the capability to learn
 283 features that are spatially invariant and complex in nature. Based on a review of existing architectures for
 284 semantic segmentation, a naïve decoder architecture was chosen over an encoder-decoder. The results
 285 reported by Guo, et al. (2020) in particular showed that the models that used bilinear upsampling were
 286 better suited to water segmentation tasks. The model chosen for this task was the DeepLabV3-ResNet101
 287 model (Figure 9). This is a state-of-the-art CNN with that is currently ranked the third highest performing
 288 network on the PASCAL-VOC 2012 *test dataset* for semantic segmentation and the second highest
 289 performing naïve decoder (Hoeser and Kuenzer, 2020). The segmentation of water bodies is possible
 290 primarily due to the atrous spatial pyramid pooling section of the network combined with 6.091×10^7
 291 trainable parameters, enabling the CNN to understand features at depth, across multiple scales.



292
293 Figure 9. Model training flow chart.

294

295 2.3.4 Models

296 This study presents four models that were trained and evaluated within this study (figure 7), train:
297 DeepLabV3 was retrained with the 33,331 samples collected in step one.

298 1. DeepLabV3_Global:

299 The DeepLabV3 model was loaded in a ‘untrained’ form. The hyperparameters were adjusted and the model
300 was retrained with all 33,311 training samples. Intended for water segmentation tasks independent of
301 location.

302 2. DeepLabV3_Florida:

303 DeepLabV3_global was retrained with 2314 samples from the state of Florida. intended to complete water
304 segmentation tasks in Florida.

305 3. DeepLabV3_New_York:

306 DeepLabV3_global was retrained with 1912 training samples within New York and Pennsylvania. Intended
307 to perform water segmentation tasks in the New York area.

308 4. DeepLabV3_Shanghai:

309 DeepLabV3_global was retrained with 2345 training samples within Shanghai and small surrounding cities.
310 Intended to perform water segmentation tasks in the New York area.

311 2.3.5 Dataset Manipulation

312

313 The non-test dataset was randomly split so 80% of the training samples were training data, and the
314 remaining 20% of the data was validation data. A train-loss and a validation-loss was computed for each
315 batch. Where the train-loss exceeded the validation loss, the model was considered to be underfitting, and
316 if they validation loss exceeded to the train loss then the model was considered to be overfitting.

317 The dataset was augmented to enhance the size, quality, and diversity of the training data set. This acts as
318 a regularizer to reduce overfitting. Synthetically created duplicates of the training set were created by
319 combinations of horizontally and vertically flipped images before splicing. Further augmentation was
320 done with ‘salt and pepper’ noise and by blurring the samples. By training the network on deliberately
321 noisy data, it was hoped that the model would better generalise when tested on noisy data.

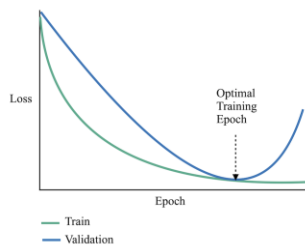
322 The model hyperparameters were ‘fine-tuned’ to find an optimal trade-off between bias and variance. The
323 model training was performed on a smaller subset of the dataset containing 1000 samples. The training
324 loss and validation loss were recorded at each batch and plotted in training logs. The relationship between
325 train-loss and validation loss was used to fine tune the hyperparameters. The final training
326 hyperparameters have been summarised in table 4.

327 Trial runs of model training with learning rates of $\times 10^{-2}$, $\times 10^{-3}$, $\times 10^{-4}$, and $\times 10^{-5}$ were tested. A learning
 328 rate of $\times 10^{-3}$ was chosen for all the training. This learning rate was found to be an optimal trade-off
 329 between large gradient descent step sizes that fail to identify global minima by overshooting, and gradient
 330 descent step sizes cause convergence on local minima and require impractical training time periods.

331 The number of epochs refers to the number of times an algorithm will train through the entire dataset.
 332 Gradient descent is an iterative optimisation algorithm, therefore, requires more than one epoch. Each
 333 epoch is comprised of at least on batch. The size of a batch is determined by how many training samples
 334 are present within the batch and the number of iterations is defined at the number of batches required to
 335 complete one epoch (Smith, et al., 2017; Masters and Luschi, 2018). Figure 10 demonstrates the trajectory
 336 of the train-loss and validation-loss outputs when a model is trained for too many epochs. Each model
 337 was trained for between 60 and 90 epochs. At a specific point during the training process, the validation
 338 loss would start to increase. This was the indicator that the model was overfitting to the training data. This
 339 point was identified, and the model was configured to stop training at the identified epoch (Figure 10).
 340 The loss function used for the model was mean squared error (MSE).

341

342



343

344 Figure 10. Visual representation of the how the optimal training epoch is chosen.

345 2.4. Transfer Learning

346 The resulting model from step three was retrained three times with smaller subsets of data from Florida,
 347 New York and Shanghai. Once the DeepLabV3 Global was successfully trained and evaluated, the model
 348 was loaded with the weights and re-trained with the smaller region-specific datasets. The Florida samples
 349 made up of sites limited to the state of Florida, the New York samples were limited to the sites within the
 350 State of New York and Pennsylvania. Shanghai training samples were limited to Shanghai and
 351 neighbouring cities Suzhou and Nantong.

352 **Table 4:** Summary of Model Hyperparameters, all parameters were set using the validation score.

Model name	Optimiser	Learning rate	Loss function	Number of training sample	Epochs stop	Max epoch	Training-validation split	Batch size
DeepLabV3 Global	Adam	1×10^{-3}	MSE	33331	10	60	80/20	20
DeepLabV3 Florida	Adam	1×10^{-3}	MSE	2314	25	90	80/20	20
DeepLabV3 New York	Adam	1×10^{-3}	MSE	1912	25	90	80/20	20
DeepLabV3 Shanghai	Adam	1×10^{-3}	MSE	2345	25	90	80/20	20

353

354

355

356 2.5. Model Evaluations

357 The trained models were used to generate water mask predictions for each test site. The predictions were
358 compared to the ground truths and the water mask predictions made using the SIs developed in step two.
359 The predictions were quantitatively and qualitatively analysed.

360 The evaluation of a semantic segmentation output is conventionally done using metrics. However, the
361 evaluation can benefit from a parallel qualitative analysis to visually identify the relationships and
362 patterns that may exist.

363 To evaluate quantitatively the CNN's performance against spectral water segmentation methods, the F1-
364 Scores and mIoU results were computed for all the CNN and SI predictions on all three test sites. Each
365 metric was calculated by comparing the prediction to the manually generated ground truths. Comparisons
366 are made between the CNN's and the SI benchmark.

367 To determine whether transfer learning improved the results with respect to a specific region, the
368 performance of DeepLabV3_Florida, DeepLabV3_New_York and DeepLabV3_Shanghai was compared
369 to DeepLabV3_Global for each test site.

370 Qualitative observations were made about the overall prediction quality and how the CNN's responded to
371 contextually dependent features. Special attention was focused on transient features such as boats and
372 intertidal zones or wetlands. This was achieved by visually identifying specific sources of false positives
373 and false negatives for all the segmentation methods. The characteristics that were exclusive to each
374 method or common to both methods were noted.

375 It is important to note the following assumptions made within this study:

- 376 1. The evaluations made within this study are made on the assumption that the test sites were labelled within an
377 error margin of 10 m (1 pixel). The use of high-resolution third-party validation data combined with a manual
378 pixel-wise classification enabling accuracy to be maximised, however there is no existing benchmark to
379 validate the accuracy of the test site labels.
- 380 2. The additional scenes used to expand the data, was assumed to have equal water body limits as the original
381 labelled sample.
- 382 3. The spectral benchmark algorithms presented within this study are assumed to be the best possible
383 representation of what can be achieved using SI.

384 3. Results

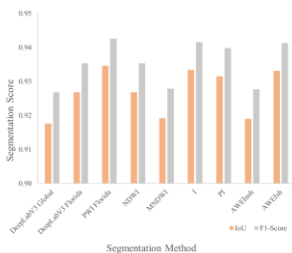
385

386 3.1 Test Site: Jacksonville, Florida

387 The results of the segmentation predictions for DeepLabV3_Global, DeepLabV3_Florida and, $PWI_{Florida}$
388 have been displayed in Figure 11 and table 5. All SIs performed better than the CNN predictions (Figure
389 12). The mIoU results for DeepLabV3_Global, DeepLabV3_Florida and the $PWI_{Florida}$ were 0.913, 0.918,
390 0.93, respectively with F1-Scores of 0.923, 0.927, 0.943, respectively. Retraining DeepLabV3_Global on
391 the 2314 sample Florida dataset increased the mIoU result and F1-Score by 0.004 respectively.

392 Predictions made on the Jacksonville Florida test site yielded the lowest segmentation performance of all
 393 the test sites for all segmentation methods.

394



395

396 Figure 11. Bar chart to compare the F1-Score and the IoU scores for all water segmentation methods for test site
 397 Jacksonville Florida.

398 **Table 5:** Florida Segmentation Results, the best performing

	mIoU	F₁ Score
DeepLabV3_Global	0.91759	0.9268
DeepLabV3_Florida	0.9268	0.9353
PWI_{Florida}	0.9346	0.9426

399

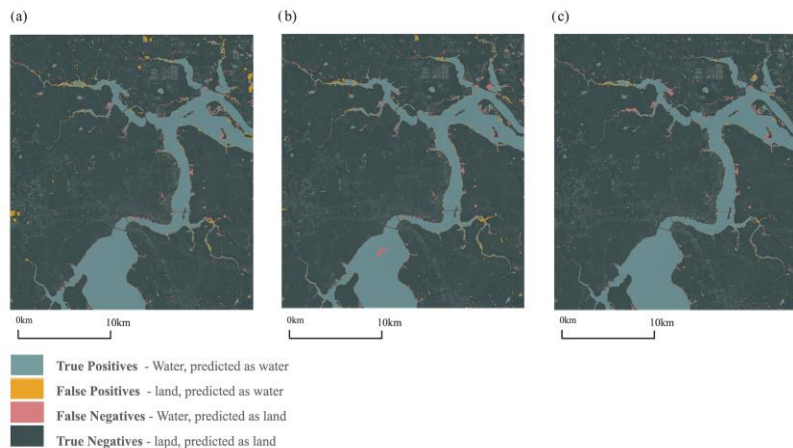
400 Both DeepLabV3_Global and DeepLabV3_Florida identified all the boats (supplemental materials figure
 401 1) within the scene as water, however both CNN's were unable to segment the protruding structures such
 402 as pontoons, jetties and harbours. The boundaries were spatially consistent but appeared to be generalised.

403 The PWI_{Florida} predictions were sharper around the land-sea interface. However, large sediment-rich water
 404 in bays, inlets, and rivers as land was segmented as land.

405 In the case of DeepLabV3_Global small sections of false positives occurred in areas of forest vegetation.
 406 This error was reduced by retraining on the local dataset, and therefore not present in the
 407 DeepLabV3_Florida prediction, however sections of false negative predictions occurred larger water
 408 bodies where they had not occurred in the DeepLabV3_Global predictions. For example, a section of
 409 water at the mouth of the estuary, classified as land by DeepLabV3_Florida.

410 A common characteristic of all the predictions was the misidentification of riverine wetlands and failure
 411 to identify narrow (< 10 m) rivers and small (< 30 m) water bodies.

412

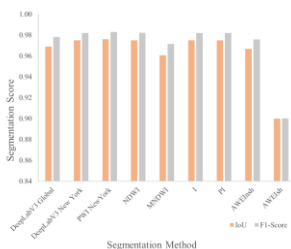


413
414
415
416
417

Figure 12. Plots to compare segmentation results of the DeepLabV3 Global Model trained on all 33,311 training samples, the DeepLabV3_Florida model retrained on imagery from the surrounding area and the PWIFlorida spectral index, fine-tuned specifically for the test site. The mask has been overlaid on the original true colour sample at 83% opacity.

418 3.2 Test Site: New York

419 The results of the segmentation predictions for DeepLabV3_Global, DeepLabV3_New_York and the
420 PWIFlorida have been displayed in Figure 13 and table 6. The CNN's were outperformed by all SIs with the
421 exception of the index I and AWEIsh (Figure 14). The mIoU results for DeepLabV3_Global,
422 DeepLabV3_New_York and the PWI_{NewYork} predictions were 0.969, 0.975, 0.976 respectively with F1-
423 Scores of 0.978, 0.982, 0.983 respectively. The quality of the water segmentations in the New York test
424 site were best of all the test sites, for all segmentation methods.



425
426
427

Figure 13. Bar chart to compare the F1-Score and the IoU scores for all water segmentation methods for test site New York.

428 Re-training of DeepLabV3_Global to samples local to New York increased the mIoU result and F1-Score
429 by 0.006, 0.005 and respectively, this was the largest improvement of all the test sites.

430 **Table 6:** New York: Segmentation Results

	mIoU	F₁ Score
DeepLabV3_Global	0.969	0.978
DeepLabV3_New_York	0.975	0.982
PWI_{NewYork}	0.976	0.983

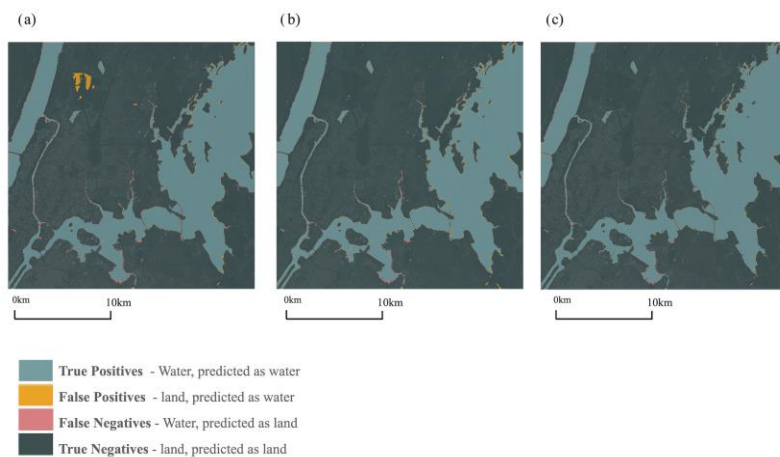
431

432 DeepLabV3_Global predictions in New York mirrored those of the Florida test site with very large false
433 positive occurring over the forest areas, particularly in Van Cortlandt Park. This was reduced in the
434 DeepLabV3_New_York predictions by retraining on localised imagery.

435 Both DeepLabV3_Global and DeepLabV3_New_York were unable to accurately label sections of river
436 that were less than 30 m wide. The CNN’s generalised across complex sections of the land-water
437 interface, which mirrored the Florida test site results.

438 The DeepLabV3_New_York prediction demonstrated a reduced ability to detecting bridges and features
439 compared to the DeepLabV3_Global prediction.

440

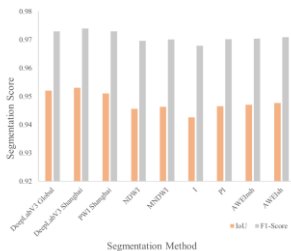


441

442 Figure 14. Plots to compare segmentation results of the DeepLabV3 Global Model trained on all 33,311 training
443 samples, the DeepLabV3_New York model retrained on imagery from the surrounding area and the PWINewYork
444 spectral index, fine-tuned specifically for the test site. The mask has been overlaid on the original true colour sample
445 at 83% opacity.

446 3.3 Test Site: Shanghai

447 The results of the segmentation predictions for DeepLabV3_Global, DeepLabV3_Shanghai and the
448 PWI_{Shanghai} have been displayed in Figure 15 and table 7. Both CNN predictions outperformed all SIs
449 (Figure 15). The mIoU results for DeepLabV3_Global, DeepLabV3_Shanghai and the PWI_{Shanghai}
450 predictions were 0.952, 0.953, 0.951 respectively with F1-Scores of 0.973, 0.974, 0.973 respectively.
451 The transfer learning process increased the mIoU result and F1-Score by 0.001 respectively, indicating
452 that the transfer learning process was the least effective in Shanghai than at any of the three test sites.



453

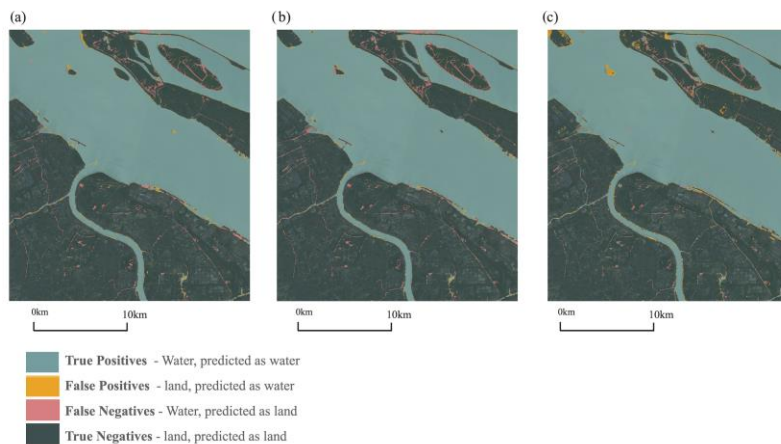
454 Figure 15. Bar chart to compare the F1-Score and the IoU scores for all water segmentation methods for test site
 455 Shanghai.

456 **Table 7:** Segmentation Results: Shanghai

	mIoU	F_1 Score
DeepLabV3_Global	0.952	0.973
DeepLabV3_Shanghai	0.953	0.974
PWIShanghai	0.951	0.973

462 PWIShanghai predictions produced false positives at the
 463 locations of white roofs and solar panels. False negatives in areas of sediment rich water, particularly
 464 around features that are intertidally submersed. The PWIShanghai algorithm was unable to predict at about
 465 5.6×10^5 m² intertidal island in the middle of the coastal estuarine system, this island was however
 identified by DeepLabV3_Global. A zoomed in perspective of this has been displayed in Figure 16.

466 The contextually dependent features like the boats and the residual turbid water of the boats was mapped
 467 as land by the PWIShanghai index (supplemental materials figure 1). The CNN’s demonstrated a
 468 capability to identify the boats, turbid water and map them as water. This has been demonstrated in closer
 469 detail in Figure 16.



470
 471 Figure 16. Plots to compare segmentation results of the DeepLabV3 Global Model trained on all 33,311 training
 472 samples, the DeepLabV3_Shanghai model retrained on imagery from the surrounding area and the PWIShanghai
 473 spectral index, fine-tuned specifically for the test site. The mask has been overlaid on the original true colour sample
 474 at 83% opacity.

475 A common characteristic of all segmentation predictions was the inability to identifying the narrow,
 476 sediment rich tributaries within the intertidal zone and the piers protruding into the estuary. This result
 477 was common to the CNN predictions of all test sites, Table 8.

478 **Table 8:** Results for all water segmentation predictions at all test sites. Numbers in bold are the highest scoring results.

Florida		New York		Shanghai	
mIoU	F_1 Score	mIoU	F_1 Score	mIoU	F_1 Score

DeepLabV3 Global	0.91328	0.9228	0.96106	0.9719	0.95165	0.9731
DeepLabV3 Fine Tuned	0.91759	0.9268	0.96873	0.9775	0.95303	0.9738
PW _I _{LOCATION}	0.9346	0.9426	0.9759	0.9827	0.9512	0.9731
NWDI	0.9268	0.9353	0.975	0.9821	0.9456	0.9696
MNDWI	0.9192	0.9279	0.9605	0.9714	0.9463	0.9701
<i>I</i>	0.9334	0.9415	0.975	0.982	0.9426	0.9679
PI	0.9315	0.9398	0.975	0.982	0.9465	0.9702
AWEInsh	0.919	0.9277	0.9667	0.9758	0.947	0.9704
AWEIsh	0.9331	0.9413	0.8998	0.8998	0.9477	0.9709

479

480 Interestingly, despite an improved mIoU and F1-score, the DeepLabV3_Shanghai was less able to
481 accurately identify islands in the middle of the estuary than DeepLabV3_Global (supplemental materials
482 figure 2).

483 4. Discussions

484 From the obtained results, it was observed that CNN's are capable of outperforming SIs for water
485 segmentation tasks. The CNN's demonstrated an ability to identify contextual features such as boats,
486 turbid water and sediment rich intertidal water bodies. It was shown in all test cases that re-training the
487 neural network to localised datasets improved prediction accuracy. This section explains these results and
488 the associated successes and limitations. The results were placed within the context of existing literature
489 with additional recommendations for further developments. The potential impact of these results on the
490 field of earth observation will be discussed.

491 4.1 DeepLab_Global

492 The Shanghai test site results showed that the CNN's were capable of outperforming all available SIs.
493 This was driven primarily by the intrinsic failures of spectral methods and CNN's ability to process
494 context at multiple scales of an image.

495 The Shanghai test site was characterised by a large intertidal zone, sediment-rich water, and a high marine
496 traffic volume. Suspended sediment within water bodies increases reflectance in NIR and SWIR radiation
497 (Pham, et al., 2018). The gradient of reflectance between the VL and NIR and SWIR wavelengths was
498 reduced causing large misclassification errors in the PW_I_{Shanghai} predictions in the East China Sea
499 intertidal zone. The solar panels and white roofs were incorrectly mapped as water, this is likely attributed
500 to the increased reflectance of VL, which in turn increases the gradient between VL and NIR/SWIR.

501 SI classify pixels on an individual basis without considering the context of the image. This accounted for
502 the high quality prediction observed on the New York test site where the majority of water bodies are
503 deep and there is a clear water-to-land interface due to the relatively small 0.5 m tidal range (Bowman,
504 1976). However, transient features such boats and turbid water are classified as non-water bodies. This an
505 intrinsic error that could not be resolved through spectral methods.

506 The CNN approach is very different. The CNN's learned combinations of characteristics that make up a

507 water body. These include edges, shapes, colour gradients and textural features (Zeiler and Fergus, 2014).
508 The DeepLabV3 network used ASPP to examine convolutional feature layers with filters at multiple
509 sample rates and fields of view (Chen, et al., 2017). This enabled the network to capture the various
510 spatial contexts associated with water detection. The CNN was trained with 33,311 samples, this was a
511 sufficient volume to develop a deep and rich contextual understanding of combinations of characteristics
512 to make accurate prediction of sediment rich water bodies, boats and turbid water. The ability of the CNN
513 to distinguish contextual features was the main driver of success when evaluated against a spectral
514 methods of water segmentation at the Shanghai test site.

515 The CNN's were unable to accurately classify narrow meandering inland rivers, smaller water bodies (<
516 3 m) and complex structures protruding into the water. This was most observable in the Florida test case.
517 The CNN's had a tendency to generalise across complex features, decreasing the prediction quality. This
518 accounted for the poor overall segmentation predictions for the Florida test site where an extensive and
519 complex land-sea interface exists. This was only partially reduced by retraining the model on region
520 specific subsets of data, implying that the detection capabilities were partially limited by the CNN
521 architecture. The 'black box' nature of deep neural networks makes drawing comparisons between results
522 and network architecture difficult and speculative. However, it is clear that the 'smoothing' effect of the
523 DeepLabV3 model precludes the model from achieving the same level of pixel-precision that is present in
524 a SI segmentation. It is important to note that the DeepLabV3 model was built for computer vision tasks
525 from terrestrial, close-range, side-view perspectives. The overhead perspective of EO imagery results in
526 clustering and random distribution of features across which is very different from typical computer vision
527 images.

528 4.2 Transfer Learning Performance

529 During the initial training stages, the DeepLabV3_Global model was shown a globally diverse range of
530 water body typologies. The characteristics of these water bodies are heavily influenced by interdependent
531 variables such as local geomorphology, weather patterns and human activities. These variables are often
532 homogenous to a region. As an example, Florida is characterised by a porous plateau of karst limestone
533 that allows water to move freely forming large wetlands and an extraordinary number of small lakes
534 (Beck, 1986). The DeepLabV3_Global model learned the features to predict water bodies at all three sites
535 based upon the generic characteristics of water bodies. Re-training the network with a small number of
536 local samples reinforced the correct predications made by the DeepLabV3_Global model, while
537 extracting characteristics that are specific to the local region and transferred the knowledge into the new
538 network. This was particularly successful when applied to the New York test site; Large areas of forest
539 were predicted as water by the DeepLabV3_Global model, but resolved in the DeepLabV3_New_York
540 model prediction. It could be speculated that the DeepLabV3_Global model had fitted to the green texture
541 rich water bodies in the Florida dataset and when knowledge was extracted and transferred to the
542 DeepLabV3_New_York network, the error was eliminated.

543 In all test cases, the retraining of the networks resulted in some new errors that did not occur in the
544 DeepLabV3_Global predictions. The most notable, unexpected error was the patch of water identified as
545 land in the Florida test case. Deep learning models are known to be robust to label noise that is evenly
546 distributed across a large dataset, yet highly sensitive to label noise that is concentrated within the dataset
547 (Karimi, et al., 2020). Errors most likely arose from concentrated label noise within the smaller subsets of
548 data. This noise would also be amplified in the augmentation process and transferred to the retrained
549 network.

550 4.3 Comparisons with other CNN performances.

551 The Sentinel-2 mission has been in operation since 2015, which is still a relatively short time frame. As a
552 result of this, the majority of studies covering water segmentation utilised different data sources. It is
553 difficult to compare the performance of CNN's across different image resolutions. The most recent and
554 closest matching study was the segmentation of water bodies within Sentinel-2 imagery exclusively in
555 Germany by Wieland and Martinis, (2020). The results of the current study marginally improved upon
556 this with more diverse and challenging urban test sites. The improved results could be attributed to the use
557 of a more modern and sophisticated CNN that utilises ASPP. The use of complex and contextually rich
558 training samples could also contribute to the small improvement in segmentation accuracy.

559 The findings of this study support the growing consensus that CNN's are becoming more capable than
560 traditional SIs for land classification tasks. It is widely acknowledged that Deep Learning will be
561 instrumental to sustainability and automation in the future (Gulati and Sharma, 2020). It can be expected
562 that the development of network architectures will continue to improve and the subsequently, the quality
563 segmentation tasks will follow.

564 4.4 Implications for the field

565 At a global scale Pekel et al. (2016) showed how long-term changes of water coverage are difficult to
566 map and represent a societal challenge due to the documented reduction of inland water occurred in the
567 Middle East, Central Asia, Australia and the USA. This is linked to drought and anthropic factors (Pekel
568 et al., 2016) and reference therein. The different algorithms reported in the state-of-the-art were based on
569 ML (Acharya et al., 2019), noise suppression methods used in order to mitigate the effect of landforms
570 shadows and solid water forms (Jiang et al., 2018, 2020), and model fusion (Wagle et al., 2020).

571 For a long time, satellite imagery has been expensive and difficult to access for both individuals and
572 organisations (Turner, et al., 2015). The barrier to entry has dropped significantly in recent years with the
573 introduction high performance computing systems and large scale cloud- based computing frameworks,
574 most notably 'Google Earth Engine' (GEE) (Gorelick, et al., 2017). However, to achieve reliable, high
575 quality water mapping with SIs, expertise is required to select and optimise a SI. The CNN's developed
576 within this study are easily deployable to cloud-based platforms. Very little skill is required to use a CNN
577 within a platform like GEE. This could help broaden scope of the possibilities available to individuals and
578 organisations who wish to use satellite imagery for water management.

579 SIs generally require EM radiation in the VL and the NIR and SWIR range. This adds a computational
580 cost for image processing chains and a dependence on satellite sensors' multispectral capabilities. This
581 study shows that state-of-the-art CNN's capabilities match and outperform SIs, potentially precluding the
582 need for NIR and SWIR channel for water segmentation tasks. Alongside a large body of parallel
583 research, this study could contribute to the development of streamlined satellite processing chains

584 5. Conclusion

585 Better results could be achieved through a redesign of the CNN architecture to better suit EO imagery.
586 This could involve adjusting the dilation rates of the atrous convolution kernels to better suit the clustered
587 nature of the water bodies. Alternatively, the use of an encoder-decoder network like DeepLabV3+ has
588 the potential to improve segmentation performance. The incorporation of additional skip connections
589 from the entry and middle blocks of the DeepLabV3+ encoder has been shown to sharpen segmentation
590 outputs (Prabha, et al., 2020). Experimenting with this technique could make it possible to detect and
591 localise very small, narrow and complex water bodies. Some recent studies have swapped RGB input

592 channels for alternatives (Jain, et al., 2020). The performance of water segmentation with CNNs could be
593 improved by replacing the RGB channels with band ratios or outputs of an existing spectral water index.

594 A further enhancement of the transfer learning aspect of this study could involve retraining
595 DeepLabV3_Global to identify specific water typologies rather than geographic locations. For example,
596 re-training DeepLabV3_Global on images collected in areas of karst limestone, instead of samples limited
597 to Florida. CNN's trained to capture the characteristics of specific typologies would enable broader usage
598 than a CNN retrained specifically to geographic location.

599 This study has shown that CNN's are an effective tool for the segmentation of water bodies in medium
600 resolution satellite imagery. This was done by training the DeepLabV3-ResNet101 network with
601 manually labelled Sentinel-2 imagery.

602 Three main conclusions can be made based upon this research:

- 603 i) CNN's can be applied to medium resolution true-colour satellite imagery to effectively map water
604 bodies on a global scale.
- 605 ii) Water segmentation using CNN's on medium resolution true colour satellite imagery can
606 outperform multispectral water segmentation indices.
- 607 iii) Transfer learning with small geographically localised datasets can improve the performance of
608 CNN water segmentation in specific geographic regions.

609 Further developments of the study could include adjusting the network to improve segmentation
610 sharpness and feature localization in EO imagery. Results could be improved by replacing the RGB input
611 channels with alternatives such as band ratios or SI outputs. Additionally, the model presented within this
612 study could be 'fine-tuned' for specific water body typologies.

613 The results of this study could help broaden and streamline the use of EO imagery for water management
614 by improving the efficiency of EO processing chains and lowering the skill barrier.

615 6. Acknowledgment

616 We thank the Academic Editor, Professor Jida Wang and the anonymous reviewers for the helpful
617 comments on the drafts of the manuscript.

618 7. References

619 Acharya, T. D., Subedi, A. and Lee, D. H.: Evaluation of Machine Learning Algorithms for Surface Water
620 Extraction in a Landsat 8 Scene of Nepal, *Sensors* 2019, Vol. 19, Page 2769, 19(12), 2769,
621 doi:10.3390/S19122769, 2019.

622 Beck, B., 1986. A generalized genetic framework for the development of sinkholes and Karst in Florida, U.S.A..
623 *Environmental Geology and Water Sciences*, 8(1), pp. 5-18.

624 Belward, A. and Skoien, J., 2015. Who launched what, when and why; trends in global land-cover observation
625 capacity from civilian earth observation satellites.. *ISPRS Journal of Photogrammetry and Remote Sensing*, Issue
626 103, pp. 115-128.

627 Bowman, M., 1976. The tides of the East River, New York. *Journal of Geophysical Research*, 81(9), pp. 1610-
628 1616.

- 629 Chen, L. C. et al., 2018. Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation.
630 *Computer Vision – ECCV 2018*, pp. 833-851.
- 631 Chen, L.-C. et al., 2016. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous
632 Convolution, and Fully Connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
633 40(4), pp. 834-848.
- 634 Chen, L.-C., Papandreou, G., Schroff, F. & Adam, H., 2017. Rethinking Atrous Convolution for Semantic Image
635 Segmentation.. *arxiv.org*..
- 636 Chen, Y. et al., 2018. Extraction of Urban Water Bodies from High-Resolution Remote-Sensing Imagery Using
637 Deep Learning. *Water*, 10(5), p. 585.
- 638 Drusch, Del Bello U., Carlier S., Colin O., Fernandez V., Gascon F., Hoersch B., Isola C., Laberinti P.,
639 Martimort P., Meygret A., Spoto F., Sy O., Marchese F., Bargellini P., 2012 Sentinel-2: ESA's Optical High-
640 Resolution Mission for GMES Operational Services. Remote Sensing of Environment, *Remote Sensing of*
641 *Environment*, Volume 120, pp. 25-36.
- 642 Ferral, A., Luccini, E., Aleksinkó, A. and Scavuzzo, C. M.: Flooded-area satellite monitoring within a Ramsar
643 wetland Nature Reserve in Argentina, *Remote Sens. Appl. Soc. Environ.*, 15, 100230,
644 doi:10.1016/j.rsase.2019.04.003, 2019.
- 645 Feyisa, G., Meilby, H., Fensholt, R. & Proud, S., 2014. Automated Water Extraction Index: A new technique for
646 surface water mapping using Landsat imagery. *Remote Sensing of Environment*, Volume 140, pp. 23-35.
- 647 Gascon, F.; Bouzinac, C.; Thépaut, O.; Jung, M.; Francesconi, B.; Louis, J.; Lonjou, V.; Lafrance, B.; Massera,
648 S.; Gaudel-Vacaresse, A.; Languille, F.; Alhammoud, B.; Viallefont, F.; Pflug, B.; Bieniarz, J.; Clerc, S.; Pessiot,
649 L.; Trémas, T.; Cadau, E.; De Bonis, R.; Isola, C.; Martimort, P.; Fernandez, V. 2017. Copernicus Sentinel-2A
650 Calibration and Products Validation Status. *Remote Sensing*, 9(6), p. 584.
- 651 Gorelick, N. Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., Moore R., 2017. Google Earth Engine:
652 Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment*, Volume 202, pp. 18-27.
- 653 Gulati, S. and Sharma, S., 2020. Challenges and Responses Towards Sustainable Future Through Machine
654 Learning and Deep Learning. *Data Visualization and Knowledge Engineering*, Volume 32, pp. 151-169.
- 655 Guo, H.; He, G.; Jiang, W.; Yin, R.; Yan, L.; Leng, W. 2020. A Multi-Scale Water Extraction Convolutional
656 Neural Network (MWEN) Method for GaoFen-1 Remote Sensing Images *ISPRS International Journal of Geo-*
657 *Information*, 9(4).
- 658 Hoeser, T. and Kuenzer, C., 2020. Object Detection and Image Segmentation with Deep Learning on Earth
659 Observation Data: A Review-Part I: Evolution and Recent Trends.. *Remote Sensing*, 12(10), p. 1667.
- 660 Hughes, L.H., Marcos, D., Lobry, S., Tuia, D., Schmitt, M., 2020. A deep learning framework for
661 matching of SAR and optical imagery. *ISPRS J. Photogramm. Remote Sens.* 169, 166–179.
662 <https://doi.org/10.1016/j.isprsjprs.2020.09.012>
- 663 Isikdogan, F., Bovik, A., Passalacqua, P., 2017. Surface Water Mapping by Deep Learning. *IEEE Journal of*
664 *Selected Topics in Applied Earth Observations and Remote Sensing*, 10(11).

665 Jain, P., Schoen-Pelan, B., Ross, R., 2020. Automatic flood detection in Sentinel-2 images using deep
666 convolutional neural networks. *SAC '20: Proceedings of the 35th Annual ACM Symposium on Applied*
667 *Computing*, p. 617–623.

668 Jiang, W., He, G., Long, T., Ni, Y., Liu, H., Peng, Y., Lv, K. and Wang, G.: Multilayer Perceptron Neural Network
669 for Surface Water Extraction in Landsat 8 OLI Satellite Images, *Remote Sens.*, 10(5), 755,
670 doi:10.3390/rs10050755, 2018.

671 Jiang, W., He, G., Pang, Z., Guo, H., Long, T. and Ni, Y.: Surface water map of China for 2015 (SWMC-2015)
672 derived from Landsat 8 satellite imagery, *Remote Sens. Lett.*, 11(3), 265–273,
673 doi:10.1080/2150704X.2019.1708501, 2020.

674 Jones, J.: Improved Automated Detection of Subpixel-Scale Inundation—Revised Dynamic Surface Water Extent
675 (DSWE) Partial Surface Water Tests, *Remote Sens.*, 11(4), 374, doi:10.3390/rs11040374, 2019.

676 Karimi, D., Dou, H., Warfield, S., Gholipour, A., 2020. Deep learning with noisy labels: exploring techniques
677 and remedies in medical image analysis. *Medical Image Analysis*.

678 Long, J., Shelhamer, E., Darrell, T., 2017. Fully Convolutional Networks for Semantic Segmentation. *IEEE*
679 *transactions on pattern analysis and machine intelligence*, 39(4), p. 640–651..

680 Masters, D. and Luschi, C., 2018. Revisiting Small Batch Training for Deep Neural Networks. *arXiv:1804.07612*
681 *[cs, stat]*.

682 McFeeters, S., 1996. The use of the Normalized Difference Water Index (NDWI) in the delineation of open water
683 features.. *International Journal of Remote Sensing*, 17(7), pp. 1-9.

684 Mishra, K., Prasad, P., 2015. Automatic Extraction of Water Bodies from Landsat Imagery Using Perceptron
685 Model. *Journal of Computational Environmental Sciences*, p. 9.

686 Molden, D., 2007. *Water for Food Water for Life: A Comprehensive Assessment of Water Management in*
687 *Agriculture*. London, Earth Scan.

688 Pekel, J. F., Cottam, A., Gorelick, N. and Belward, A. S.: High-resolution mapping of global surface water and
689 its long-term changes, *Nature*, 540(7633), 418–422, doi:10.1038/nature20584, 2016.

690 Pham, Q.V., Ha, N.T.T., Pahlevan, N., Oanh, L.T., Nguyen, T.B., Nguyen, N.T. Using Landsat-8 Images for
691 Quantifying Suspended Sediment Concentration in Red River (Northern Vietnam). *Remote Sensing*, 10(11), p.
692 1841.

693 Prabha, R., Tom, M., Rothermel, M., Baltasvias, E., Leal-Taixe, L., and Schindler, K. 2020. Lake Ice Monitoring
694 With Webcams And Crowd-Sourced Images, *ISPRS Ann. Photogramm. Remote Sens. Spatial Inf. Sci.*, V-, 549–
695 556.

696 Sasaki, Y., 2007. *The truth of the F-measure*, Manchester: School of Computer Science, University of
697 Manchester.

698 Schanze, J., Zeman, E., Marsalek, J., 2006. *Flood Risk Management: Earth and Environmental Sciences*.
699 Dordrecht: Springer.

700 Smith, S., Kindermans, P.-J., Ying, C. & Le, Q. V., 2017. *Don't Decay the Learning Rate, Increase the Batch*
701 *Size*. s.l., s.n.

702 Sun, W., Du, B. and Xiong, S.: Quantifying Sub-Pixel Surface Water Coverage in Urban Environments Using
703 Low-Albedo Fraction from Landsat Imagery, *Remote Sens.* 2017, Vol. 9, Page 428, 9(5), 428,
704 doi:10.3390/RS9050428, 2017.

705 Tarpley, J., Hirota, N., Kato, M. & Arakane, S., 2015. Combined Effect of an Atmospheric River and a Cut-off
706 Low in Hiroshima Flooding Event on August 19, 2014. *AGUFM*, p. A53F–06.

707 Turner, W. Rondinini, C., Pettorelli, N., Mora, B., Leidner, A.K., Szantoi, Z., Buchanan, G., Dech, S., Dwyer,
708 J., Herold, M., Koh, L.P., Leimgruber, P., Taubenboeck, H., Wegmann, M., Wikelski, M., Woodcock, C.
709 2017. Free and open-access satellite data are key to biodiversity conservation. *Biological Conservation*, Volume
710 182, pp. 173-176.

711 Wang, Y., Li, Z., Zeng, C., Xia G. -S., Shen H. 2020. "An Urban Water Extraction Method Combining Deep
712 Learning and Google Earth Engine," in *IEEE Journal of Selected Topics in Applied Earth Observations and*
713 *Remote Sensing*, vol. 13, pp. 769-782.

714 Wagle, N., Acharya, T. D., Kolluru, V., Huang, H. and Lee, D. H.: Multi-Temporal Land Cover Change Mapping
715 Using Google Earth Engine and Ensemble Learning Methods, *Appl. Sci.*, 10(22), 8083,
716 doi:10.3390/app10228083, 2020.

717 Wieland, M. and Martinis, S., 2020. Large-scale surface water change observed by Sentinel-2 during the 2018
718 drought in Germany. *International Journal Of Remote Sensing*, Volume 41, pp. 4742-4756.

719 Xu, H., 2006. Modification of normalised difference water index (NDWI) to enhance open water features in
720 remotely sensed imagery. . *International Journal of Remote Sensing*, 27(14), p. 3025–3033.

721 Zeiler, M. and Fergus, R., 2014. Visualizing and Understanding Convolutional Networks. *Computer Vision –*
722 *ECCV 2014*, pp. 818-833.

723 Zha, Y., Gao, J. & Ni, S., 2004. Use of normalized difference built-up index in automatically mapping urban areas
724 from TM imagery.. *International Journal of Remote Sensing*, 24(3), p. 583–594.

725 Zhou, Y., Dong, J., Xiao, X., Xiao, T., Yang, Z., Zhao, G., Zou, Z., Qin, Y. 2017. Open Surface Water Mapping
726 Algorithms: A Comparison of Water-Related Spectral Indices and Sensors. *Water*, 9, 256.

727

728