

Adopting deep learning methods for airborne RGB fluvial scene classification.

5 **Patrice E Carbonneau¹, Stephen J Dugdale³, Toby P Breckon², James T Dietrich⁴,
Mark A Fonstad⁵, Hitoshi Miyamoto⁶, and Amy S Woodget⁷**

* Corresponding author:

1: Department of Geography, Durham University, Mountjoy Site, Durham, DH1 3LE, UK.
patrice.carbonneau@durham.ac.uk

10

2: Department of Computer Sciences, Durham University, Mountjoy Site, Durham, DH1
3LE, UK. toby.breckon@durham.ac.uk

15

3: School of Geography, University of Nottingham, University Park, Nottingham, NG7 2RD
UK. stephen.dugdale@nottingham.ac.uk

4: Department of Geography, University of Northern Iowa, 1227 West 27th Street Cedar
Falls, IA 50614, USA. james.dietrich@uni.edu

20

5: Department of Geography, University of Oregon, 1251 University of Oregon, Eugene
OR 97403-1251 USA. fonstad@uoregon.edu

6: Department of Civil Engineering, Shibaura Institute of Technology, 3-7-5 Toyosu, Koto-
ku, Tokyo 135-8548, Japan. miyamo@shibaura-it.ac.jp

25

7: Department of Geography and Environment, Loughborough University, Epinal Way,
Loughborough, Leicestershire, LE11 3TU, UK. a.woodget@lboro.ac.uk

30

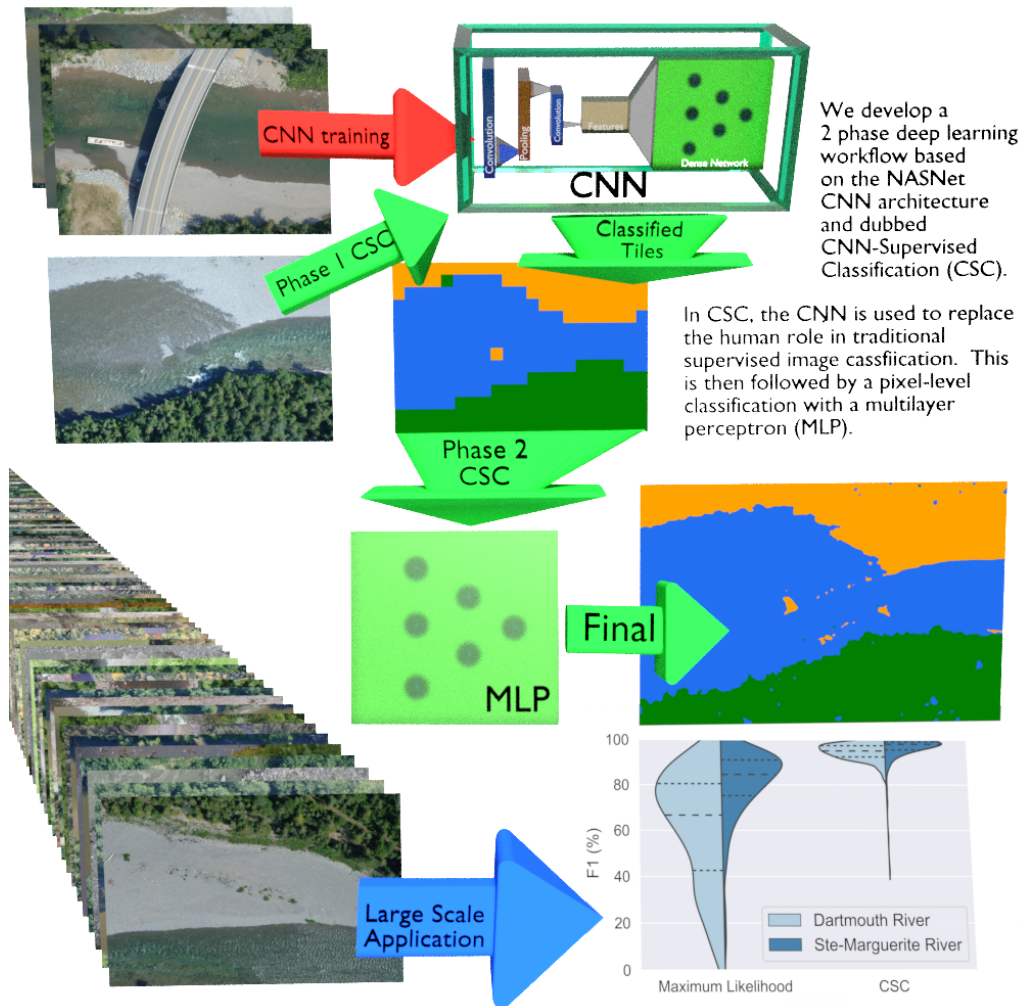
Abstract

Rivers are among the world's most threatened ecosystems. Enabled by the rapid development of drone technology, hyperspatial resolution (<10 cm) images of fluvial environments are now a common data source used to better understand these sensitive habitats. However, the task of image classification remains challenging for this type of imagery and the application of traditional classification algorithms such as maximum likelihood, still in common use among the river remote sensing community, yields unsatisfactory results. We explore the possibility that a classifier of river imagery based on deep learning methods can provide a significant improvement in our ability to classify fluvial scenes. We assemble a dataset composed of RGB images from 11 rivers in Canada, Italy, Japan, the United Kingdom, and Costa Rica. The images were labelled into 5 land-cover classes: water, dry exposed sediment, green vegetation, senescent vegetation and roads. In total, >5 billion pixels were labelled and partitioned for the tasks of training (1 billion pixels) and validation (4 billion pixels). We develop a novel supervised learning workflow based on the NASNet convolutional neural network (CNN) called 'CNN-Supervised Classification' (CSC). First, we compare the classification performance of maximum likelihood, a multilayer perceptron, a random forest, and CSC. Results show median F1 scores (a commonly used quality metric in machine learning) of 71%, 78%, 72% and 95%, respectively. Second, we train our classifier using data for 5 of 11 rivers. We then predict the validation data for all 11 rivers. For the 5 rivers that were used in model training, median F1 scores reach 98%. For the 6 rivers not used in model training, median F1 scores are 90%. We reach two conclusions. First, in the traditional workflow where images are classified one at a time, CSC delivers an unprecedented mix of labour savings and classification F1 scores above 95%. Second, deep learning can predict land-cover classifications (F1 = 90%) for rivers not used in training. This demonstrates the

potential to train a generalised open-source deep learning model for airborne river surveys suitable for most rivers 'out of the box'. Research efforts should now focus on further development of a new generation of deep learning classification tools that will encode
60 human image interpretation abilities and allow for fully automated, potentially real-time, interpretation of riverine landscape images.

Graphical Abstract

65



Highlights

- Deep Learning can classify RGB river imagery to 90%-99% F1.
- This result exceeds the state-of-the-art in fluvial scene classification.
- 70 • Deep Learning models can encode river features that transfer to new rivers.
- Hyper- and multispectral data are not required.
- We provide open source GIS integration via PyQGIS.

75 Introduction

Freshwater environments and the flora and fauna they contain are among the most threatened ecosystems on the planet (Carrizo et al., 2017; Strayer and Dudgeon, 2010; WWF, 2018). Of these habitats, rivers in particular have been the focus of intensive research and conservation initiatives (e.g. Linke et al., 2007; Nel et al., 2009; Ormerod, 80 2009) due to the combined threats of impoundments and flow alteration (Rosenberg et al., 2000; Vörösmarty et al., 2010), land-use modification (Rogger et al., 2017; Zhang and Schilling, 2006), and climate change (Arnell and Gosling, 2016; van Vliet et al., 2013). In tandem with this increasingly intensive 'applied' research focus, recognition has grown that the improved conservation of river environments will naturally stem from a deeper 85 understanding of patterns and processes in physical river habitats (e.g. Palmer et al., 2010; Ward et al., 2001; Wohl et al., 2005) and their linkages to aquatic organisms. Indeed, this concept is central to the *riverscapes* paradigm (Fausch et al., 2002), which dictates that a spatially continuous view of the river is key to understanding and conserving stream biota. The collection and assembly of high-resolution data pertaining to river 90 environments is therefore a fundamental first step in protecting these critically endangered global ecosystems (Vannote et al., 1980).

The sinuous, dendritic nature of rivers, coupled with the difficulty of conducting spatially-intensive sampling in aquatic environments, has led researchers to increasingly turn to 95 remote sensing to provide the spatially continuous data necessary to yield improved fundamental and applied understanding of river environments. The sub-discipline of fluvial remote sensing tends to be divided in 2 principal areas depending on the acquisition platform: spaceborne and airborne. Spaceborne river remote sensing tends to address continental and global scale issues (e.g. Allen and Pavelsky, 2018; Durand et al., 2016;

100 Gleason and Smith, 2014; Smith, 1997). Airborne fluvial remote sensing (hereafter abbreviated as airborne FRS) focusses on local, sub-meter scale features and parameters that can be algorithmically retrieved and generally depend on the much higher spatial resolutions, typically cm-scale, made possible by airborne acquisitions (Carbonneau and Piégay, 2012a). Earlier airborne FRS work (e.g. Seto et al., 2002; Winterbottom and
105 Gilvear, 1997; Yang et al., 1999)(eg. Seto et al. 2002; Winterbottom & Gilvear, 1997; Yang et al. 1999)(eg. Seto et al. 2002; Winterbottom & Gilvear, 1997; Yang et al. 1999)(eg. Seto et al. 2002; Winterbottom & Gilvear, 1997; Yang et al. 1999) highlighted the utility of multi/hyperspectral satellite and airborne platforms for mapping fluvial environments. While these coarser spatial resolution data continue to be useful for monitoring rivers, particularly
110 planform change, hydrometry or water quality; (e.g. Bjerklie et al., 2003; Kuhn et al., 2019; Langat et al., 2020)(eg Bjerklie et al., 2003; Kuhn et al., 2019; Langat et al., 2020), algorithms for quantifying in-stream metrics such as grain size, water depth, or sub-meter scale temperature topography (e.g. Black et al., 2014; Dietrich, 2016; Willis and Holmes, 2019) require the acquisition of very high resolution (usually <10 cm) RGB images which
115 are not available from any satellite platform. Carbonneau and Piégay (2012a) define 'hyperspatial' resolution threshold as <10 cm and state that such images have increasing value in the analysis of river systems. Downing et al. (2012) also estimate that 97% of the world's rivers by length have a width below 30m. Similarly, Allen and Pavelsky (2018) estimate that 369 000 km² of the earth's surface are occupied by rivers smaller than 90m.

120 The prevalence of small streams therefore creates a niche for hyperspatial data capable of resolving even narrow channels with hundreds to thousands of pixels per average river width. Indeed, supported by the explosion of drone-based remote sensing techniques over the last 10 years (e.g. Woodget et al., 2017; Woodget and Austrums, 2017), the extraction of river habitat data from hyperspatial RGB imagery is fast becoming a mature

125 and accepted technique in the river sciences (Bagheri et al., 2015; Black et al., 2014; Carbonneau et al., 2012; Carbonneau and Piégay, 2012; Dugdale et al., 2019; Hamshaw et al., 2017; Kalacska et al., 2019; Michez et al., 2016; Tamminga et al., 2015; Woodget et al., 2016, 2015). A similar body of work also describes the use of RGB images acquired from terrestrial platforms for extracting a range of fluvial characteristics

130 (e.g. Ashmore and Sauks, 2006; Butler et al., 2001; Chandler et al., 2002; Ghaffarian et al., 2020; B. MacVicar and Piégay, 2012; MacVicar et al., 2012; B. J. MacVicar and Piégay, 2012; Purinton and Bookhagen, 2019). Additionally, the increasingly ubiquitous use of structure from motion (SfM) photogrammetry in river remote sensing (e.g. Carrivick and Smith, 2019; Hemmelder et al., 2018; Seitz et al., 2018) - a technique reliant on sub-

135 decimeter (RGB) imagery – means that hyperspatial RGB imagery acquired from airborne platforms is a widely exploited river remote sensing data type and allows for small scale investigations that are not possible with orbital sensors.

In studies involving the quantification of river habitat from remote sensing, it is often

140 necessary to first distinguish the wetted channel from other land cover types prior to the application of algorithms to extract hydromorphic metrics (e.g. Carbonneau et al., 2012, 2006). However, this basic task of image classification remains challenging for RGB hyperspatial imagery where the extremely fine spatial detail and relatively low number of spectral bands means that traditional statistical learning classification algorithms (e.g.

145 maximum likelihood or k-means clustering) that are still widely-used among river remote sensing practitioners (Brigante et al., 2017; Spada et al., 2018; Wang et al., 2016) have difficulty correctly allocating image pixels to semantic classes that are radiometrically similar to one another. For example, riparian vegetation and river water often share dark green colours that make them very difficult to distinguish on a purely spectral basis in the

150 visible RGB range of the spectrum. Shadows, both in amongst the vegetation and cast by
vegetation, compound the problem and make classification even more difficult. In certain
cases, deeply shaded sediment can even be spectrally similar to shallow water. The end
result is that classification of river imagery is a very challenging problem and progress in
the field has been somewhat limited. Indeed, despite rapid advances in image
155 classification within other fields (e.g. computer vision), river remote sensing studies do not
achieve classification accuracies above 90%, (e.g. Boruah et al., 2008; Casado et al.,
2015; Demarchi et al., 2020; Gilvear et al., 2008; Legleiter and Goodchild, 2005; Marcus et
al., 2012; Rusnák et al., 2018; Smikrud et al., 2008; Wang et al., 2016). This is largely
because at meter-scale and centimeter-scale resolution, the assumption that a semantic
160 class can be described by a set of unimodal distributions of brightness values is not
necessarily valid. Furthermore, the incredible global variety of rivers means that
classification techniques solely based on radiometric properties are unlikely to be
successful when applied to other, less radiometrically variable land-use types. This
reliance on the use of outdated algorithms and the resulting difficulty in classifying riverine
165 imagery is a pressing problem in the airborne FRS community. A prime example is the
lack of an automated workflow that can approach human performance when identifying the
wetted perimeter. Indeed, not only does this classification difficulty currently prohibit the
easy application of advanced image processing algorithms for the extraction of physical
habitat data (Carbonneau et al., 2012), but also severely limits our ability to explore
170 patterns and processes in channel morphology at riverscape scales.

In the area of airborne FRS, previous efforts to solve the challenges of sub-meter
resolution image classification have been dominated by hardware approaches involving
the use of multi- or hyperspectral sensors (Demarchi et al., 2017, 2016; Laliberte et al.,

2011; Legleiter et al., 2004, 2002; Marcus et al., 2003; Olmanson et al., 2013; Tian et al., 2010; Zhong and Zhang, 2012). The main finding of this body of literature is that the addition of spectral detail, including information from non-visible, infrared wavelengths, greatly enhances classification performance. This improvement occurs because the inability of near-infrared wavelengths to penetrate water render the wetted channel easy to segment from terrestrial features which otherwise have similar spectral signatures in visible wavelengths. Indeed, using such multispectral data, Marcus et al. (2003) report accuracies as high as 86% for the classification of a fluvial landscape. Demarchi et al. (2020) use the infrared band and a DEM layer in an object-based approach to reach an overall performance of 89%. However, validation of these studies is typically carried out by visual labelling, often using RGB images. Given that a trained human observer is readily capable of delimiting land-cover classes in RGB imagery, so-called ‘Artificial Intelligence’ methods potentially could solve this classification problem without the need for costly multi- and hyperspectral sensors. Such methods hold great promise for raising the classification accuracy of river remote sensing data to the >90% levels currently considered the state-of-the-art in computer vision and related fields (e.g. Barré et al., 2017; Debats et al., 2016; Hernández-Serna and Jiménez-Segura, 2014)

Chollet (2017) defines artificial intelligence (AI) as ‘*the effort to automate intellectual tasks normally performed by humans*’. The author then introduces the terms: Machine Learning (ML) and Deep Learning (DL) with mutually inclusive sets:

$$AIC \supset MLC \supset DL \quad (1)$$

Machine learning is therefore a subset of artificial intelligence methods where any algorithm is capable of learning and encoding prediction rules from data (Chollet, 2017;

200 Goodfellow et al., 2016). Deep learning methods, a subset of machine learning methods, distinguish themselves by their ability to encode multiple layers of features learned from large datasets (LeCun et al., 2015). In practice, deep learning relies on convolutional neural network (CNN) (Goodfellow et al., 2016; Lecun et al., 1998) architectures, whereby a locally-weighted operator performs a variety of de-noising, feature extraction, and data
205 reduction operations by varying only the weights of the convolution operator itself (Solomon and Breckon, 2011). Such deep learning architectures essentially offer a huge parametrisation space which is then tuned to recover an optimal set of feature extraction/classification parameters as a set of neural network operations (i.e. image in; classification out).

210

Advances in deep learning (e.g. Zhang et al., 2016) have started to show great potential for the classification and segmentation of diverse landscape features from remote sensing data. Convolutional neural networks are being used increasingly for large-scale satellite image classification (e.g. Chen et al., 2019; Kussul et al., 2017; Romero et al., 2016; 215 Zhong et al., 2017), enabling the segmentation of imagery into broad classes (e.g. trees, grassland, soil, roads, water) with accuracy substantially greater than conventional classification techniques. . However, in the specific context of rivers, applications of deep learning are sparse. Casado et al. (2015) used a non-convolutional artificial neural network (often called a multilayer perceptron) to identify hydromorphic units in a river reach. Daigle
220 et al. (2013) demonstrated a similar perceptron-based approach to detect river ice from fixed RGB imagery. More recently, Isikdogan et al (2018) and Ling et al. (2019) highlighted the utility of deep learning for extracting channel characteristics from satellite imagery, and Buscombe and Ritchie (2018) have applied DEEPLAB (Chen et al., 2018) to landscapes (including rivers), demonstrating that river corridor classification in high resolution imagery

225 with deep learning is possible. However, while successful in isolation, the uptake of these methods has been slow, and the lack of a deployable, repeatable and accurate classifier for river corridors remains a crucial issue in river remote sensing.

The complexity of implementing deep learning approaches partially accounts for the lack
230 of uptake among river scientists and managers who are not trained in computer vision. However, the specificity of deep learning methods has potentially also prohibited their wider application in the airborne FRS domain. While previous research using CNNs has demonstrated an ability to achieve extremely high classification accuracy when deployed in a target recognition sense (e.g. Foody et al., 2019; Guo et al., 2018; Li et al., 2017), the
235 transferability of these approaches across diverse riverine landscapes and remote sensing systems/platforms remains untested. Unlike relatively homogeneous landscape types (e.g. urban environments or forest canopies; Khan et al., 2017; Mahdianpari et al., 2018; Pouliot et al., 2019) where deep learning has previously seen success, rivers are extremely heterogeneous. This diversity implies that the development of a fully transferable classifier
240 for river corridors from environments as disparate as the tropics or alpine regions is an extremely complex problem. Furthermore, despite the number of CNN-based landscape/land-use classification approaches documented in the literature, there is a relative absence of examples that are ready for deployment in an environmental management context. Given that one of the key factors precluding the use of advanced
245 image processing techniques in the applied river sciences is the lack of coding or scientific computing expertise among environmental management communities, these issues create a compelling need for the development of a high quality, transferable and easy to use image classifier for use in the river sciences.

250 Current options for deep learning approaches in commercial and/or open-source remote sensing packages are limited (Table 1). Indeed, with the exception of the *Orfeo* open source toolbox, only high-end, high-cost, commercial products currently have built-in implementations of deep learning. Not only are these packages rarely available to river management organisations, they also offer very limited flexibility to adapt algorithms to
255 specific cases such as hyperspatial imagery. Another common issue with all machine learning algorithms deployed in software is that the requirement for training and validation. In the case of deep learning, the need for large labelled sets of data makes software implementation even more problematic. Indeed, the dominant paradigm in classification of Earth Observation (EO) data is that the user manually draws polygons on-screen in order
260 to form labelled pixels for supervised classification training. In the case of deep learning this is very problematic since the human effort required to generate a sufficient sample size is very considerable. We argue that implementation of deep learning in research fields with lower levels of computer vision expertise would be greatly facilitated if pre-trained, freely available, deep convolutional networks could be called upon to classify new image
265 data without the need for the labour intensive and time-consuming generation of new training labels. Such a facility would not only be a substantial boon for the classification and interpretation of new airborne FRS data, but would also greatly enhance the extraction of river habitat data from archival aerial photography acquired over the past 20 years during the emergence of the airborne FRS sub-discipline. This would aid the rapid
270 detection and analysis of river habitat change in the context of land-use and climate modification, and allow for improved testing of prevailing theories regarding hydromorphic processes and the linkages between river habitats and ecosystems.

275 **Table 1. Supervised classification workflows currently available within remote sensing software packages.**

Software Package	Machine Learning	Deep Learning	Access Type
eCognition	✓ Yes (e.g. decision trees, random forests, support vector machines)	✓ Yes (uses Google TensorFlow library, including trainable convolutional neural network models)	Commercial
ERDAS Imagine Professional	✓ Yes (e.g. random forests, support vector machines, CART)	✓ Yes (e.g. Faster regional-based convolutional neural networks)	Commercial
ENVI	✓ Yes (Interactive data language framework: e.g. support vector machines, SoftMax, Feed Forward Neural Network-based classifications)	✓ Yes (Deep learning module built on Google TensorFlow)	Commercial
ESRI ArcPro	✓ Yes (e.g. Random trees, support vector machines)	Support for export to third party deep learning tools	Commercial
QGIS + GRASS	✓ Yes (e.g. Gaussian mixture models, random forests, support vector machines)	✗ No	Open Source
SAGA	✓ Yes (e.g. support vector machines)	✗ No	Open Source
Orfeo Toolbox	✓ Yes (e.g. Support vector machines, Bayes, Random forests)	✓ Yes (e.g. <i>otbtf</i> module built on Google's TensorFlow)	Open Source

280 The overarching aim of this work is therefore to examine the potential of machine learning and deep learning in the specific context of hyperspatial airborne FRS. We do not claim to advance the field of deep learning, and we recognize that ‘cutting edge’ computer vision approaches involve even more advanced algorithms than those considered here (e.g. Long et al., 2015). Rather, our intention is to advance the state-of-the-art in fluvial scene

285 classification by quantifying accuracy improvements possible with deep learning approaches that are sufficiently mature and established to be accessed and manipulated by non-specialists and, ultimately, integrated into a GIS workflow. Furthermore, we wish to understand if deep learning classifiers can mimic a human expert and consistently classify riverine land-cover types in hyperspatial (<10cm) resolution colour imagery to higher levels

290 of accuracy (>90% F1 and above) than those common to past and present river remote sensing studies typically performing in the 70%-90% range (e.g. Boruah et al., 2008; Casado et al., 2015; Feng et al., 2018; Gilvear et al., 2008; Legleiter and Goodchild, 2005; Marcus et al., 2012; Rusnák et al., 2018; Smikrud et al., 2008; Wang et al., 2016). Indeed, given that river habitats are highly complex environments characterised by gradients and

295 discontinuities (Fausch et al., 2002), the ability to improve classification accuracy above current norms is crucial for accurately identifying small discontinuous habitat features that may have a disproportionate role in key ecosystem processes. In this manner, even relatively incremental increases in classification accuracy (e.g. from ~80% to >90%) have the potential to yield major advances in our understanding of fluvial forms and dynamics by

300 yielding a fuller picture of spatial patterns in key habitat features that might have been misclassified by less advanced techniques. Our study therefore has three specific objectives: First, we compare the performance of a range of land-cover classifier algorithms (maximum likelihood, Random Forests, depth-limited Neural Networks, and Convolutional Neural Networks) in order to demonstrate the potential of deep learning

305 methods to fluvial scientists and river managers. Second, we evaluate the potential of a deep learning workflow called CNN-Supervised Classification to transform current practice in river land-cover classification where classifiers are trained to predict land-cover for single rivers, one at a time. Third, we critically assess the future potential of CNN-Supervised Classification as a transferable classifier eventually capable of river corridor
310 segmentation without the need for further model training. Finally, we demonstrate GIS integration and direct readers to an open-source code repository ready for deployment.

Methods

Hardware and software

315 We use capable but modest resources accessible to most researchers. The data presented here were processed with two laptops. The main unit had a 4-core Intel i7-6820 CPU clocked at 3.4 Ghz with 32 Gb RAM and an NVIDIA GTX 1070 GPU with 8Gb of memory and 1920 CUDA cores available for parallel processing. The secondary unit had a 4-core Intel i7-4700MQ CPU clocked at 3.4 Ghz but with 24 Gb RAM and an NVIDIA
320 Quadro K1100M GPU with 2 Gb of memory and 384 CUDA cores. With these laptops and using the data volumes described below, training times for the deep networks ranged from 1 to 5 hours. Classification of a single image required 2-5 minutes. Whilst these are moderately high specifications for laptops, equivalent desktops are readily available.

All software used in this work is open-source. Core deep learning work was undertaken in
325 Python 3.6 using the Anaconda distribution. We use the *scikit-learn* library for classification metrics and for the random forest machine learning algorithm (Pedregosa et al., 2011). We use *scikit-image* for basic image import/export and more advanced processing and filtering (Walt et al., 2014). For dense and convolutional neural networks, we use the *Keras* API (Chollet, 2017) running the GPU-enabled version of TensorFlow

330 v1.14 (Abadi et al., 2016). The *Pandas* library is used for basic tabular data storage, manipulation and management (McKinney, 2010). Visualisation is delivered with the *Seaborn* library. Spyder (Scientific PYthon Development EnviRonment) was used as the main integrated development environment (IDE) for coding and debugging. We deliberately avoided CNN architectures that are closer to the research frontier and instead
335 sought a deep learning architecture that is established and ready for deployment. Within the *Keras* API, we selected the pre-existing NASNet Large CNN model (Zoph et al., 2017) because it has the highest prediction accuracy according to the *Keras* documentation. Furthermore, we also decided to test the NASNet Mobile architecture to explore whether a smaller version of the NASNet architecture could deliver good results with less
340 computational overhead. For digitising and GIS tasks, we use QGIS 3.4 *Long Term Release* distributed with an integrated version of GRASS GIS 7.6. GRASS GIS is used to perform maximum likelihood classification. GIS integration is achieved by installing all the libraries listed above in the QGIS python environment. PyQGIS can then be used to geocode the classification outputs and run the entire process from the QGIS Python
345 console.

Data preparation

We use existing data and have compiled a database of hyperspatial resolution imagery. Our objective was to compile a database with several billion labelled pixels that included a
350 wide range of rivers from diverse morpho-sedimentary settings. Another essential criteria was that imagery be available under an open source license and made freely available as part of this paper. We found that availability and absence of intellectual property or licensing restriction was the major limiting factor. Ultimately, we assembled a database

with imagery from 11 rivers in Canada (Quebec and Alberta), Italy, Japan, the UK and
355 Costa Rica (Figure 1).



Figure 1. Location map for the 11 study rivers.

360

We argue that this is a state-of-the-art dataset which is more varied than anything previously presented in the high resolution airborne FRS literature. Within this subset of the remote sensing literature, we recognize that there are a small number of publications with datasets that exceed our own in terms of sheer number of pixels (e.g. Black et al.,
365 2014; Carbonneau et al., 2004). However, such studies are usually focussed on a single river and we find no other report in the peer-reviewed literature with in excess of 5 billion labelled pixels distributed among 11 rivers spanning the Americas, Europe, and Asia. Our images were acquired between 2002 and 2017 from both piloted aircraft and unpiloted

aircraft systems (UAS). The images are composed of what might be termed a standard
370 view in airborne FRS, with the channel roughly in the centre of the scene and with
vegetated areas and frequent occurrences of exposed sediment on either side. Most
images are dominated by green vegetation but some sets have a frequent occurrence of
different types of senescent vegetation that ranges from the dry grasses of the Scottish
Highlands to bright autumn foliage in eastern Canada. Water colour varies substantially
375 and also contains instances of sun glint, white water and shadows. Man-made features
are rare but sometimes present, mostly in the form of (paved) roads. Sediment type also
varies. The Dartmouth, Ste-Marguerite, Kananaskis, Ouelle and Pacuare rivers are single
thread channels with sediment bars. The reach of the Eamont used here is a typical
English lowlands river with dark peaty waters, pasture banks and relatively few sediment
380 bars. The Kingie river is a Highland river with very dark peaty water, banks dominated by
senescent grasses but punctuated by fir trees and very coarse and angular sediment. The
Dora di Veny river is our only Alpine river. The Sesia river is an anastomosing channel
with a relatively high and coarse sediment load and, finally, the Kurobe and Kinu rivers are
braining channels in a densely populated area with significant occurrences of roads.
385 Figures 2 and 3 give names and thumbnail examples of each river with basic
characteristics and a sample of final classification outputs from the results discussed
below. Most of the imagery was available as original single frames. However, the Kurobe
and Kinu rivers were only available as large image mosaics. These were separated into
tiles of 2250X2250 pixels in order to match the format of the other images and allow for a
390 common data management and processing scheme.

Some of the imagery has an existing classification available for usage, from a variety of
sources such as manual classification and both semi- and fully-automated classification

methods. Specifically, the images of the Dartmouth river in Canada had an existing
395 classification derived from eCognition software (then known as Definiens) in 2008. Pre-
existing classifications of the Ste-Marguerite and Ouelle rivers were achieved using the
approach of Carbonneau et al. (2004), whereby a semi-automated classification method
using a combination of thresholding and extensive manual editing was applied. We began
by a manual inspection of each of these pre-classified images in order to insure that the
400 labels were accurate. For the other rivers, no existing classification data was available.
We therefore use QGIS 3.4 to manually label portions of each available image. The
objective of this classification was not to derive detailed classification labels for each pixel
in each image. Rather, our objective was to rapidly develop an overall dataset with a large
number of labelled pixels suitable for training and validation of machine learning
405 classifiers. Prior to classification, we examined the imagery to derive a parsimonious
classification system that would encompass the main elements present in the dataset. We
decided to establish five training classes: water (class 1), sediment (class 2), green
vegetation (class 3), senescent vegetation (class 4) and paved roads (class 5). We also
observed that shadows appeared in many images but decided to classify shadow patches
410 as per the underlying land-cover type; i.e. shaded water was classified as water, shaded
sediment was classified as sediment, etc...The QGIS digitising tools were used to classify
portions of each image according to this scheme, and the resulting vector polygons
rasterised to derive class rasters of the same spatial resolution and extent as the
associated image. The QGIS graphical modeller was used to batch process this
415 rasterisation operation. For the Ste-Marguerite, Dartmouth, and Ouelle rivers where class
label data already exists, we recoded the data to conform to our classification scheme. All
classification was conducted in such a way that classes contained a representative
coverage of all pixels within a semantic class; for example, 'water' contained pixels

including shadows, sun glint and white water, so as not to present a biased 'best case'
420 classification scenario.




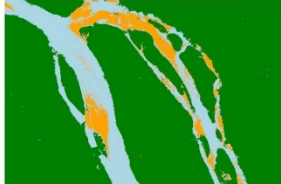





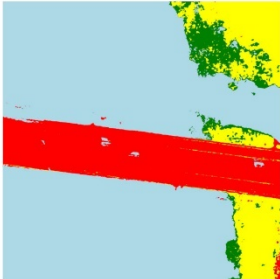
As stated above, this work aims to build on the resources developed in computer vision in order to develop a state-of-the-art method for fluvial scene classification. One of the key datasets that has driven progress in image recognition tasks is the ImageNet database of
425 millions of images (Deng et al., 2009). The images tend to be of common categories such as 'cat' or 'dog'. However, the database is also organised in a hierarchical manner with each class having subdivisions such as 'persian' and 'maincoon' as cat breeds and 'poodle' and 'labrador' as dog breed. In total there are over 1000 classes in the ImageNet databases. We have therefore constructed a classification scheme which is also
430 hierarchical and borrows ideas from the field of hierarchical segmentation (e.g. Li et al., 2011; Poggi et al., 2005). We do not assume that a given semantic class has similar radiometric properties in different image sets. For example, we do not assume that the vegetation in the Ste-Marguerite river data generates pixels of similar radiometric response to that of the Pacuare or Kurobe rivers, owing to both a) real differences in the vegetation's
435 spectral signature and b) variations in recorded pixel values owing to the use of different cameras which are not radiometrically calibrated. Furthermore, differences in vegetation species and structure, in water colour due to local conditions and in sediment texture due to local geology mean that the image properties of a given class can diverge significantly for different study sites. Therefore, in order to work with multiple classes across multiple
440 rivers, we developed a micro-class labelling procedure for training machine learning algorithms. In cases where we train a classifier with data from multiple rivers, semantically identical classes from multiple rivers are transformed to unique micro-classes after manual labelling. For example, when working with a single river, we use 5 classes labelled 1 to 5

as defined above. If we work with 2 rivers, the classes of the second river are shifted to
445 values of 6 to 10, thus resulting in 5 semantic classes (or macro-classes) and 10 micro-
classes. A classification key then records that classes 1 and 6 are water; 2 and 7 are
sediment; 3 and 8 are green vegetation, etc. This process can be extended to as many
rivers / micro-classes as required. At the end of the classification process, the information
in the classification key is used to collapse the classes back to the 5 unique semantic
450 classes.

Following image labelling, we divided our data into training and validation data sets based
on the actual number of labelled pixels. Result validation in deep machine learning follows
the principle of reporting statistical performance on a randomly selected subset of the
455 available data set used for the study. It is established practice, to randomly split the
dataset into either 70%/30% or 80%/20% subsets with the smaller set being used for
testing (evaluation) and hence statistical reporting of results in the literature (Bishop,
2006). Normally, algorithm (DL CNN / deep net) training is performed using the larger of
the two subsets (70% or 80%). This is established practice defined by the leaders of the
460 deep learning field (LeCun et al., 2015b). Given the size of our database, we decided to
use a 20%/80% split in order to reserve more pixels for validation and reduce
computational loads to manageable levels during the training phase. Smaller volumes of
data at the training stage also allows for simpler deep learning code that loads the training
data into available RAM memory.

465 As a basic design criterion, we aim to classify image patches composed of mostly pure
classes. In order to classify patches, it then becomes necessary to tile the input image.
Preliminary-experiments indicated that a 50 x 50-pixel tile size in the NASNet convolutional
neural network architectures represented an optimal balance in terms of processing time

and classification quality. We only used image tiles that were 90% occupied by a pure
470 class. We also decided to retain the same number of training tiles for each river rather
than having a variable number of tiles/pixels available for validation. This approach
resulted in approximately 38 000 training tiles for each river. Table 2 gives full details of the
volumes of available data. In total, our training data is composed of 405 768 labelled,
single-class tiles of 50x50 pixels that could be used to train predictor models that could in
475 turn be validated with up to 861 images having 4.36 billion labelled pixels. Table 3 details
the population of classes in training and validation sets.

River	Thumbnail image	Thumbnail class raster	Size [pix] and GSD [cm]
Ste-Marguerite			3008 X 1908 3 cm
Kananaskis			5184 x 3456 3 cm
Kingie			4000 X 3000 2 cm
Sesia			4000 x 3000 2 cm
Kinu			2250 x 2250 3 cm

480 **Figure 2. Image and Classification samples for 5 of 11 rivers in the image dataset. The classification sample is taken from results of the third experiment described in this paper. In the classification rasters, blue denotes water, green denotes fresh vegetation, yellow denotes senescent vegetation, orange denotes exposed sediment and red denotes paved roads.**


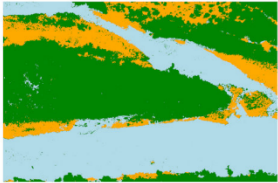

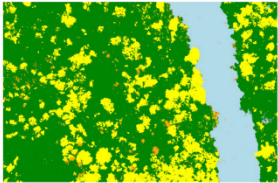

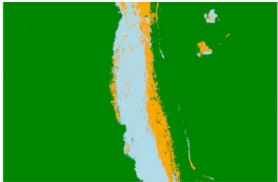

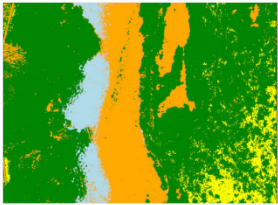


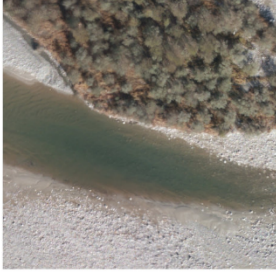
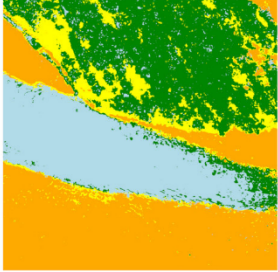
River	Thumbnail image	Thumbnail class raster	Size [pix] and GSD [cm]
Dartmouth			4288 X 2848 3 cm
Ouelle			5184 x 3456 3 cm
Pacuare			5184 x 3456 10 cm
Eamont			4000 x 3000 2cm
Dora di Veny			4000 x 3000 2cm
Kurobe			2250 x 2250 3 cm

Figure 3. Image and Classification samples for the remaining 6 rivers with the same classification key as in figure 2.

490 **Table 2. Data Availability. Readers should note the loose correlation between the number of images and training tiles. In cases where the image scenes were composed of large uniform areas, it was easier to rapidly classify large areas. In such cases, like the Dora di Veny river, fewer images are required to assemble the ca. 37K tiles. Conversely, in complex landscapes such as the Kinu river, a larger**

495 **number of images was required to reach the required volume of training data. We also include two large orthomosaics of 1 km reaches of the St-Marguerite and Kurobe rivers with a spatial resolution of 7.5cm that will be used to demonstrate GIS integration. * denotes rivers used in the training data for the second experiment.**

500

River	Training Data			Validation Data	
	# of Images	# labelled tiles	# labelled pixels	# of Images	# labelled pixels
St-Marguerite, Canada*	44	38,041	109,993,375	224	955,482,438
St-Marguerite, orthomos.	0	0	0	1	29,687,610
Ouelle, Canada	29	37,396	94,797,100	117	424,805,106
Dartmouth, Canada	17	36,443	100,823,415	243	1,671,866,288
Kananaskis, Canada*	16	37,010	104,521,527	34	419,790,696
Pacuare, Costa Rica	25	37,271	100,746,483	38	150,388,739
Sesia, Italy*	26	37,337	101,222,965	21	80,943,299
Dora di Veny, Italy	10	36,696	98,080,466	28	249,874,235
Kingie, UK*	24	35,315	95,272,952	15	50,634,616
Eamont, UK	23	36,991	100,538,759	9	42,543,651
Kinu, Japan*	53	37,057	102,751,306	54	107,686,602
Kurobe, Japan	38	36,211	98,641,597	78	206,807,563
Kurobe, orthomos.	0	0	0	1	88,471,766
TOTAL	305	405,768	1,107,389,945	863	4,390,510,843

505

Table 3. Class representation across training and validation datasets.

Class	Training Data		Validation Data	
	# labelled pixels	% Total	# labelled pixels	% Total
Water	419,580,438	38	2,007,533,862	39
Sediment	269,759,643	24	575,600,093	11
Green Vegetation	343,385,982	31	2,408,766,446	47
Senescent Vegetation	72,698,683	7	96,807,737	2
Paved Roads	8,173,202	1	11,237,600	0

510

CNN-Supervised classification approach

CNN-Supervised classification (CSC) is a novel two-phase workflow that chains a deep convolutional neural network with a multilayer perceptron in order to deliver pixel-level classification in a deep learning workflow based on convolutional architectures. In phase 1, the input image is tiled with multiple tiles stored as a single 4D tensor with dimensions of (Tiles, X, Y, RGB Bands) and fed into a pre-trained CNN. This is analogous to having a single video file which is a 4D temporal sequence of RGB images. The use of a pre-trained CNN as the first phase is crucial because it allows for a local association between a class and predictive features such as local brightness, local texture and even local geometric structures (eg branches, boulders). In phase 2, the resulting CNN predictions in the form of labelled tiles are rasterised and re-assembled into the shape of the original image. For example, if the CNN has predicted the class of each tile of 50x50 pixels, then

520

525 each class prediction is converted into a small raster of 50x50 pixels with a uniform value corresponding to the class. These small 50x50 rasters are reassembled into the shape of the original image with zeros used to pad edges. This CNN-derived class prediction raster is used as labelled pixels and, along with RGB features, is then fed into a multilayer perceptron (MLP) in order to train a model specific to the input image. Finally, this MLP
530 (detailed in Table 4), is used to predict the class of each pixel in the original image. Our intent is to mimic the traditional supervised land-cover classification workflow in which a human operator outlines training areas of desired classes, which are then fed into a machine learning algorithm. In CSC, the CNN replaces the human operator, with a MLP used as the specific machine learning algorithm. We demonstrate the benefit gained from
535 the MLP's characteristic robustness to noise in the training data (in this case, the CNN predictions).

CSC requires a pre-trained CNN. In the work presented here, the CNN is trained with our own data, presented below. Goodfellow et al. (2016) suggest that *ca.* 10 million samples
540 are required to train a deep learning algorithm to the point of matching human performance. Therefore, as in Buscombe and Ritchie (2018), we decided to use a transfer learning procedure whereby initial model weights are imported from an existing dataset in order to allow the CNN to train with a smaller dataset. We use the initial weights as derived from the ImageNet database. This database is an archive composed of in excess
545 of 1 million tiles and serves as a benchmark for AI performance. For the NASNet CNN architectures, we freeze all the weights except those of the top 4 convolutional layers. This results in 11,515,046 *trainable* parameters out of a total of 89,079,512 parameters for NASNet Large. For NASNet Mobile, we have 1,484,986 trainable parameters out of a total

of 3,902,580. Figure 4 shows the generic workflow of CNN-supervised classification
550 inclusive of the pre-training of the CNN and the two-phase classification workflow.

In addition to the CNN base architecture, a CNN classifier requires a ‘top’ neural network to convert the features detected by the CNN into classes represented as integer numbers. In this case the densely connected top is composed of 3 additional layers: a dense layer of 256 nodes, a drop out layer (Szegedy et al., 2015), a dense layer with 128 nodes, and,
555 finally, the usual softmax layer with the same number of nodes as classes. This layer functions by returning the final probability that an image tile is a member of each class. By convention, the final attributed class is the one with the highest probability of membership. For both dense layers, we use kernel L2 regularization in order to inhibit over-training (Goodfellow et al., 2016). The overall CNN-supervised classification process can be seen
560 in Figure 4. Once CNN model training is complete (the upper part of Figure 4), the resulting CNN can be re-used for multiple images. In the experiments described below, we examined increasingly ambitious scenarios up to the point where the process was expected to classify entirely new rivers never seen by the classifier.

565

570

575

580

585

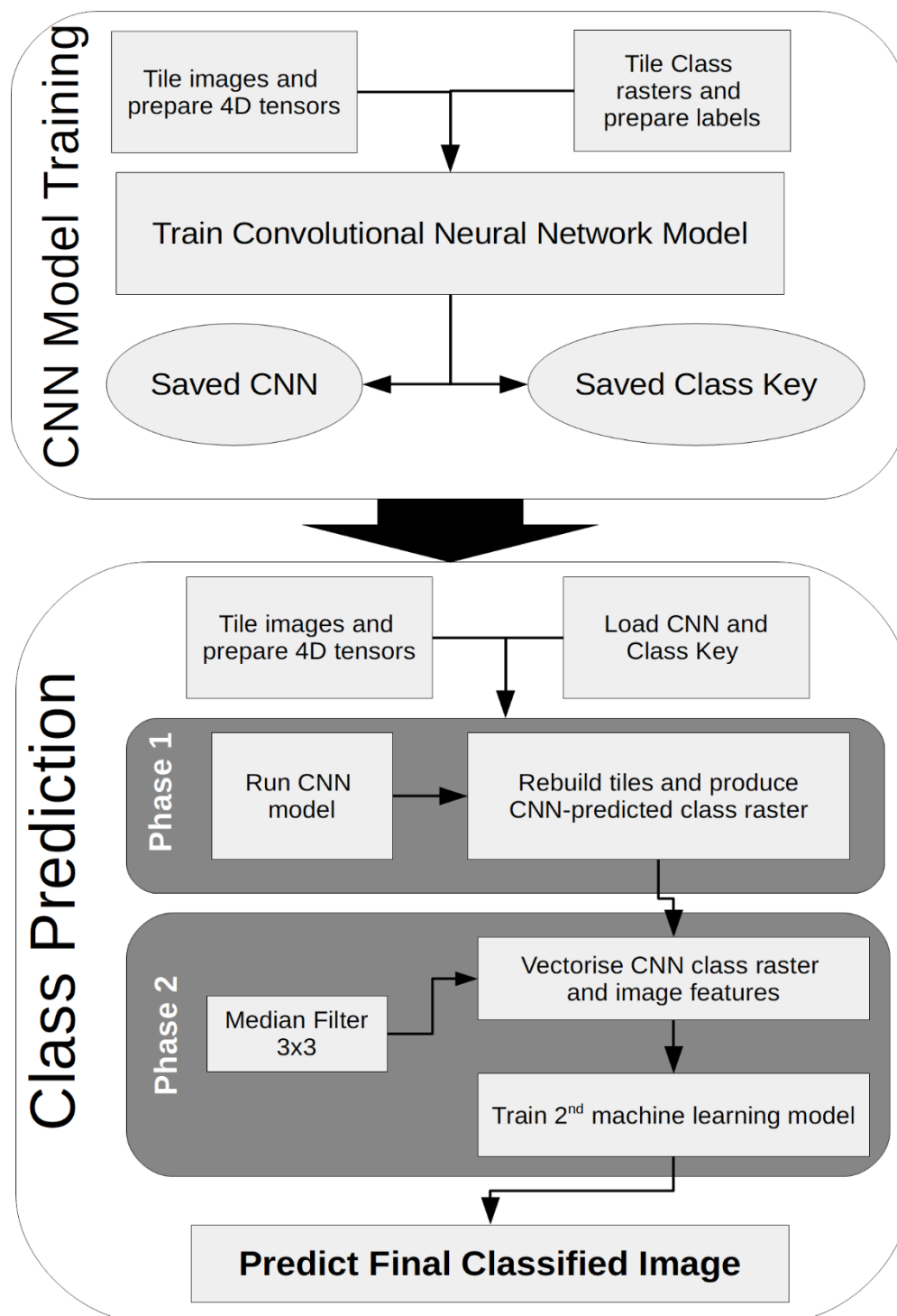
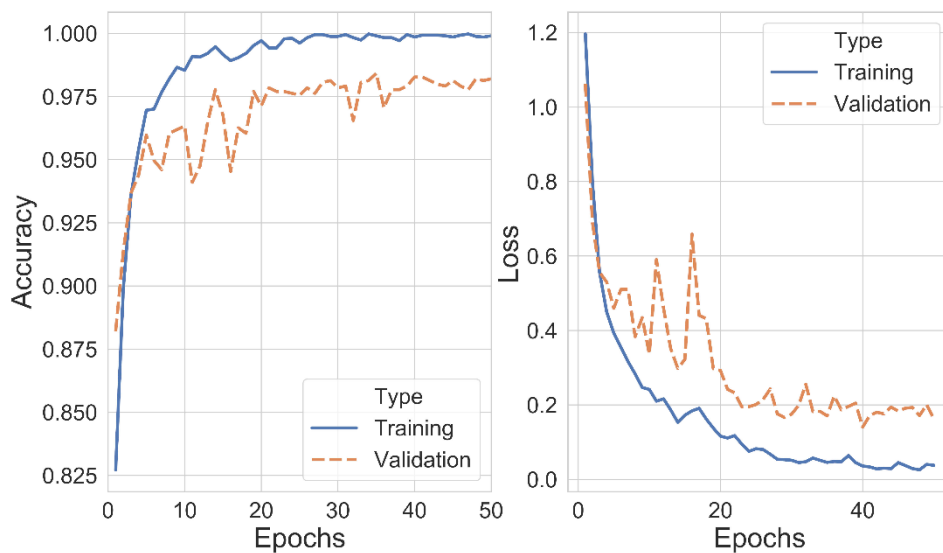


Figure 4. Workflow chart for the generic process of CNN-Supervised Classification.

CNN training is known to be sensitive to the number of training epochs (i.e. the number of training iterations). Here we tune our model training with a train-test-split procedure from Chollet (2017). The initial training data is split with 20% of the data set aside for internal validation (note that this procedure does not use the data we have set aside for additional validation as detailed in Table 2). The model is then trained for a full 50 epochs. At each epoch, we save the training loss, the validation loss, the training accuracy (% correct tiles) and the validation accuracy. These can then be plotted as a function of the training epochs (Figure 5). When the validation and the training results diverge, optimal training has been reached. This final stage is key in avoiding network over-fitting which occurs when a model learns the noise in the data and loses the ability to generalise to new out-of-sample data (i.e. data not in the training set). In Figure 5, for example, we can see that the trends for training and validation diverge at 6 epochs, thus indicating the optimal training length specific to this example.



610 **Figure 5. Example of the tuning procedure used to determine the appropriate number of training epochs for CNN architectures from Chollet (2017). Here we see the divergence between training data (lines) and validation data (dashed) after 6 epochs as visible in both accuracy (right) and loss (right) data.**

615

Table 4. Multilayer Perceptron (MLP) used in phase 2 of CNN-supervised classification and as a pixel-based classifier in experiment 1.

Layer	Description
1	Dense layer, 256 nodes with L2 regularisation
2	Dropout layer, drop 50% of nodes
3	Dense layer, 128 nodes with L2 regularisation
4	Dense layer, same number of nodes as micro-classes, softmax activation to get class.

620

Experimental Design.

We conducted three experiments to address our research objectives. We first compare our approach to accepted statistical and machine learning classifiers: the maximum likelihood algorithm, the random forest algorithm, and a pixel-based multilayer perceptron. Second, we assess if a CNN-based approach such as CSC is capable of 1) simultaneously learning features for several rivers and 2) if such learned features can transfer to new rivers. We proceed by training a single CNN with data from only 5 of our 11 rivers and subsequently testing its performance on all 11 rivers. Third, we assess if training a CNN on relatively few samples (<40,000) from a single river can classify the remaining images for that same river. After experimentation, we demonstrate the GIS integration of the method by processing and displaying class maps for the orthomosaics of

630

the Ste-Marguerite river (part of the CNN training) and the Kurobe river (not part of the CNN training).

635 We begin the first experiment by creating a spatial composite image from thumbnail samples extracted from the Ste-Marguerite and Dartmouth data. The general appearance of these rivers is similar and thus allows us to assume that we do not need to use micro-classes and we therefore consider, for this experiment only, that a semantic class is identical across the whole patchwork image composed of data from 2 rivers. This image is
640 6000 x 9000 pixels in the usual RGB bands. The image is composed of 24 thumbnails of 1500 x 1500 pixels, 12 each from the St-Marguerite and Dartmouth arranged in a 4x6 grid. The associated training data for each thumbnail was carried over in order to construct a training raster, also of 6000 x 9000 pixels. In total, this raster had 33,752,194 labelled pixels and an associated 12,149 labelled tiles of 50x50 pixels. This data can now be used
645 as a testing ground for established statistical and machine learning algorithms. Maximum likelihood is arguably the most deployed classification algorithm, has served the Earth Observation community for decades (e.g. Erbek et al., 2004; Otukei and Blaschke, 2010; Strahler, 1980) , and is the most commonly available technique in classification software. Random Forest classification is a powerful ensemble method that uses random sampling
650 to produce a large number of classification trees (Belgiu and Drăguț, 2016; Pal, 2005). It is frequently deployed in remote sensing and GIS software (Table 1) and has been noted for strong performance in the remote sensing literature (e.g. Chen et al., 2017; Feng et al., 2015; Stumpf and Kerle, 2011). A multilayer perceptron (MLP), alternatively referred to as a Densely Connected Neural Network or an Artificial Neural Network, is a classic network
655 of weighted and connected nodes that can be used for regression and classification problems (Foody, 1995; Jain et al., 1996). These three established methods will therefore

be compared to our proposed deep-learning methods based on convolutional neural networks. For maximum likelihood, we used the GRASS *r.maxlik* routine as implemented in QGIS 3.4. Other algorithms were coded in Python with the libraries described above.

660 After training the algorithms, we use two validation cases. First, we validate the results by using the full validation datasets for the Ste-Marguerite and Dartmouth rivers as described in Table 2. Second, we apply the trained models to all the validation images of the remaining nine rivers: Kurobe, Kinu, Sesia, Dora di Veny, Eamont, Kingie, Pacuare, Ouelle and Kananaskis. Given that many of these rivers have patches of senescent vegetation

665 which are not present in the Ste-Marguerite and Dartmouth rivers, we have excluded the senescent vegetation class from this experiment.

In our second experiment, we aim to produce a single classifier that can transfer to all of our rivers. We therefore train the NASNet CNNs with the 184 760 labelled tiles from the

670 Ste-Marguerite, Kananaskis, Kingie, Sesia and Kinu rivers (Figure 2). After training, we conduct two separate validation tests. First, we validate the CSC outcomes with the validation images from rivers shown in Figure 2 (i.e. rivers used in training but where specific training images are not used in validation). Second, we validate the outcomes with the images from rivers shown in Figure 3 (i.e. rivers not used in training). For NASNet

675 Large, there is a clear divergence after 7 epochs of training. The case of NASNet Mobile was more ambiguous and it was determined to train up to 25 epochs where the performance seems to stop improving.

In our third experiment, we aim to test our new CSC approach in the context of current practice where data is most often acquired for a single (or few) rivers or catchments.

680 Within this workflow, data acquisition will typically result in hundreds of thousands of images. Normally, the task of classifying this data to within a reasonable accuracy can

takes months and our collective experience shows that more time is spent in manual editing of errors after a first-pass classifications are produced. For example, in Carboneau et al (2004) the first author used a set of ~2500 hyperspatial images of the Ste-Marguerite river, some of which are used here. After a first classification based on basic thresholding with Otsu's method, a month's full-time work was required to manually edit the classification mistakes and get a high quality dataset. This experiment therefore aims to assess if a deep learning approach could deliver a classification that is immediately useable and obviates the need for manual editing of classification errors. We focus on the NASNet large architecture and we will use the data for each single river to train a bespoke model specific to this river, we then validate this model against the validation images for the specific river. We repeat this for all 11 rivers.

Validation

We primarily use the F1 score as a validation metric (Burkov, 2019; Chinchor, 1992; Goodfellow et al., 2016). The F1 score, sometimes called the F-measure (Hripcsak and Rothschild, 2005), is defined as the harmonic mean of precision (P) and recall (R):

$$F1 = \frac{2 \times PR}{(P + R)} \quad (1)$$

where:

$$P = \frac{Tp}{Tp + Fp} \quad (2)$$

$$R = \frac{Tp}{Tp + Fn} \quad (3)$$

In (1), (2) and (3), P and R, the precision and recall, respectively, are defined with the concepts of true positives (Tp), false positives (Fp) and false negatives (Colquhoun, 2017). True positives are correct observations, in this case a correctly classified pixel. False

positives are observations of a factor which are incorrect, in this case a mistakenly classified pixel. Conversely, false negatives are incorrect failures to detect an observation. For example, all actual water pixels classified as vegetation are false negatives. With these quantities, precision is defined as the ratio of true positives to the sum of true positives and false positives as in equation (2). Recall is the ratio of true positives to the sum of false negatives and true positives as in equation (3) (Buckland and Gey, 1994; Burkov, 2019). The precision metric is internal to each class, it only considers correct (true positives) and incorrect (false positives) for each class. However, recall gives a measure of class confusion. The inclusion of recall therefore makes the F1 metric sensitive to class imbalance and therefore a better metric in our case than traditional accuracy (Labatut and Cherifi, 2012). In the case of very high classification qualities, accuracy and F1 are nearly identical. A perfect classification will have a F1 and accuracy scores of 100%. However for lower quality classifications, the recall parameter in the F1 score will mitigate the importance of class imbalance and values of F1 can be either higher or lower than corresponding accuracy. We strongly encourage readers to adopt this new quality metric which is in fact standard in the wider field of machine learning. In order to facilitate this transition, we provide additional information in the supporting information data where readers will find a scatter plot of F1 vs the traditional accuracy metric as well as some key results expressed as accuracy instead of F1. Additionally, we use Cohen's Kappa statistic to account for random true positives in the results (Cohen, 1960; Smeeton, 1985). The Kappa statistic ranges from -1 to 1 and should not be interpreted as a percentage of 'correct' outcomes or as a correlation. Rather, it compares the agreement between two operators, in this case the human-based validation and the machine learning classifier. The resulting measurement of agreement needs to be interpreted within established boundaries. Landis and Koch (1977) propose that Kappa values < 0 indicate

no agreement; Kappa values from 0 to 0.2 indicate slight agreement; values from 0.2 to 0.4 indicate fair agreement; values from 0.4 to 0.6 indicate moderate agreement; values from 0.6 to 0.8 indicate substantial agreement and values above 0.8 indicate almost perfect agreement. Similarly, Fleiss et al (2013) suggest that a kappa value below 0.4 indicate a poor agreement, from 0.4 to 0.75 a good agreement and above 0.75 indicate excellent agreement.

For each experiment, we calculate F1 and Kappa for each resulting classification and we compile the results to create distributions. The individual observations in these distributions are the classification metric (F1 or Kappa) for single images. In the case of the F1 score evaluation for the second experiment, we disaggregate the score for each class and can therefore produce additional distributions of F1 for each class where each observation is the classification metric for a single class in a single image. We present the results by using violin plots (Hintze and Nelson, 1998) to visualise the distributions and use the median and mean values of the distributions as summary statistics. Additionally, the supporting information document presents a large-scale validation where single values of F1 and kappa are calculated based on the aggregation of the entire set of relevant validation pixels in each experiment.

We make a careful distinction when categorising the data as in-sample or out-of-sample. Strictly speaking, machine learning practitioners define in-sample data as data that was used in training and out-of-sample data as not used in training (Chollet, 2017). It is therefore expected that in-sample data always gives strong results at the validation stage since the classifier has been trained specifically to this data. Conversely, out-of-sample data is expected to have a lower quality in validation because it has never been seen by

the classifier. We argue that this distinction is not as clear-cut in the case of our data. For the type of airborne data used here where all images from any given river were collected on the same day and with the same sensor, the resulting imagery has very similar properties across the entire image set. We therefore expect that a classifier trained on a portion of this data will perform well on the rest of the data even if it has never seen this data in training. We therefore adopt a slightly more stringent definition of in-sample and out-of-sample data. In this work, we never validate a classifier with the same data that has been used for training. Rather, we define in-sample data as image data from a river that the classifier has seen in training, but where the specific validation images have not been used in training. Out-of-sample data is therefore defined simply as data from a river never seen by the classifier in the training stage. This notion of in- and out-of-sample is crucial because the most important goal of this study is to explore the transferability potential of deep learning classifiers across multiple rivers.

770

Results

First Experiment: Classifier comparison

Figure 6 shows the outcome of the first experiment. Overall, we see that the pixel-based approaches, i.e. those that predict classification of a given pixel solely based on the radiance values of that single pixel (Maximum Likelihood, Random Forests, MLP), reach similar performances on the order of ~70%-80% F1. We also show the outcomes of both phases of the CSC process (CNN and CNN+MLP). In Figure 6, the CNN results correspond to the tiled predictions of the pre-trained CNN when re-formed as an image and validated against labelled pixels. The CNN+MLP results are therefore the final outcome of the CSC workflow where the CNN predictions become the training labels for the MLP phase. Figure 6a shows that CNN and final CNN+MLP (the final CSC result) results are markedly better than the Maximum Likelihood, Random Forest and MLP classifiers, with the CNN and CNN+MLP approaches yielding F1 scores of 92% and 95%, respectively. Overall, maximum likelihood exhibits a stronger difference in performance between the Dartmouth vs the Ste-Marguerite datasets than do the other methods. The violin plot distributions also show that the maximum likelihood classifier is generally much less reliable than other approaches, with many occurrences of classifications below 60% and some even as low as 40%. MLP and Random Forest algorithms have a low incidence of classifications below 60% and almost no instances of results below 40%. However, we note that for the Ste-Marguerite River, maximum likelihood actually outperformed the MLP and the random forest. However, the key result is the good performance of the CNN-based CSC method, with a particularly encouraging F1 score of 95%.

In figure 6b, we see the outcomes of the application of the trained classifiers obtained above to the remaining nine rivers (i.e. those not used in training). Outcomes have

degraded markedly. Maximum likelihood is strongly bimodal indicating that for some rivers, performance was good but for others, poor. Median F1 score is 52%. The pixel-based MLP and random forest algorithms had extremely variable performances with many instances of very poor performance with median F1 scores of 62% and 55%, respectively. The CNN performs somewhat better than the pixel-based approaches, but not markedly so with a median F1 score of 72%. This indicates that even the CNN tiled predictions suffered from significant error. Contrary to these results, the outcome of our novel CSC workflow (CNN+MLP) is generally encouraging; despite none of these rivers being included in the training data, the median F1 score was 89%. The senescent vegetation class was removed from this analysis because no senescent vegetation was present in the training data. We also note that the lower quartile was only 54% F1 which indicates a marked tail of poor results within this distribution.

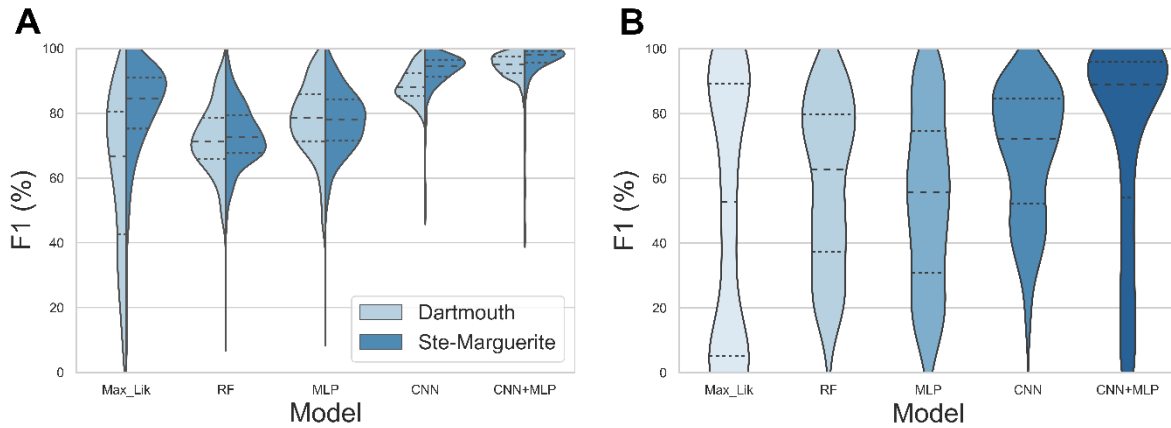


Figure 6. Results of the first experiment displayed as violin plots. Each plot is a distribution of weighted F1 scores for individual images, smoothed with a kernel density estimate. A) Results for the 2 in-sample rivers (Ste-Marguerite and Dartmouth) with vertical partition in each violin distinguishing data from each river used in the experiment. The CNN result corresponds to the first phase of the CSS process. For each violin, the number of images/samples (n) is 467. B) Results for the remaining nine out-of-sample rivers (Kurobe, Kinu, Sesia, Dore di Veny, Kingie, Eamont, Pacuare, Ouelle and Kananaskis). For each violin in B), $n=394$. In both A) and B), dotted lines give the upper and lower quartiles and the dashed line gives the median. Note that the horizontal width of these plots is scaled for maximum visibility and is not proportional to the number of samples in the data.

825 *Second Experiment: CSC Model Transferability*

The second experiment used a pre-trained CNN based on the 184,760 labelled tiles extracted from the five rivers shown in Figure 2. Table 5 and Figures 7 and 8 show the results of the second experiment. In the case of in-sample data (five rivers used in CNN training; Figure 2) and for NASNet Large (both CNN and CNN+MLP phases), we obtain
830 extremely high median (pixel weighted mean) classification accuracies of 98% (96%) and 99% (97%). In the case of NASNet Mobile (CNN and CNN+MLP phases), we obtain slightly lower but nonetheless impressive median (pixel weighted mean) values of 97% (96%) and 98% (95%) respectively. When compared to Figure 6, the larger training dataset (12K vs 184K tiles) used in the second experiment has reduced error at the CNN
835 phase thus allowing the second MLP phase to attain exceptional performance levels. The per-class disaggregation (Figures 7b and 8b) shows a similar pattern with green vegetation and water performing well but with the other classes having lower quartiles below 80% F1 (Table 5). We note that classes with poor performance are less well represented in the validation data (e.g. sediment/senescent veg/paved roads: Table 3) and
840 that this is accompanied by a degradation in performance as we move from phase 1 to phase 2 of the CSC process. However, overall, we note that once again the second phase MLP delivers an improvement on the first stage CNN (Figure 7a).

845

850

855

In-Sample Data				
Class	NASNet Large		NASNet Mobile	
	CNN	CNN+MLP	CNN	CNN+MLP
Water	97 (93)	98 (93)	96 (93)	98 (94)
Sediment	79 (69)	83 (67)	77 (68)	84 (66)
Green Vegetation	99 (98)	99 (96)	99 (96)	99 (96)
Senescent Vegetation	96 (84)	96 (79)	92 (82)	97 (80)
Paved Roads	94 (80)	93 (66)	92 (75)	93 (64)
ALL F1	98 (96)	99 (97)	97 (96)	99 (96)
ALL Kappa	0.94(0.90)	0.96(0.92)	0.93(0.89)	0.96(0.92)
Out-of-Sample Data				
Water	79 (72)	89 (76)	74 (68)	86 (72)
Sediment	68 (62)	85 (73)	67 (62)	81 (67)
Green Vegetation	84 (78)	90 (83)	83 (76)	89 (82)
Senescent Vegetation	75 (68)	80 (70)	57 (52)	77 (66)
Paved Roads	75 (64)	73 (62)	67 (63)	64 (56)
ALL F1	82 (79)	90 (83)	80 (76)	88 (80)
ALL Kappa	0.57(0.54)	0.74(0.66)	0.49(0.49)	0.71(0.62)

860

Table 5. Disaggregated results for the second experiment and for both in-sample and out-of-sample validation data. Values correspond to Median (Mean) % F1 scores. The median and mean are calculated based on each instance of a class in each image.

865

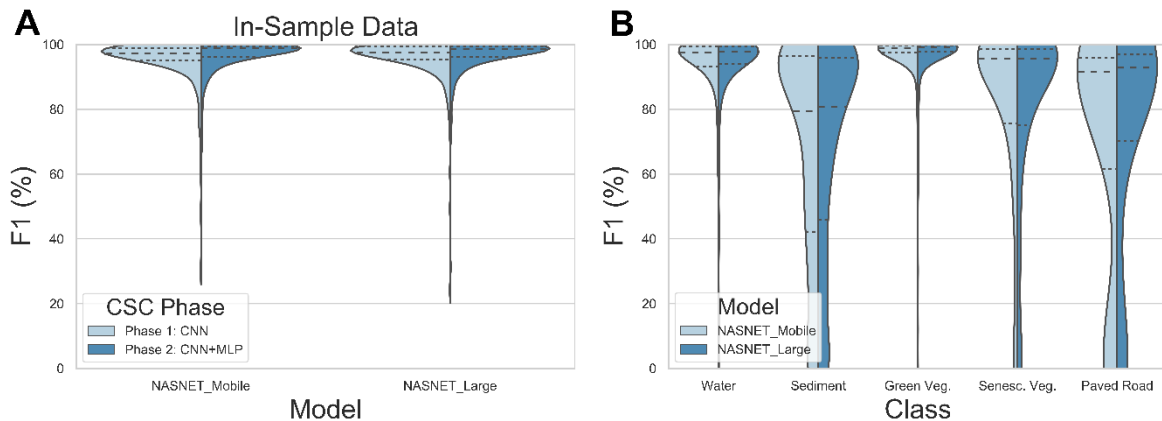


Figure 7. CSC performance for the second experiment validated with in-sample data only. A) Overall performance for each CNN model. The violins are split according to phase 1 (CNN) and phase 2 (MLP) of the CSC process. B) Final CSC performance (MLP phase 2) for the second experiment disaggregated over individual classes. For both A) and B), $n=348$. Violin plots are split according to the CNN model used. . Note that the x-axis in both plots is non-linear. The width of each violin plot is scaled for maximum visibility with each violin having the same width. Relative number of samples in each violin cannot be inferred from this figure.

880 In the case of the out-of-sample rivers in Figures 3–8, we see a degradation of performance at the initial CNN stage followed by a marked improvement at the CNN+MLP stage with respect to the in-sample data in Figure 7. In the case of NASNet Large and for the CNN and CNN+MLP phases, we obtain median (pixel weighted mean) values of 82% (79%) and 90% (83%), respectively (Figure 8A). In the case of NASNet Mobile and for the
885 CNN and CNN+MLP phases, we obtain median (pixel weighted mean) values of 80% (76%) and 88% (80%), respectively. In Table 5, we see that all classes except Paved Road have significantly improved after running an MLP on CNN outputs. Figure 8b shows an improvement in the classification of several classes with green vegetation, water and sediment achieving F1 scores above 80% in the CNN+MLP column. Furthermore, we
890 note that the lower quartile for the final MLP classification using the NASNet Large model is 81% showing that the expanded training of the CNN model has stabilised the final outcome when compared to Figure 6b where the lower quartile F1 score was 54%.

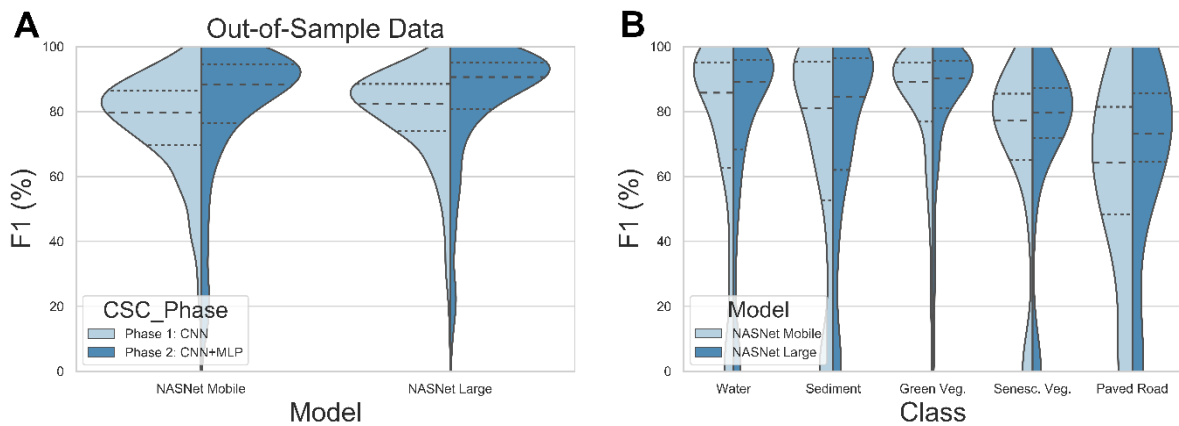


Figure 8. CSC performance for the second experiment validated with out-of-sample data only. A) Overall performance for each CNN model. The violins are split according to phase 1 (CNN) and phase 2 (MLP) of the CSC process.. B) Final CSC performance (MLP phase 2) for the third experiment disaggregated over individual classes. For both A) and B) $n=513$. Violin plots are split according to the CNN model used. Dotted lines give the quartiles. Note that the x-axis in both plots is non-linear. The width of each violin plot is scaled for maximum visibility with each violin having the same width. Relative number of samples in each violin cannot be

900

905 **inferred from this plot**

Third Experiment: Multiriver deployment

In the third experiment we examine the results when CSC is deployed to multiple rivers in a conventional workflow with training data provided for each river. Note that we do not consider the large orthomosaics used to demonstrate GIS integration. Table 6 and figure 9 show the outcomes. Since not all classes are present in all rivers, some rivers, (eg the Eamont) have 3 classes. Given that it is easier to classify an image with fewer classes, we report Cohen's kappa statistic for each river which is given as the mean kappa obtained from the kappa score for each image of every given river. Confusion matrices are available in the supporting information document (figures S2 to S21). In table 6 we see very strong performance with the weakest performance being a median classification F1 score of 95% and 93% for the rivers Dartmouth and Kananaskis, respectively. Kappa scores are generally above 0.8 with the exception of the Kanaskis river results with 0.75 for the phase 1 CNN and 0.72 phase 2 CNN+MLP. These results would only be qualified as 'good' in the interpretation of the Kappa score (Cohen, 1960; Smeeton, 1985). Across all the images, the median F1 score was 98% with a mean of 95%. In figure 9, the poorest performance for lower quartiles is 90%. When we examine specific error sources, they can be traced to specific problems. First, there is still a tendency to classify very deeply shaded sediment as water (e.g. figure 2, Kinu river). Second, very bright white water and strong sun glint can be classified as sediment (e.g. figure 3, Pacuare and Eamont river. Third, shallow water with a deep green hue and sometimes having algae can sometimes be classified as green vegetation (e.g. figure 3, kurobe river). Nevertheless, as per table 6, all mean and median values of F1 are above 90%. We note that 633 images of 861 (73.5%) were classified with an F1 score above 95%. Of these, 330 returned an F1 score of 99% (38.3%). However, figure 9 does show tails to the distributions and we note

instances of poor performance. In total, we find 10 of 861 images (1.2%) with $50\% < F1 < 80\%$ and 7 images of 861 images (0.8%) with $0\% < F1 < 50\%$. Examination of the data shows that this is caused by the misclassification of sun glint over water. Nevertheless, overall these results exceed any classification performance reported in the airborne FRS literature. Within the wider perspective of the whole Earth Observation literature, it is only deep learning methods that have reported this level of performance over a similarly wide number of samples.

Table 6. Results of CNN-supervised classification for experiment 3. Outcomes are given as median F1 [%] / mean F1 [%] / mean Kappa [-1 to 1]. The last 2 lines report median/mean for F1 and kappa. The number of validation images per river (n) is reproduced from table 2. ANOVA testing indicates that there is no correlation between F1 scores and sample size ($p=0.05$).

NASNet Large			
River	CNN	CNN+MLP	n
Dartmouth	93/92/0.83	95/93/0.85	243
Kananaskis	95/94/0.75	95/93/0.72	34
Ouelle	97/96/0.87	98/97/0.89	117
Ste-Marguerite	97/96/0.90	99/97/0.94	224
Pacuare	99/97/0.92	98/96/0.91	38
Dora diVeny	98/97/0.93	97/96/0.90	28
Sesia	98/98/0.85	99/99/0.93	21
Kinu	97/93/0.85	99/93/0.89	54
Kurobe	99/95/0.89	99/93/0.89	78
Eamont	98/96/0.88	98/96/0.91	9
Kingie	98/97/0.94	98/95/0.93	15
ALL F1	97/94	98/95	861
ALL Kappa	0.91/0.85	0.93/0.87	861

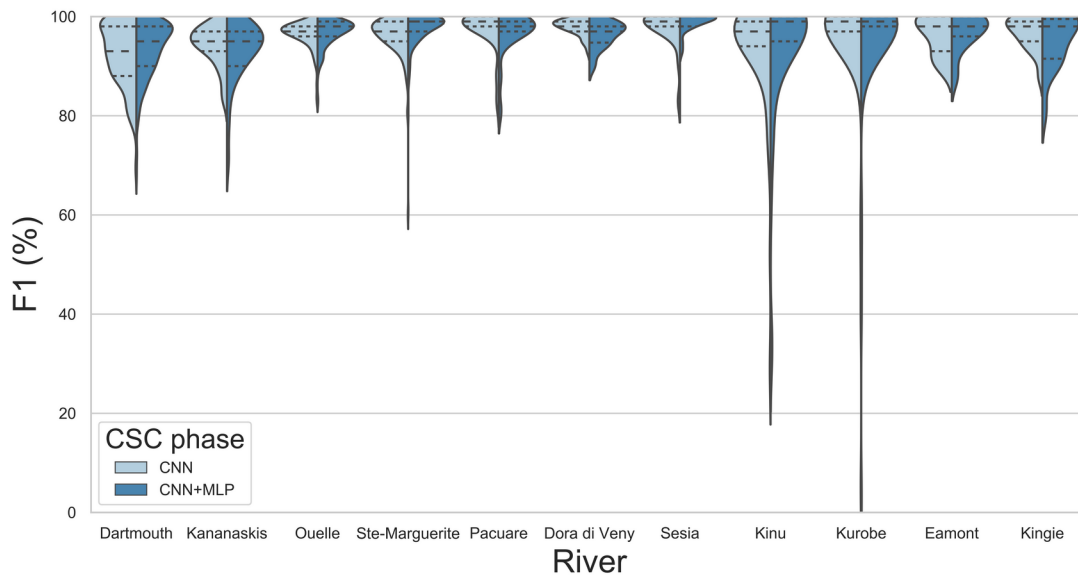


Figure 9. CSC performance for the third experiment. The violins are split according to phase 1 (CNN) and phase 2 (CNN+MLP) of the CSC process. Note that the x-axis in both plots is non-linear. The width of each violin plot is scaled for maximum visibility with each violin having the same width. Relative number of samples in each violin cannot be inferred from this figure but are given in table 6.

GIS integration

Figures 10 and 11 demonstrate GIS integration and show larger examples of mapped classification outputs. We show the original orthomosaic, the phase 1 CNN output, reformed as an image, and the final CSC classification with the phase 2 MLP. In figure 10, we show a classification for an orthomosaic of a 1km stretch of the Ste-Marguerite River that was included in CNN training (in-sample). Notably, the first phase (CNN) of classification has a significant number of errors where several patches of senescent vegetation, absent from this river reach, were falsely identified. The second stage MLP classification, using the CNN data as a training input, delivered a significant improvement,

with a final F1 score of 97%. Figure 11 follows the same pattern but we use a 1km stretch of the Kurobe river. This river was never seen by the pre-trained CNN and the F1 score is 87%. This case is a good example of the use of a pre-trained CNN in a CSC workflow to train a newly acquired orthomosaic in a fully automated fashion. The resulting accuracy is 970 unprecedented in fluvial scene classification with the major advance being the complete absence of user intervention to provide further training data. Furthermore, this was achieved with standard RGB imagery, without the need for near-infrared multispectral data.

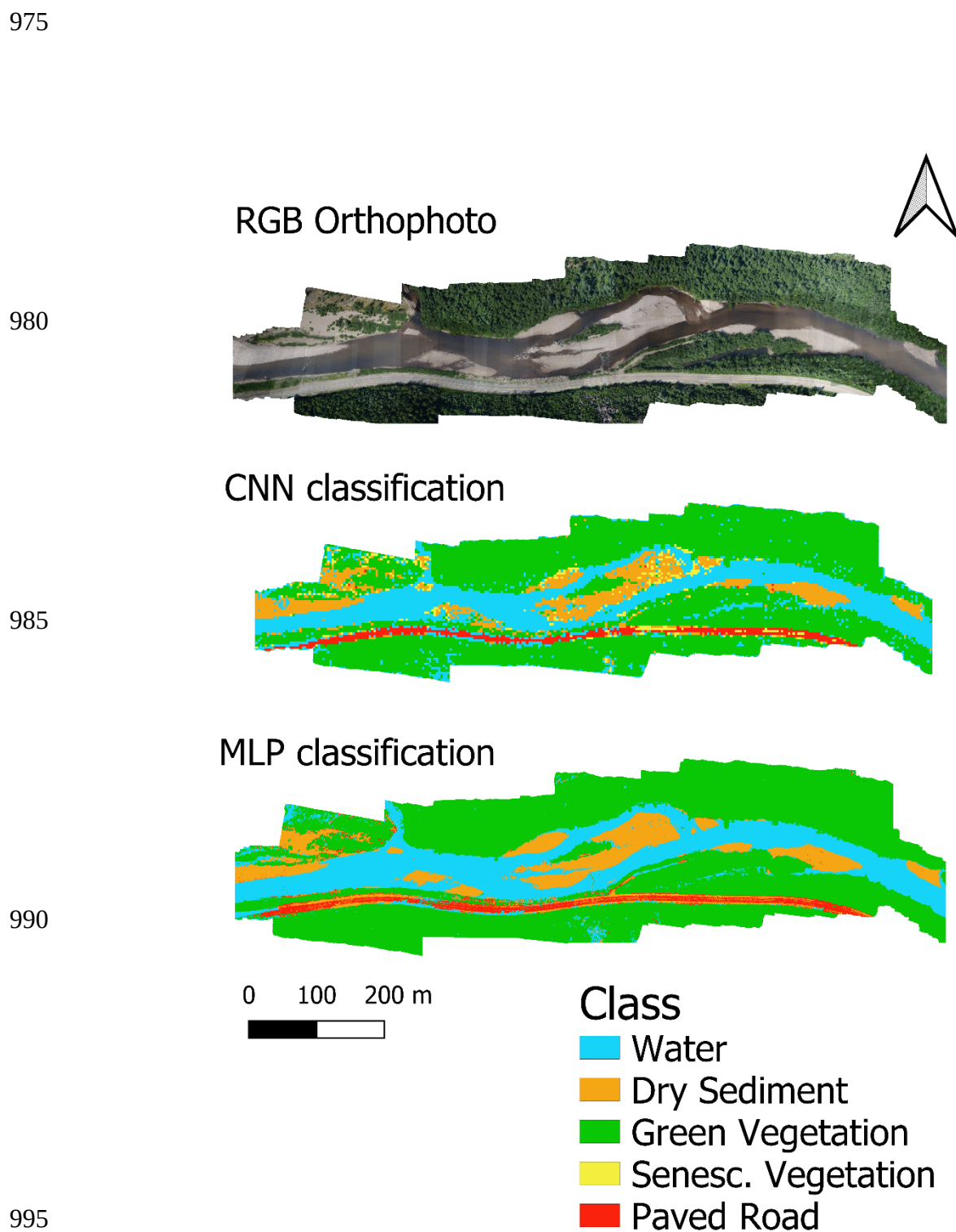
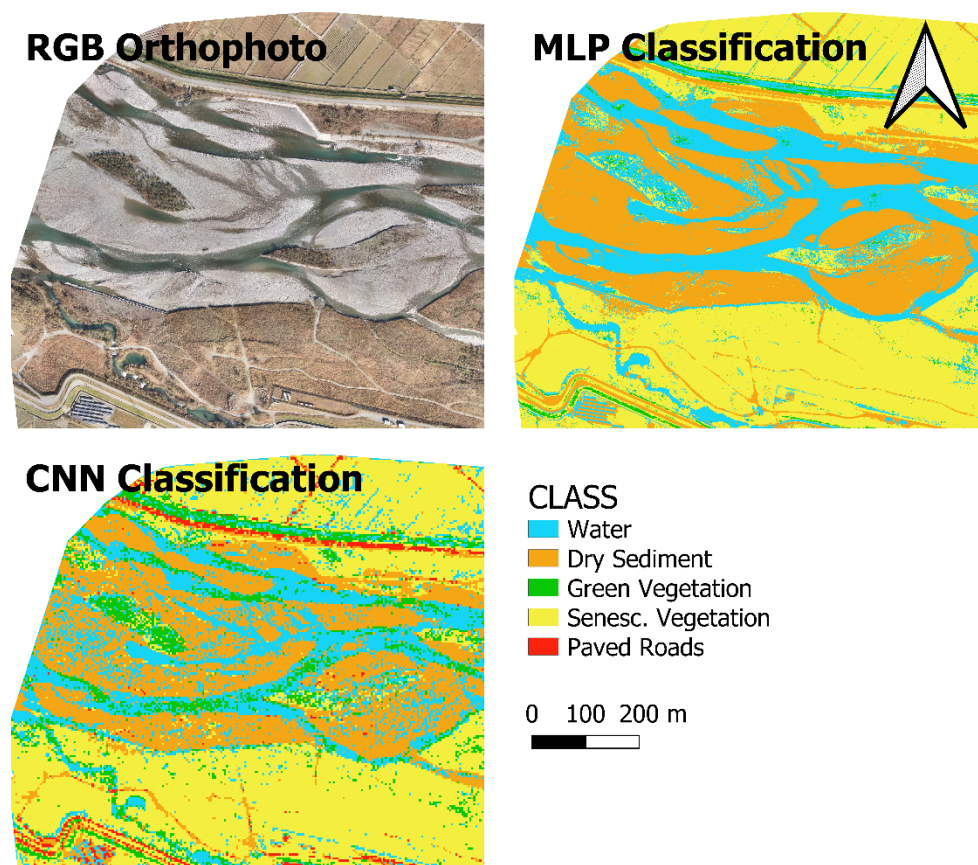


Figure 10. Mapping outputs for an orthoimage showing a 1km reach of the Ste-Marguerite at a spatial resolution of 7.5 cm. Geocoded outputs for both the CNN and MLP phases of the CSC workflow are shown. The final pixel-weighted accuracy of the MLP classification is 97% F1.

1000



1005

Figure 11. Mapping outputs for an orthoimage showing a 1km reach of the Kurobe river at a spatial resolution of 7.5 cm. Geocoded outputs for both the CNN and MLP phases of the CSC workflow are shown. Data from this river was not included in the CNN training sample. This classification output is fully automated and has not required additional training data or human-operator intervention. The final pixel-weighted accuracy of the MLP classification is 87% F1.

Discussion

Classification quality

1010 The quality results presented here substantially exceed the current state-of-the-art for fluvial scene classification. We have demonstrated that a trained deep learning classifier using our CNN-supervised classification (CSC) workflow can reach extremely high F1 scores of 99%. Our first experiment clearly shows that traditional methods do not match the performance of our deep learning approach. It also shows that with a relatively small
1015 label dataset of 12,000 tiles, our CSC approach can classify new images for the Dartmouth and Ste-Marguerite rivers to a median F1 of 95%. When applied to the remaining 9 rivers, most methods deteriorate markedly, but the final CSC result gives a median F1 of 89%. This is the first explicit demonstration within the context of airborne FRS whereby a classifier can deliver a good performance on rivers not included in the training set. The
1020 failure of the maximum likelihood, random forest and pixel-based multilayer perceptron (Figure 6b) also demonstrates that older methods cannot transfer to new rivers which illustrates how deep learning methods can reset the accepted state-of-the-art in image classification. Here we note that some of the rivers in the validation set were markedly different to those in the training set. The Kurobe and Kinu rivers in Japan share few
1025 similarities with the Ste-Marguerite and Dartmouth rivers in Quebec, Canada. The success of this experiment is entirely due to the phase 1 pre-trained CNN. The older methods are all reliant on pixel-level data only. But the CNN, trained on patches of 50x50 pixels, has learned other associated features such as texture and geometry. These learned contextual features mean that it is able to predict the class of a patch even when
1030 image brightness values are slightly different. It is therefore able to transfer well to other images (much more strongly in in-sample images). The second stage MLP fit then uses these purely image-specific brightness values to derive a classifier bespoke to a single

image, without the need for a human user to supervise the process and provide labels for each single image. The second experiment demonstrates the performance increase associated with a larger training dataset. Here we show that when 5 rivers are included in the training with a total number of tiles of 190,000, the resulting classifications are even more robust within the remaining in-sample images. Here we reach median performances of 96%-99% F1. This sets a new state-of-the-art for classification performance for hyperspatial river imagery acquired from airborne platforms. In the second part of experiment 2, we find that even when challenged with our most difficult task, the classification of six rivers never seen by the pre-trained CNN model, our best results still achieve a median F1 score of 90% (Table 5) with a lower quartile performance of 81% (figure 8a). In a specific case (Figure 11), we show that CSC can classify an orthoimage never seen by the pre-trained CNN to an F1 score of 87%. At first glance, this result might be considered equivalent to the previous state-of-the-art. However, our approach also represents a major improvement in terms of time and labour efficiency because it does not require any user intervention, user label production, or deep network training. The value of this finding is further evidenced by Figure 6b, which shows that transferring trained models, even a CNN, to river imagery not seen in training does not necessarily deliver good results. In our third experiment, our method, tested over 11 rivers, delivers an overall average of 93% F1 with 73.5% of the tested images achieving F1 scores above 95% and only 0.8% of images failing to exceed a 50% F1 score. We even note numerous instances (38.3%) of near perfect outcomes with F1 scores of 99%. We argue that this is the most readily applicable finding of our work. With the rise of drones as an affordable and easy to use airborne platform for hyperspatial image acquisitions, our method offers a step-change in the potential quality for the classification of such data at minimal time and effort. Given a day's data labelling work by a moderately skilled GIS user (~a number of pixels

equivalent to 40k training samples), our CSC method will be able to classify an entire dataset consisting of several thousands of images to extremely high ($\geq 90\%$ F1) accuracy. Indeed, we find that 73.5% of our tested imagery has a classification outcome above 95% F1 and argue that at this level of quality, no manual editing is required. For the 2% of images that yield an F1 score below 50%, the manual editing/classification work necessitated by these is a fraction of that previously necessitated by 'conventional' classification algorithms. Overall, the performance levels we report here are not matched in the airborne FRS literature. Even for recent methods using Object Based Image Analysis, Demarchi et al. (2020) report their best accuracies as 89% for the classification of meter-scale RGB imagery with the addition of a DEM layer as a 4th predictive feature and using what we define here as in-sample data for validation. However our results show that the level of detail present in hyperspatial imagery can be leveraged by deep learning and produce un-equalled classification performance.

We have also demonstrated the value and novelty of our CNN-supervised workflow. Examination of Figures 10 and 11 shows that significant errors occur when a CNN classifier is used in isolation. In contrast, the second phase (CNN+MLP) classification recovers many of these errors leading to a pixel-level classification that is more accurate than the phase 1 CNN-only classification. This effect can also be seen in Figures 6b and 8a, where the CNN+MLP violins plots show improved performance with respect to the CNN alone. Overall, our results show that deep learning methods have greatly outperformed statistical and machine learning methods and should now be adopted in airborne FRS as a standard classification tool. In order to facilitate adoption by other users, all the methods here are based on open source code available on GitHub and implemented using PyQGIS scripts to deliver mapping capabilities via QGIS.

Comparison to fluvial image classification 'state-of-the-art'

1085 We find that our results compare favourably to similar recent works leveraging deep learning techniques. Casado et al. (2015) use a pixel-based MLP classifier on a short river reach with an accuracy of 81%. This is comparable to results in Figure 6a. However, pixel-based classifiers have limited potential and our results, along with those of Buscombe and Ritchie (2018), show that convolutional neural networks are the way
1090 forward. Buscombe and Ritchie (2018) apply the DeepLab method (Chen et al., 2018) and present a similar two-stage workflow to CSC where the first phase of CNN classification is followed by a pixel-based classification based on conditional random fields. They report results similar to ours with mean F1 scores ranging from 88% to 98%. Detailed examination shows a pattern similar to our results where the quality statistics increase with
1095 greater data aggregation. When disaggregated, Buscombe and Ritchie (2018) find some poor results as low as 30% mean F1. Interestingly, their data do not show that the second stage of pixel-level classification, performed with conditional random fields, can improve on the performance of the phase 1 CNN. However, this might be because the authors have not attempted to highlight this behaviour and/or that they have more severe class
1100 imbalance problems. We argue that our approach goes beyond that of Buscombe and Ritchie (2018), by achieving both a) higher classification accuracies across datasets of substantially increased size; and b) demonstrating the viability and transferability of our approach across several hundreds of rivers from a range of geographically-diverse river locations. In another example of a 'chained' classification approach, Zhang et al. (2018)
1105 also combine CNNs and MLPs to perform pixel-level classification, albeit with a different workflow, in an urban/semi-urban context and with a considerably smaller dataset of

approximately 11 million pixels. Similar to our results, these authors report accuracies of 74% to 95%. W

1110 We further note that the 90%-99% F1 scores reported here are slightly better than the hyperspectral fluvial scene classification results reported by Marcus et al. (2003). We therefore argue that our results, supported by those of Buscombe and Ritchie (2018) and Zhang et al. (2018) indicate that available deep learning workflows are now capable of obviating the use of multi- and hyperspectral sensors for image classification. While these
1115 sensors retain a crucial function in advanced applications requiring airborne imaging spectroscopy capabilities (e.g. Candiago et al., 2015; Pölönen et al., 2013; Vanegas et al., 2018), their extra cost is no longer justified in any application where the final objective of image acquisition is land-cover classification of the scale described within this work. Our findings could have a significant impact on the drone industry, where we note intense
1120 commercial pressure to expand the market for multi- and hyperspectral sensors. We argue that the scientific rationale for this expansion needs re-examination.

Implications of findings for airborne fluvial remote sensing science and practise

Our results suggest an avenue for future research allowing for the inclusion of deep
1125 learning in GIS software. In Table 1, we show that at present, the inclusion of deep learning tools within GIS packages is embryonic, and indeed largely absent from open source software options. We argue that training data availability, and associated processing power requirements, pose a significant access barrier that may explain this situation. For most users, the task of image classification remains focused on a relatively
1130 small volume of data (e.g. images from a specific river reach). Therefore, in most cases, the required volume of data needed to train a deep network from scratch is not available.

We have demonstrated that the features developed by a pre-trained CNN can transfer to rivers not seen at the training stage. The accuracy of CNN predictions does decrease on transfer to unseen rivers, but in this case we have shown that the use of a chained MLP
1135 pixel-level classifier can recover some of these errors and deliver state-of-the-art classification performance (Figures 6b, 8a and 11). We therefore envisage a workflow where a classification routine embedded in a GIS could use orthoimage metadata to select and load a pre-trained CNN according to a proximity criteria (e.g. space and season). The software could then execute CNN-supervised classification and deliver a truly
1140 automated semantic classification with identified land-cover types. Optionally, users that require performance at the 95% level could add a limited selection of training areas and use transfer learning to retrain a river-specific CNN and adapt it to their specific imagery with relatively little expenditure in personnel time. Both these scenarios could function with modest processing power; throughout this work we used laptops with single processors
1145 and single, mid-range, GPUs. However, the main challenge to this vision would be the assembly of the required banks of pre-trained CNNs. Despite the fact that hyperspatial resolution aerial imagery is now available from most environments on the Earth, thanks to an explosion in the use of drones, there is still no global database of such imagery.

1150 In addition to these highly encouraging results regarding the classification of airborne FRS data, CNNs also hold a great deal of promise for addressing fundamental questions in the river sciences. For example, one interesting perspective is the possibility of using a deep CNN as an objective tool for investigating ontological issues in river morphology cataloguing. Considerable efforts have been deployed to categorise fluvial forms in a way that is both scientifically accurate and useable in
1155 a management context (Brierley et al., 2013; Brierley and Fryirs, 2000; Fryirs and Brierley, 2018; Gurnell et al., 2016). Most of these efforts rely on a mix of knowledge from fluvial geomorphology and other related sciences and they often rely on visual image interpretation, with a very high level

of expert knowledge, in order to assign their respective categories and nomenclatures to fluvial form (e.g. Fryirs and Brierley, 2018). However, in the case of surface flow features, Woodget et al. 1160 (2016) have shown that physical characteristics attributed via visual identification can suffer from ontology issues which lead to a questioning of the intrinsic existence of certain natural river features (when categorised through a conceptual process). We therefore argue that CNN-based feature classification approaches could be used to clarify the ontology of fluvial forms and serve as a testable benchmark, a 'reality check' of sorts, applied to the ontology of human-conceived 1165 features. The approach in this case would be to re-orient the classification system towards an explicit labelling of fluvial forms (point bars, braided channels, etc). If, after training, CNN-predicted labelling of these forms in validation imagery agrees with human expert knowledge, then this confirms the ontology of the given fluvial structure and the CNN can then be further used as an objective method for wider scale deployment of a given fluvial classification scheme. Such an 1170 approach would be required to robustly make the subtle transition from fluvial land-cover classification, as done in this work, to fluvial habitats (i.e. land-use by flora and fauna). Such work could make fundamental contributions to our understanding of fluvial forms that go beyond the functional requirement to classify imagery and make objective cataloguing of fluvial habitats a practical reality. However, this idea does have important technical implications. For example, if the 1175 training data labels fluvial forms, then the tiling procedure must move away having tiles 100% occupied by a single, pure, class label (as seen in this paper). For example, if we seek to train a CNN to identify point bars, then suitable labelled tiles must have the entire bar AND a portion of surrounding water. This is therefore somewhat similar to the classic case of CNN image identification where a photograph of a subject must be identified and thus the image tile contains 1180 pixels that are not semantically part of the subject to be identified. However, in the case of natural forms, issues of scalar and rotational invariance must also be considered. Fluvial forms can occur in any orientation and can vary in size by orders of magnitude. Whilst there is a body of work reporting approaches to transform invariance in the context of deep learning (Cabrera-Vives et al., 2017; Cheng et al., 2019; Dieleman et al., 2015; Srivastava and Grill-Spector, 2018), this work 1185 remains closer to the research frontier and more challenging to apply.

Method Limitations: Class imbalance and hyperparameter tuning

Class imbalance is a problem arising when training data has a large disparity in the number of samples in each class. It is the focus of significant research both in pure machine learning (e.g. Buda et al., 2018; Krawczyk, 2016; Lemaitre et al., 2016) and, to a lesser extent, Earth observation (Kampffmeyer et al., 2016; Stumpf and Kerle, 2011). This effect has an impact on our results. As visible in Table 3, one of the near-impossibilities of data preparation was to ensure equal class representation in both the training and validation data across all classes. Typical airborne remote sensing images of fluvial scenes are dominated by vegetated areas and the water. Sediment might be prominent in certain rivers but less so in others. Some images might have large sediment bars, while others only have small patches of exposed sediment. There also might be man-made features in the imagery. Ultimately, having an engineered balance, in terms of pixel numbers, for all classes is not possible unless we greatly under-sample all the better-represented classes to unacceptable levels. At a smaller scale, we observe that in cases where the phase 1 CNN predictions have a small minority in a single class, this class can be eliminated by the MLP if the training achieves minimal loss simply by predicting that a class is absent in an image. A good example is Figure 11, where we see that the paved roads class, occupying a very small percentage of pixels, has been eliminated and classed as sediment in the final MLP classification. Similarly, vegetation patches in this image, again with a small surface coverage in the image, have often been confused with water. In an attempt to address this problem, we investigated mitigation methods for class imbalance (Batista et al., 2004; Chawla et al., 2002; Lemaitre et al., 2016). We tested the Synthetic Minority Oversampling TEchnique (SMOTE). The SMOTE technique works by creating new samples of synthetic data to strengthen the minority sets in training data.

Specifically, it interpolates between inliers and outliers. This strengthens the signal of smaller samples and prevents the classifier from reaching a minimal loss solution by totally ignoring the minority class. However, in our case, we found that the application of SMOTE severely degraded performance. By interpolating between inliers and outliers, the SMOTE method amplified the erroneous CNN predictions beyond the point where the MLP predictions could mitigate against them. Consequently, we find that our workflow of CNN-supervised classification is most suited to applications where the major land-cover types need to be accurately classified and quantified. For applications where smaller features in the landscape need to be identified, we would recommend alternative approaches geared towards feature recognition as opposed to semantic classification and using a CNN to identify these small-scale local features.

One of the most problematic aspects of work such as that presented here is the very high number of CNN parameters and design decisions that we did not investigate but undoubtedly influenced our results. While we have made efforts to provide some basis for parameter selection (e.g. the tuning procedure for the NASNet architectures), it was not computationally possible to conduct a deep parameter space investigation through brute-force modelling; even Monte-Carlo approaches of random sampling within the parameter space carried an overly large computational overhead. We made efforts to justify parameter choices, but clear advice regarding hyperparameter tuning for deep neural networks is not always readily available and new users are often left with a bewildering number of choices to test. In this case, we faced several choices. At the outset, the use of a transfer learning approach requires the user to fix the weights on certain deeper layers in the CNN architecture. With a network architecture as large as NASNet, the choice of layers to fix was based on limited trial and error. Our results are satisfactory, but we

recognise that an alternative structure of fixed/trainable parameters might deliver improvements. Another issue is the size of tile to use. The selection of tile size must allow the training design to deliver a large number of labelled images. There is a trade-off between the smaller sizes/larger numbers and the information content of each tile.

1240 Buscombe and Ritchie (2018) use a tiles size of 75x75, but here we found that 50x50 gave better results. Overall, we made an effort to minimise tunable parameters in this work but we recognize that the work had a significant number of parameters chosen and tuned solely based on experience and/or minimal preliminary experiments. Exploring these parameters quantitatively might clarify small details about the overall process but at

1245 significant cost in terms of computation. We therefore advocate the use of optimisation approaches (e.g. Zheng and Wang, 1996) to identify parameter combinations that yield further improvements on our results. However, crucially, we argue that while our results might be improved upon, this does not change or invalidate our findings, namely, that the application of deep learning methods such as those outlined in this paper have delivered

1250 state-of-the-art results in hyperspatial fluvial scene classification.

1255

Conclusion

This paper uses a state-of-the-art dataset to demonstrate that deep learning methods are
1260 now ready for a wider uptake by the airborne fluvial remote sensing community,
transforming the fundamental task of supervised classification. We have shown that
replacing the conventional classifiers (eg. maximum likelihood) with deep convolutional
neural networks can substantially increase classification performance and set a new
benchmark for expected performance in RGB fluvial scene classification using a
1265 supervised workflow. With CNN-Supervised Classification, users proficient in GIS now
only need to manually label 4-8 RGB images of 12-20 Mpix in order to generate the
training data (~ 37k tiles) required to classify an entire river with hundreds or even
thousands of images to a very high standard ($F1 > 95\%$) with training data that can
manually be generated in less than 1 person/day and without the need for costly multi- or
1270 hyperspectral sensors. Finally, our results show that an advanced convolutional network
architecture such as NASNet can effectively learn a visual classification scheme for fluvial
scenes that can transfer to other rivers never seen in training. This shows a way forward
where large pre-trained CNN might be capable of classifying rivers on regional/national
scales thus truly minimising the need for human supervision. However, such work will
1275 require a coordinated effort in order to pool, organise and label the large volume of
hyperspatial river imagery that already exists but is scattered in the community. Indeed, the
wider uptake of deep learning by the airborne fluvial remote sensing community is now
somewhat dependent on improving the level of cooperation and coordination among
scientists working with hyperspatial resolution airborne imagery in order to compile and
1280 generate the so-called Big Data that drives the training of deep neural networks.

1285 **Code and data access**

Core Python scripts and usage instructions for CNN-supervised classification are available from the following GitHub repository: <https://github.com/geojames/CNN-Supervised-Classification> and can be cited as Carbonneau and Dietrich (2020). All the image and label data used in this work is also available for download from [this](#) institutional repository
1290 and can be cited as Carbonneau et al. (2019).

Acknowledgements

The authors would like to thank several funding bodies who have supported image acquisition. Images of the Ste-Marguerite and Dartmouth rivers were funded by the
1295 GEOSALAR project, part of the GEOIDE network of centres of excellence. The authors thank Professor Normand Bergeron for use of the Ouelle and Kananaskis river imagery; these data were collected as part of the NSERC/CRSNG Collaborative Research and Development Grant CRDJ 379745-08 in partnership with the Ouranos consortium on regional climatology and adaptation to climate change and also as part of the
1300 NSERC/CRSNG HydroNet Strategic Network Grant. Images of the Eamont, Sesia and Kingie rivers were funded by the AMBER project, grant number 689682, part of the EU Horizon 2020 program. The images of the Dora di Veny river were acquired thanks to support from Dr Catriona Fyffe, the University of Worcester and the British Society for Geomorphology. The images of the Kinu and Kurobe rivers were funded by the KAKENHI
1305 program of the Japanese Society for the Promotion of Science, grant number JP16H04422. The authors would like to thank Dr Pollyanna Lind for the images of the Pacuare river, funded by the National Science Foundation, the Tokyo foundation and the University of Oregon. Finally, we would like to thank three anonymous reviewers for

thorough and constructive comments that have greatly improved the clarity and impact of
1310 this work.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X., 2016. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. arXiv:1603.04467v2.
- Allen, G.H., Pavelsky, T.M., 2018. Global extent of rivers and streams. *Science* 361, 585–588. <https://doi.org/10.1126/science.aat0636>
- Arnell, N.W., Gosling, S.N., 2016. The impacts of climate change on river flood risk at the global scale. *Climatic Change* 134, 387–401. <https://doi.org/10.1007/s10584-014-1084-5>
- Ashmore, P., Sauks, E., 2006. Prediction of discharge from water surface width in a braided river with implications for at-a-station hydraulic geometry. *Water Resources Research* 42. <https://doi.org/10.1029/2005WR003993>
- Bagheri, O., Ghodsian, M., Saadatseresht, M., 2015. Reach scale application of UAV+SFM method in shallow rivers hyperspatial bathymetry, in: ISPRS - International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences. Presented at the WG I/4
 International Conference on Sensors & Models in Remote Sensing & Photogrammetry - 23–25 November 2015, Kish Island, Iran, Copernicus GmbH, pp. 77–81. <https://doi.org/10.5194/isprsarchives-XL-1-W5-77-2015>
- Barré, P., Stöver, B.C., Müller, K.F., Steinhage, V., 2017. LeafNet: A computer vision system for automatic plant species identification. *Ecological Informatics* 40, 50–56. <https://doi.org/10.1016/j.ecoinf.2017.05.005>
- Batista, G.E.A.P.A., Prati, R.C., Monard, M.C., 2004. A Study of the Behavior of Several Methods for Balancing Machine Learning Training Data. *SIGKDD Explor. Newsl.* 6, 20–29. <https://doi.org/10.1145/1007730.1007735>
- Belgiu, M., Drăguț, L., 2016. Random forest in remote sensing: A review of applications and future directions. *ISPRS Journal of Photogrammetry and Remote Sensing* 114, 24–31. <https://doi.org/10.1016/j.isprsjprs.2016.01.011>
- Bishop, C., 2006. *Pattern Recognition and Machine Learning*, Information Science and Statistics. Springer-Verlag, New York.
- Bjerklie, D.M., Lawrence Dingman, S., Vorosmarty, C.J., Bolster, C.H., Congalton, R.G., 2003. Evaluating the potential for measuring river discharge from space. *Journal of Hydrology* 278, 17–38. [https://doi.org/10.1016/S0022-1694\(03\)00129-X](https://doi.org/10.1016/S0022-1694(03)00129-X)
- Black, M., Carbonneau, P., Church, M., Warburton, J., 2014. Mapping sub-pixel fluvial grain sizes with hyperspatial imagery. *Sedimentology* 61, 691–711. <https://doi.org/10.1111/sed.12072>
- Boruah, S., Gilvear, D., Hunter, P., Sharma, N., 2008. Quantifying channel planform and physical habitat dynamics on a large braided river using satellite data - The Brahmaputra, India. *River Research and Applications* 24, 650–660. <http://dx.doi.org/10.1002/rra.1132>
- Brierley, G., Fryirs, K., Cullum, C., Tadaki, M., Huang, H.Q., Blue, B., 2013. Reading the landscape: Integrating the theory and practice of geomorphology to develop place-based understandings of river systems. *Progress in Physical Geography: Earth and Environment* 37, 601–621. <https://doi.org/10.1177/0309133313490007>

- Brierley, G.J., Fryirs, K., 2000. River Styles, a Geomorphic Approach to Catchment Characterization: Implications for River Rehabilitation in Bega Catchment, New South Wales, Australia. *Environmental Management* 25, 661–679. <https://doi.org/10.1007/s002670010052>
- Brigante, R., Cencetti, C., Rosa, P.D., Fredduzzi, A., Radicioni, F., Stoppini, A., 2017. Use of aerial multispectral images for spatial analysis of flooded riverbed-alluvial plain systems: the case study of the Paglia River (central Italy). *Geomatics, Natural Hazards and Risk* 8, 1126–1143. <https://doi.org/10.1080/19475705.2017.1300607>
- Buda, M., Maki, A., Mazurowski, M.A., 2018. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks* 106, 249–259. <https://doi.org/10.1016/j.neunet.2018.07.011>
- Burkov, A., 2019. *The Hundred-Page Machine Learning Book* by Andriy Burkov. Self-Published.
- Buscombe, D., Ritchie, A., 2018. Landscape Classification with Deep Neural Networks. *Geosciences* 8, 244. <https://doi.org/10.3390/geosciences8070244>
- Butler, J.B., Lane, S.N., Chandler, J.H., 2001. Automated extraction of grain-size data from gravel surfaces using digital image processing. *Journal of Hydraulic Research* 39, 519–529. <https://doi.org/10.1080/00221686.2001.9628276>
- Cabrera-Vives, G., Reyes, I., Förster, F., Estévez, P.A., Maureira, J.-C., 2017. Deep-HiTS: Rotation Invariant Convolutional Neural Network for Transient Detection. *ApJ* 836, 97. <https://doi.org/10.3847/1538-4357/836/1/97>
- Candiago, S., Remondino, F., De Giglio, M., Dubbini, M., Gattelli, M., 2015. Evaluating Multispectral Images and Vegetation Indices for Precision Farming Applications from UAV Images. *Remote Sensing* 7, 4026–4047. <https://doi.org/10.3390/rs70404026>
- Carbonneau, P., Fonstad, M.A., Marcus, W.A., Dugdale, S.J., 2012. Making riverscapes real. *Geomorphology* 137, 74–86.
- Carbonneau, P.E., Dietrich, J.T., 2020. CNN-Supervised-Classification. Zenodo. <https://doi.org/10.5281/zenodo.3928808>
- Carbonneau, P.E., Dugdale, S.J., Miyamoto, H., Woodget, A.S., Fonstad, M.A., Dietrich, J.T., Breckon, T.P., 2019. Self-Supervised Image Classification [dataset].
- Carbonneau, P.E., Lane, S.N., Bergeron, N., 2006. Feature based image processing methods applied to bathymetric measurements from airborne remote sensing in fluvial environments. *Earth Surf. Process. Landforms* 31, 1413–1423. <https://doi.org/10.1002/esp.1341>
- Carbonneau, P.E., Lane, S.N., Bergeron, N.E., 2004. Catchment-scale mapping of surface grain size in gravel bed rivers using airborne digital imagery. *Water Resour. Res.* 40, W07202. <https://doi.org/10.1029/2003WR002759>
- Carbonneau, P.E., Piégay, H., 2012a. *Fluvial Remote Sensing for Science and Management*. John Wiley & Sons.
- Carbonneau, P.E., Piégay, H., 2012b. Introduction: The Growing Use of Imagery in Fundamental and Applied River Sciences, in: *Fluvial Remote Sensing for Science and Management*. John Wiley & Sons, Ltd, pp. 1–18. <https://doi.org/10.1002/9781119940791.ch1>
- Carrivick, J.L., Smith, M.W., 2019. Fluvial and aquatic applications of Structure from Motion photogrammetry and unmanned aerial vehicle/drone technology. *WIREs Water* 6, e1328. <https://doi.org/10.1002/wat2.1328>
- Carrizo, S.F., Jähnig, S.C., Bremerich, V., Freyhof, J., Harrison, I., He, F., Langhans, S.D., Tockner, K., Zarfl, C., Darwall, W., 2017. Freshwater Megafauna: Flagships for

- Freshwater Biodiversity under Threat. *BioScience* 67, 919–927. <https://doi.org/10.1093/biosci/bix099>
- Casado, M.R., Gonzalez, R.B., Kriechbaumer, T., Veal, A., 2015. Automated Identification of River Hydromorphological Features Using UAV High Resolution Aerial Imagery. *Sensors* 15, 27969–27989. <https://doi.org/10.3390/s151127969>
- Chandler, J., Ashmore, P., Paola, C., Gooch, M., Varkaris, F., 2002. Monitoring River-Channel Change Using Terrestrial Oblique Digital Imagery and Automated Digital Photogrammetry. *Annals of the Association of American Geographers* 92, 631–644. <https://doi.org/10.1111/1467-8306.00308>
- Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research* 16, 321–357. <https://doi.org/10.1613/jair.953>
- Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L., 2018. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 834–848. <https://doi.org/10.1109/TPAMI.2017.2699184>
- Chen, W., Xie, X., Wang, J., Pradhan, B., Hong, H., Bui, D.T., Duan, Z., Ma, J., 2017. A comparative study of logistic model tree, random forest, and classification and regression tree models for spatial prediction of landslide susceptibility. *CATENA* 151, 147–160. <https://doi.org/10.1016/j.catena.2016.11.032>
- Chen, Y., Ming, D., Lv, X., 2019. Superpixel based land cover classification of VHR satellite image combining multi-scale CNN and scale parameter estimation. *Earth Science Informatics*. <https://doi.org/10.1007/s12145-019-00383-2>
- Cheng, G., Han, J., Zhou, P., Xu, D., 2019. Learning Rotation-Invariant and Fisher Discriminative Convolutional Neural Networks for Object Detection. *IEEE Transactions on Image Processing* 28, 265–278. <https://doi.org/10.1109/TIP.2018.2867198>
- Chinchor, N., 1992. Muc-4 evaluation metrics, in: *In Proceedings of the Fourth Message Understanding Conference*. pp. 22–29.
- Cohen, J., 1960. A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement* 20, 37–46. <https://doi.org/10.1177/001316446002000104>
- Colquhoun, D., 2017. The reproducibility of research and the misinterpretation of p-values. *R Soc Open Sci* 4. <https://doi.org/10.1098/rsos.171085>
- Daigle, A., Bérubé, F., Bergeron, N., Matte, P., 2013. A methodology based on Particle image velocimetry for river ice velocity measurement. *Cold Regions Science and Technology* 89, 36–47. <https://doi.org/10.1016/j.coldregions.2013.01.006>
- Debats, S.R., Luo, D., Estes, L.D., Fuchs, T.J., Caylor, K.K., 2016. A generalized computer vision approach to mapping crop fields in heterogeneous agricultural landscapes. *Remote Sensing of Environment* 179, 210–221. <https://doi.org/10.1016/j.rse.2016.03.010>
- Demarchi, L., Bizzi, S., Piégay, H., 2017. Regional hydromorphological characterization with continuous and automated remote sensing analysis based on VHR imagery and low-resolution LiDAR data. *Earth Surface Processes and Landforms* 42, 531–551. <https://doi.org/10.1002/esp.4092>
- Demarchi, L., Bizzi, S., Piégay, H., 2016. Hierarchical Object-Based Mapping of Riverscape Units and in-Stream Mesohabitats Using LiDAR and VHR Imagery. *Remote Sensing* 8, 97. <https://doi.org/10.3390/rs8020097>

- Demarchi, L., van de Bund, W., Pistocchi, A., 2020. Object-Based Ensemble Learning for Pan-European Riverscape Units Mapping Based on Copernicus VHR and EU-DEM Data Fusion. *Remote Sensing* 12, 1222. <https://doi.org/10.3390/rs12071222>
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. ImageNet: A large-scale hierarchical image database. 2009 IEEE Conference on Computer Vision and Pattern Recognition. <https://doi.org/10.1109/CVPRW.2009.5206848>
- Dieleman, S., Willett, K.W., Dambre, J., 2015. Rotation-invariant convolutional neural networks for galaxy morphology prediction. *Mon Not R Astron Soc* 450, 1441–1459. <https://doi.org/10.1093/mnras/stv632>
- Dietrich, J.T., 2016. Riverscape mapping with helicopter-based Structure-from-Motion photogrammetry. *Geomorphology, The Natural and Human Structuring of Rivers and other Geomorphic Systems: A Special Issue in Honor of William L. Graf* 252, 144–157. <https://doi.org/10.1016/j.geomorph.2015.05.008>
- Downing, J.A., Cole, J.J., Duarte, C.M., Middelburg, J.J., Melack, J.M., Prairie, Y.T., Kortelainen, P., Striegl, R.G., McDowell, W.H., Tranvik, L.J., 2012. Global abundance and size distribution of streams and rivers. *Inland Waters* 2, 229–236. <https://doi.org/10.5268/IW-2.4.502>
- Dugdale, S.J., Malcolm, I.A., Hannah, D.M., 2019. Drone-based Structure-from-Motion provides accurate forest canopy data to assess shading effects in river temperature models. *Science of The Total Environment* 678, 326–340. <https://doi.org/10.1016/j.scitotenv.2019.04.229>
- Durand, M., Gleason, C.J., Garambois, P.A., Bjerklie, D., Smith, L.C., Roux, H., Rodriguez, E., Bates, P.D., Pavelsky, T.M., Monnier, J., Chen, X., Di Baldassarre, G., Fiset, J.-M., Flipo, N., Frasson, R.P. d. M., Fulton, J., Goutal, N., Hossain, F., Humphries, E., Minear, J.T., Mukolwe, M.M., Neal, J.C., Ricci, S., Sanders, B.F., Schumann, G., Schubert, J.E., Vilmin, L., 2016. An intercomparison of remote sensing river discharge estimation algorithms from measurements of river height, width, and slope. *Water Resources Research* 52, 4527–4549. <https://doi.org/10.1002/2015WR018434>
- Erbek, F.S., Özkan, C., Taberner, M., 2004. Comparison of maximum likelihood classification method with supervised artificial neural network algorithms for land use activities. *International Journal of Remote Sensing* 25, 1733–1748. <https://doi.org/10.1080/0143116031000150077>
- Fausch, K.D., Torgersen, C.E., Baxter, C.V., Li, H.W., 2002. Landscapes to Riverscapes: Bridging the Gap between Research and Conservation of Stream Fishes A Continuous View of the River is Needed to Understand How Processes Interacting among Scales Set the Context for Stream Fishes and Their Habitat. *BioScience* 52, 483–498. [https://doi.org/10.1641/0006-3568\(2002\)052\[0483:LTRBTG\]2.0.CO;2](https://doi.org/10.1641/0006-3568(2002)052[0483:LTRBTG]2.0.CO;2)
- Feng, Q., Liu, J., Gong, J., 2015. UAV Remote Sensing for Urban Vegetation Mapping Using Random Forest and Texture Analysis. *Remote Sensing* 7, 1074–1094. <https://doi.org/10.3390/rs70101074>
- Feng, R., Wang, L., Zhong, Y., 2018. Least Angle Regression-Based Constrained Sparse Unmixing of Hyperspectral Remote Sensing Imagery. *Remote Sensing* 10, 1546. <https://doi.org/10.3390/rs10101546>
- Fleiss, J.L., Levin, B., Paik, M.C., 2013. *Statistical Methods for Rates and Proportions*. John Wiley & Sons.
- Foody, G.M., 1995. Land cover classification by an artificial neural network with ancillary information. *International Journal of Geographical Information Systems* 9, 527–542. <https://doi.org/10.1080/02693799508902054>

- Foody, G.M., Ling, F., Boyd, D.S., Li, X., Wardlaw, J., 2019. Earth Observation and Machine Learning to Meet Sustainable Development Goal 8.7: Mapping Sites Associated with Slavery from Space. *Remote Sensing* 11, 266.
- Fryirs, K.A., Brierley, G.J., 2018. What's in a name? A naming convention for geomorphic river types using the River Styles Framework. *PLOS ONE* 13, e0201909. <https://doi.org/10.1371/journal.pone.0201909>
- Ghaffarian, H., Piégay, H., Lopez, D., Mignot, E., MacVicar, B.J., Antonio, A., Riviere, N., 2020. Video-monitoring of wood discharge: first inter-basin comparison and recommendations to install cameras. *Earth Surface Processes and Landforms*.
- Gilvear, D.J., Sutherland, P., Higgins, T., 2008. An assessment of the use of remote sensing to map habitat features important to sustaining lamprey populations. *Aquatic Conservation: Marine and Freshwater Ecosystems* 18, 807–818. <https://doi.org/10.1002/aqc.876>
- Gleason, C.J., Smith, L.C., 2014. Toward global mapping of river discharge using satellite images and at-many-stations hydraulic geometry. *PNAS* 111, 4788–4791. <https://doi.org/10.1073/pnas.1317606111>
- Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. MIT Press.
- Guo, W., Yang, W., Zhang, H., Hua, G., 2018. Geospatial Object Detection in High Resolution Satellite Images Based on Multi-Scale Convolutional Neural Network. *Remote Sensing* 10. <https://doi.org/10.3390/rs10010131>
- Gurnell, A.M., Rinaldi, M., Belletti, B., Bizzi, S., Blamauer, B., Braca, G., Buijse, A.D., Bussettini, M., Camenen, B., Comiti, F., Demarchi, L., García de Jalón, D., González del Tánago, M., Grabowski, R.C., Gunn, I.D.M., Habersack, H., Hendriks, D., Henshaw, A.J., Klösch, M., Lastoria, B., Latapie, A., Marcinkowski, P., Martínez-Fernández, V., Mosselman, E., Mountford, J.O., Nardi, L., Okruszko, T., O'Hare, M.T., Palma, M., Percopo, C., Surian, N., van de Bund, W., Weissteiner, C., Ziliani, L., 2016. A multi-scale hierarchical framework for developing understanding of river behaviour to support river management. *Aquat Sci* 78, 1–16. <https://doi.org/10.1007/s00027-015-0424-5>
- Hamshaw, S.D., Bryce, T., Rizzo, D.M., O'Neil-Dunne, J., Frolik, J., Dewoolkar, M.M., 2017. Quantifying streambank movement and topography using unmanned aircraft system photogrammetry with comparison to terrestrial laser scanning. *River Research and Applications* 33, 1354–1367. <https://doi.org/10.1002/rra.3183>
- Hemmelder, S., Marra, W., Markies, H., De Jong, S.M., 2018. Monitoring river morphology & bank erosion using UAV imagery – A case study of the river Buëch, Hautes-Alpes, France. *International Journal of Applied Earth Observation and Geoinformation* 73, 428–437. <https://doi.org/10.1016/j.jag.2018.07.016>
- Hernández-Serna, A., Jiménez-Segura, L.F., 2014. Automatic identification of species with neural networks. *PeerJ* 2, e563. <https://doi.org/10.7717/peerj.563>
- Hintze, J.L., Nelson, R.D., 1998. Violin Plots: A Box Plot-Density Trace Synergism. *The American Statistician* 52, 181–184. <https://doi.org/10.1080/00031305.1998.10480559>
- Hripcsak, G., Rothschild, A.S., 2005. Agreement, the F-Measure, and Reliability in Information Retrieval. *J Am Med Inform Assoc* 12, 296–298. <https://doi.org/10.1197/jamia.M1733>
- Isikdogan, F., Bovik, A., Passalacqua, P., 2018. Learning a River Network Extractor Using an Adaptive Loss Function. *IEEE Geoscience and Remote Sensing Letters* 15, 813–817. <https://doi.org/10.1109/LGRS.2018.2811754>
- Jain, A.K., Jianchang Mao, Mohiuddin, K.M., 1996. Artificial neural networks: a tutorial. *Computer* 29, 31–44. <https://doi.org/10.1109/2.485891>

- Kalacska, M., Lucanus, O., Sousa, L., Vieira, T., Arroyo-Mora, J.P., 2019. UAV-Based 3D Point Clouds of Freshwater Fish Habitats, Xingu River Basin, Brazil. *Data* 4, 9. <https://doi.org/10.3390/data4010009>
- Kampffmeyer, M., Salberg, A.-B., Jenssen, R., 2016. Semantic Segmentation of Small Objects and Modeling of Uncertainty in Urban Remote Sensing Images Using Deep Convolutional Neural Networks. Presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 1–9.
- Khan, S.H., He, X., Porikli, F., Bennamoun, M., 2017. Forest Change Detection in Incomplete Satellite Images With Deep Neural Networks. *IEEE Transactions on Geoscience and Remote Sensing* 55, 5407–5423. <https://doi.org/10.1109/TGRS.2017.2707528>
- Krawczyk, B., 2016. Learning from imbalanced data: open challenges and future directions. *Prog Artif Intell* 5, 221–232. <https://doi.org/10.1007/s13748-016-0094-0>
- Kuhn, C., de Matos Valerio, A., Ward, N., Loken, L., Sawakuchi, H.O., Kampel, M., Richey, J., Stadler, P., Crawford, J., Striegl, R., Vermote, E., Pahlevan, N., Butman, D., 2019. Performance of Landsat-8 and Sentinel-2 surface reflectance products for river remote sensing retrievals of chlorophyll-a and turbidity. *Remote Sensing of Environment* 224, 104–118. <https://doi.org/10.1016/j.rse.2019.01.023>
- Kussul, N., Lavreniuk, M., Skakun, S., Shelestov, A., 2017. Deep Learning Classification of Land Cover and Crop Types Using Remote Sensing Data. *IEEE Geoscience and Remote Sensing Letters* 14, 778–782. <https://doi.org/10.1109/LGRS.2017.2681128>
- Labatut, V., Cherifi, H., 2012. Accuracy Measures for the Comparison of Classifiers. *arXiv:1207.3790*.
- Laliberte, A.S., Goforth, M.A., Steele, C.M., Rango, A., 2011. Multispectral Remote Sensing from Unmanned Aircraft: Image Processing Workflows and Applications for Rangeland Environments. *Remote Sensing* 3, 2529–2551. <https://doi.org/10.3390/rs3112529>
- Landis, J.R., Koch, G.G., 1977. The Measurement of Observer Agreement for Categorical Data. *Biometrics* 33, 159–174. <https://doi.org/10.2307/2529310>
- Langat, P.K., Kumar, L., Koech, R., Ghosh, M.K., 2020. Characterisation of channel morphological pattern changes and flood corridor dynamics of the tropical Tana River fluvial systems, Kenya. *Journal of African Earth Sciences* 163, 103748. <https://doi.org/10.1016/j.jafrearsci.2019.103748>
- LeCun, Y., Bengio, Y., Hinton, G., 2015a. Deep learning. *Nature* 521, 436–444. <https://doi.org/10.1038/nature14539>
- LeCun, Y., Bengio, Y., Hinton, G., 2015b. Deep learning. *Nature* 521, 436–444. <https://doi.org/10.1038/nature14539>
- Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86, 2278–2324. <https://doi.org/10.1109/5.726791>
- Legleiter, C.J., Goodchild, M.F., 2005. Alternative representations of in-stream habitat: classification using remote sensing, hydraulic modeling, and fuzzy logic. *International Journal of Geographical Information Science* 19, 29–50. <https://doi.org/10.1080/13658810412331280220>
- Legleiter, C.J., Marcus, W.A., Lawrence, R.L., 2002. Effects of Sensor Resolution on Mapping InStream Habitats. *Photogrammetric Engineering and Remote Sensing* 68, 801–807.
- Legleiter, C.J., Roberts, D.A., Marcus, W.A., Fonstad, M.A., 2004. Passive optical remote sensing of river channel morphology and in-stream habitat: Physical basis and

- feasibility. *Remote Sensing of Environment* 93, 493–510. <https://doi.org/10.1016/j.rse.2004.07.019>
- Lemaitre, G., Nogueira, F., Aridas, C.K., 2016. Imbalanced-learn: A Python Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning. arXiv:1609.06570.
- Li, P., Guo, J., Song, B., Xiao, X., 2011. A Multilevel Hierarchical Image Segmentation Method for Urban Impervious Surface Mapping Using Very High Resolution Imagery. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 4, 103–116. <https://doi.org/10.1109/JSTARS.2010.2074186>
- Li, W., Fu, H., Yu, L., Cracknell, A., 2017. Deep Learning Based Oil Palm Tree Detection and Counting for High-Resolution Remote Sensing Images. *Remote Sensing* 9, 22.
- Ling, F., Boyd, D., Ge, Y., Foody, G.M., Li, X., Wang, L., Zhang, Y., Shi, L., Shang, C., Li, X., Du, Y., 2019. Measuring River Wetted Width from Remotely Sensed Imagery at the Sub-pixel Scale with a Deep Convolutional Neural Network. *Water Resources Research* in press. <https://doi.org/10.1029/2018WR024136>
- Linke, S., Pressey, R.L., Bailey, R.C., Norris, R.H., 2007. Management options for river conservation planning: condition and conservation re-visited. *Freshwater Biology* 52, 918–938. <https://doi.org/10.1111/j.1365-2427.2006.01690.x>
- Long, J., Shelhamer, E., Darrell, T., 2015. Fully convolutional networks for semantic segmentation, in: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Presented at the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3431–3440. <https://doi.org/10.1109/CVPR.2015.7298965>
- MacVicar, B., Piégay, H., 2012. Implementation and validation of video monitoring for wood budgeting in a wandering piedmont river, the Ain River (France). *Earth Surface Processes and Landforms* 37, 1272–1289. <https://doi.org/10.1002/esp.3240>
- MacVicar, B.J., Hauet, A., Bergeron, N., Tougne, L., Ali, I., 2012. River Monitoring with Ground-Based Videography, in: *Fluvial Remote Sensing for Science and Management*. John Wiley & Sons, Ltd, pp. 367–383. <https://doi.org/10.1002/9781119940791.ch16>
- MacVicar, B.J., Piégay, H., 2012. Validation of video monitoring technique to measure wood transport in a river, in: *River Flow 2012 - Proceedings of the International Conference on Fluvial Hydraulics*. pp. 735–740.
- Mahdianpari, M., Salehi, B., Rezaee, M., Mohammadimanesh, F., Zhang, Y., 2018. Very Deep Convolutional Neural Networks for Complex Land Cover Mapping Using Multispectral Remote Sensing Imagery. *Remote Sensing* 10, 1119.
- Marcus, W.A., Fonstad, M.A., Legleiter, C.J., 2012. Management Applications of Optical Remote Sensing in the Active River Channel, in: *Fluvial Remote Sensing for Science and Management*. John Wiley & Sons, Ltd, pp. 19–41. <https://doi.org/10.1002/9781119940791.ch2>
- Marcus, W.A., Legleiter, C.J., Aspinall, R.J., Boardman, J.W., Crabtree, R.L., 2003. High spatial resolution hyperspectral mapping of in-stream habitats, depths, and woody debris in mountain streams. *Geomorphology, Mountain Geomorphology - Integrating Earth Systems, Proceedings of the 32nd Annual Binghamton Geomorphology Symposium* 55, 363–380. [https://doi.org/10.1016/S0169-555X\(03\)00150-8](https://doi.org/10.1016/S0169-555X(03)00150-8)
- McKinney, W., 2010. Data Structures for Statistical Computing in Python. Presented at the Proceedings of the 9th Python in Science Conference, pp. 51–56.
- Michez, A., Piégay, H., Lisein, J., Claessens, H., Lejeune, P., 2016. Classification of riparian forest species and health condition using multi-temporal and hyperspatial

- imagery from unmanned aerial system. *Environ Monit Assess* 188, 146. <https://doi.org/10.1007/s10661-015-4996-2>
- Nel, J.L., Reyers, B., Roux, D.J., Cowling, R.M., 2009. Expanding protected areas beyond their terrestrial comfort zone: Identifying spatial options for river conservation. *Biological Conservation* 142, 1605–1616. <https://doi.org/10.1016/j.biocon.2009.02.031>
- Olmanson, L.G., Brezonik, P.L., Bauer, M.E., 2013. Airborne hyperspectral remote sensing to assess spatial distribution of water quality characteristics in large rivers: The Mississippi River and its tributaries in Minnesota. *Remote Sensing of Environment* 130, 254–265. <https://doi.org/10.1016/j.rse.2012.11.023>
- Ormerod, S.J., 2009. Climate change, river conservation and the adaptation challenge. *Aquatic Conservation: Marine and Freshwater Ecosystems* 19, 609–613. <https://doi.org/10.1002/aqc.1062>
- Otukey, J.R., Blaschke, T., 2010. Land cover change assessment using decision trees, support vector machines and maximum likelihood classification algorithms. *International Journal of Applied Earth Observation and Geoinformation, Supplement Issue on “Remote Sensing for Africa – A Special Collection from the African Association for Remote Sensing of the Environment (AARSE)”* 12, S27–S31. <https://doi.org/10.1016/j.jag.2009.11.002>
- Pal, M., 2005. Random forest classifier for remote sensing classification. *International Journal of Remote Sensing* 26, 217–222. <https://doi.org/10.1080/01431160412331269698>
- Palmer, M.A., Menninger, H.L., Bernhardt, E., 2010. River restoration, habitat heterogeneity and biodiversity: a failure of theory or practice? *Freshwater Biology* 55, 205–222. <https://doi.org/10.1111/j.1365-2427.2009.02372.x>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, É., 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12, 2825–2830.
- Poggi, G., Scarpa, G., Zerubia, J.B., 2005. Supervised segmentation of remote sensing images based on a tree-structured MRF model. *IEEE Transactions on Geoscience and Remote Sensing* 43, 1901–1911. <https://doi.org/10.1109/TGRS.2005.852163>
- Pölonen, I., Saari, H., Kaivosoja, J., Honkavaara, E., Pesonen, L., 2013. Hyperspectral imaging based biomass and nitrogen content estimations from light-weight UAV, in: *Remote Sensing for Agriculture, Ecosystems, and Hydrology XV*. Presented at the Remote Sensing for Agriculture, Ecosystems, and Hydrology XV, International Society for Optics and Photonics, p. 88870J. <https://doi.org/10.1117/12.2028624>
- Pouliot, D., Latifovic, R., Pasher, J., Duffe, J., 2019. Assessment of Convolution Neural Networks for Wetland Mapping with Landsat in the Central Canadian Boreal Forest Region. *Remote Sensing* 11, 772.
- Purinton, B., Bookhagen, B., 2019. Introducing PebbleCounts: a grain-sizing tool for photo surveys of dynamic gravel-bed rivers. *Earth Surface Dynamics* 7, 859–877. <https://doi.org/10.5194/esurf-7-859-2019>
- Rogger, M., Agnoletti, M., Alaoui, A., Bathurst, J.C., Bodner, G., Borga, M., Chaplot, V., Gallart, F., Glatzel, G., Hall, J., Holden, J., Holko, L., Horn, R., Kiss, A., Kohnová, S., Leitinger, G., Lennartz, B., Parajka, J., Perdigão, R., Peth, S., Plavcová, L., Quinton, J.N., Robinson, M., Salinas, J.L., Santoro, A., Szolgay, J., Tron, S., van den Akker, J.J.H., Viglione, A., Blöschl, G., 2017. Land use change impacts on floods at the catchment scale: Challenges and opportunities for future research. *Water Resources Research* 53, 5209–5219. <https://doi.org/10.1002/2017wr020723>

- Romero, A., Gatta, C., Camps-Valls, G., 2016. Unsupervised Deep Feature Extraction for Remote Sensing Image Classification. *IEEE Transactions on Geoscience and Remote Sensing* 54, 1349–1362. <https://doi.org/10.1109/TGRS.2015.2478379>
- Rosenberg, D.M., McCully, P., Pringle, C.M., 2000. Global-Scale Environmental Effects of Hydrological Alterations: Introduction. *BioScience* 50, 746–751. [https://doi.org/10.1641/0006-3568\(2000\)050\[0746:GSEEOH\]2.0.CO;2](https://doi.org/10.1641/0006-3568(2000)050[0746:GSEEOH]2.0.CO;2)
- Rusnák, M., Sládek, J., Kidová, A., Lehotský, M., 2018. Template for high-resolution river landscape mapping using UAV technology. *Measurement* 115, 139–151. <https://doi.org/10.1016/j.measurement.2017.10.023>
- Seitz, L., Haas, C., Noack, M., Wieprecht, S., 2018. From picture to porosity of river bed material using Structure-from-Motion with Multi-View-Stereo. *Geomorphology* 306, 80–89. <https://doi.org/10.1016/j.geomorph.2018.01.014>
- Seto, K.C., Woodcock, C.E., Song, C., Huang, X., Lu, J., Kaufmann, R.K., 2002. Monitoring land-use change in the Pearl River Delta using Landsat TM. *International Journal of Remote Sensing* 23, 1985–2004. <https://doi.org/10.1080/01431160110075532>
- Smeeton, N.C., 1985. Early History of the Kappa Statistic. *Biometrics* 41, 795–795.
- Smikrud, K.M., Prakash, A., Nichols, J.V., 2008. Decision-Based Fusion for Improved Fluvial Landscape Classification Using Digital Aerial Photographs and Forward Looking Infrared Images. *Photogrammetric Engineering & Remote Sensing* 74, 903–911. <https://doi.org/doi:10.14358/PERS.74.7.903>
- Smith, L.C., 1997. Satellite remote sensing of river inundation area, stage, and discharge: a review. *Hydrological Processes* 11, 1427–1439. [https://doi.org/10.1002/\(SICI\)1099-1085\(199708\)11:10<1427::AID-HYP473>3.0.CO;2-S](https://doi.org/10.1002/(SICI)1099-1085(199708)11:10<1427::AID-HYP473>3.0.CO;2-S)
- Solomon, C., Breckon, T., 2011. *Fundamentals of Digital Image Processing: A Practical Approach with Examples in Matlab*, 1st ed. Wiley Publishing.
- Spada, D., Molinari, P., Bertoldi, W., Vitti, A., Zolezzi, G., 2018. Multi-Temporal Image Analysis for Fluvial Morphological Characterization with Application to Albanian Rivers. *ISPRS International Journal of Geo-Information* 7, 314. <https://doi.org/10.3390/ijgi7080314>
- Srivastava, M., Grill-Spector, K., 2018. The Effect of Learning Strategy versus Inherent Architecture Properties on the Ability of Convolutional Neural Networks to Develop Transformation Invariance. *arXiv:1810.13128*.
- Strahler, A.H., 1980. The use of prior probabilities in maximum likelihood classification of remotely sensed data. *Remote Sensing of Environment* 10, 135–163. [https://doi.org/10.1016/0034-4257\(80\)90011-5](https://doi.org/10.1016/0034-4257(80)90011-5)
- Strayer, D.L., Dudgeon, D., 2010. Freshwater biodiversity conservation: recent progress and future challenges. *Journal of the North American Benthological Society* 29, 344–358. <https://doi.org/10.1899/08-171.1>
- Stumpf, A., Kerle, N., 2011. Object-oriented mapping of landslides using Random Forests. *Remote Sensing of Environment* 115, 2564–2577. <https://doi.org/10.1016/j.rse.2011.05.013>
- Szegedy, C., Wei Liu, Yangqing Jia, Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A., 2015. Going deeper with convolutions, in: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Presented at the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1–9. <https://doi.org/10.1109/CVPR.2015.7298594>
- Tammaing, A., Hugenholtz, C., Eaton, B., Lapointe, M., 2015. Hyperspatial Remote Sensing of Channel Reach Morphology and Hydraulic Fish Habitat Using an

- Unmanned Aerial Vehicle (UAV): A First Assessment in the Context of River Research and Management. *River Res. Applic.* 31, 379–391. <https://doi.org/10.1002/rra.2743>
- Tian, Y.Q., Yu, Q., Zimmerman, M.J., Flint, S., Waldron, M.C., 2010. Differentiating aquatic plant communities in a eutrophic river using hyperspectral and multispectral remote sensing. *Freshwater Biology* 55, 1658–1673. <https://doi.org/10.1111/j.1365-2427.2010.02400.x>
- van Vliet, M.T.H., Franssen, W.H.P., Yearsley, J.R., Ludwig, F., Haddeland, I., Lettenmaier, D.P., Kabat, P., 2013. Global river discharge and water temperature under climate change. *Global Environmental Change* 23, 450–464. <https://doi.org/10.1016/j.gloenvcha.2012.11.002>
- Vanegas, F., Bratanov, D., Powell, K., Weiss, J., Gonzalez, F., 2018. A Novel Methodology for Improving Plant Pest Surveillance in Vineyards and Crops Using UAV-Based Hyperspectral and Spatial Data. *Sensors* 18, 260. <https://doi.org/10.3390/s18010260>
- Vannote, R.L., Minshall, G.W., Cummins, K.W., Sedell, J.R., Cushing, C.E., 1980. The River Continuum Concept. *Can. J. Fish. Aquat. Sci.* 37, 130–137. <https://doi.org/10.1139/f80-017>
- Vörösmarty, C.J., McIntyre, P.B., Gessner, M.O., Dudgeon, D., Prusevich, A., Green, P., Glidden, S., Bunn, S.E., Sullivan, C.A., Liermann, C.R., Davies, P.M., 2010. Global threats to human water security and river biodiversity. *Nature* 467, 555–561. <https://doi.org/10.1038/nature09440>
- Walt, S. van der, Schönberger, J.L., Nunez-Iglesias, J., Boulogne, F., Warner, J.D., Yager, N., Goullart, E., Yu, T., 2014. scikit-image: image processing in Python. *PeerJ* 2, e453. <https://doi.org/10.7717/peerj.453>
- Wang, C., Pavlowsky, R.T., Huang, Q., Chang, C., 2016. Channel bar feature extraction for a mining-contaminated river using high-spatial multispectral remote-sensing imagery. *GIScience & Remote Sensing* 53, 283–302. <https://doi.org/10.1080/15481603.2016.1148229>
- Ward, J.V., Tockner, K., Uehlinger, U., Malard, F., 2001. Understanding natural patterns and processes in river corridors as the basis for effective river restoration. *Regulated Rivers: Research & Management* 17, 311–323. <https://doi.org/10.1002/rrr.646>
- Willis, A., Holmes, E., 2019. Eye in the Sky: Using UAV Imagery of Seasonal Riverine Canopy Growth to Model Water Temperature. *Hydrology* 6, 6. <https://doi.org/10.3390/hydrology6010006>
- Winterbottom, S.J., Gilvear, D.J., 1997. Quantification of channel bed morphology in gravel-bed rivers using airborne multispectral imagery and aerial photography. *Regul. Rivers: Res. Mgmt.* 13, 489–499. [https://doi.org/10.1002/\(SICI\)1099-1646\(199711/12\)13:6<489::AID-RRR471>3.0.CO;2-X](https://doi.org/10.1002/(SICI)1099-1646(199711/12)13:6<489::AID-RRR471>3.0.CO;2-X)
- Wohl, E., Angermeier, P.L., Bledsoe, B., Kondolf, G.M., MacDonnell, L., Merritt, D.M., Palmer, M.A., Poff, N.L., Tarboton, D., 2005. River restoration. *Water Resources Research* 41. <https://doi.org/10.1029/2005wr003985>
- Woodget, A.S., Austrums, R., 2017. Subaerial gravel size measurement using topographic data derived from a UAV-SfM approach. *Earth Surf. Process. Landforms* 42, 1434–1443. <https://doi.org/10.1002/esp.4139>
- Woodget, A.S., Austrums, R., Maddock, I.P., Habit, E., 2017. Drones and digital photogrammetry: from classifications to continuums for monitoring river habitat and hydromorphology. *Wiley Interdisciplinary Reviews: Water* 4, e1222-n/a. <https://doi.org/10.1002/wat2.1222>

- Woodget, A.S., Carbonneau, P.E., Visser, F., Maddock, I.P., 2015. Quantifying submerged fluvial topography using hyperspatial resolution UAS imagery and structure from motion photogrammetry. *Earth Surf. Process. Landforms* 40, 47–64. <https://doi.org/10.1002/esp.3613>
- Woodget, A.S., Visser, F., Maddock, I.P., Carbonneau, P.E., 2016. The Accuracy and Reliability of Traditional Surface Flow Type Mapping: Is it Time for a New Method of Characterizing Physical River Habitat? *River Research and Applications* 32, 1902–1914. <https://doi.org/10.1002/rra.3047>
- WWF, 2018. Living Planet Report 2018: Aiming higher. World Wildlife Fund, Gland, Switzerland.
- Yang, X., Damen, M.C.J., van Zuidam, R.A., 1999. Satellite remote sensing and GIS for the analysis of channel migration changes in the active Yellow River Delta, China. *International Journal of Applied Earth Observation and Geoinformation* 1, 146–157. [https://doi.org/10.1016/S0303-2434\(99\)85007-7](https://doi.org/10.1016/S0303-2434(99)85007-7)
- Zhang, L., Zhang, L., Du, B., 2016. Deep Learning for Remote Sensing Data: A Technical Tutorial on the State of the Art. *IEEE Geoscience and Remote Sensing Magazine* 4, 22–40. <https://doi.org/10.1109/MGRS.2016.2540798>
- Zhang, Y.K., Schilling, K.E., 2006. Increasing streamflow and baseflow in Mississippi River since the 1940s: Effect of land use change. *Journal of Hydrology* 324, 412–422. <https://doi.org/10.1016/j.jhydrol.2005.09.033>
- Zheng, C., Wang, P., 1996. Parameter structure identification using tabu search and simulated annealing. *Advances in Water Resources* 19, 215–224. [https://doi.org/10.1016/0309-1708\(96\)00047-4](https://doi.org/10.1016/0309-1708(96)00047-4)
- Zhong, Y., Fei, F., Liu, Y., Zhao, B., Jiao, H., Zhang, L., 2017. SatCNN: satellite image dataset classification using agile convolutional neural networks. *Remote Sensing Letters* 8, 136–145. <https://doi.org/10.1080/2150704X.2016.1235299>
- Zhong, Y., Zhang, L., 2012. An Adaptive Artificial Immune Network for Supervised Classification of Multi-/Hyperspectral Remote Sensing Imagery. *IEEE Transactions on Geoscience and Remote Sensing* 50, 894–909. <https://doi.org/10.1109/TGRS.2011.2162589>
- Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V., 2017. Learning Transferable Architectures for Scalable Image Recognition. *arXiv:1707.07012*.

Supporting Information

1- Accuracy and F1.

1315 In the figure and tables below, we give the readers unfamiliar with the F1 score some
reference points to translate and interpret our reported F1 scores in terms of traditional
accuracy values. Distributions of accuracy were compiled in the same manner as for the
F1 scores reported in the main paper on a per-image basis. For experiments 1 and 2, we
present tables of summary statistics of median(mean) accuracy (tables S1 and S2). For
1320 experiment 3, we produce a scatter plot of F1 vs. Accuracy.

The results show that for our data accuracy and F1 are very closely correlated. Figure S1
gives a regression line of $\text{Accuracy} = 1.03\text{F1} + 4.1\%$ with an R^2 of 0.96. Importantly, an
Accuracy of 100% is the same as an F1 of 100% with no bias present at high values of F1/
Accuracy. Figure S1 shows that for very high values F1 and accuracy converge to 100%.

1325

**Table S1. Accuracy summary statistics for experiment 1. Results are shown as
median(mean) accuracy values. Here N=394 for in-sample results and N=467 for out-
of-sample results.**

In-Sample Data					
	MLIK	RF	DNN	CNN	CNN+MLP
Accuracy [%]	83(79)	67(62)	66(63)	92(91)	96(94)
Out-of-Sample Data					
Accuracy [%]	46(46)	58(55)	55(55)	63(58)	78(63)

1330

**Table S2. Accuracy summary statistics for experiment 2. Results are shown as
median(mean) accuracy values. Here N=348 for in-sample results and N=513 for out-
of-sample results.**

1335

In-Sample Data				
Class	NASNet Large		NASNet Mobile	
	CNN	CNN+MLP	CNN	CNN+MLP
Accuracy [%]	98 (96)	99 (97)	97 (95)	98 (96)
Out-of-Sample Data				
Accuracy [%]	81(78)	90(84)	80(76)	88 (82)

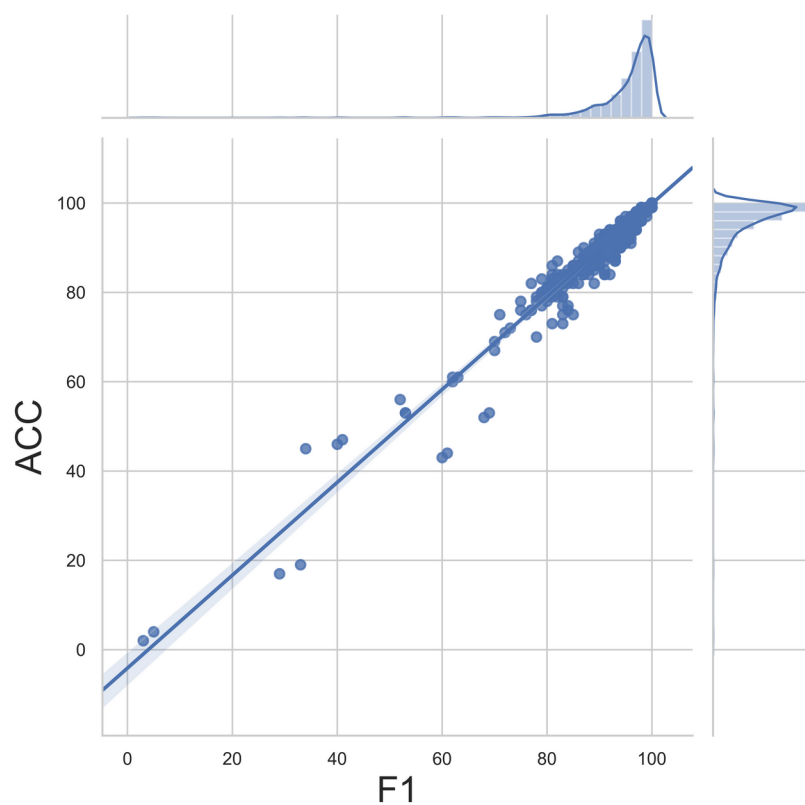


Figure S1 Accuracy (ACC) vs. F1 scatter plot. Here we use all the results from experiment 3 with the 2 phases of the CSC process (CNN and CNN+MLP) combined ($n=1724$). The regression line gives Accuracy = $1.03F1 + 4.1\%$, $R^2 = 0.96$.

2-Large scale summary of CSC results

1340

1345

1350

1355

1360

Table S3 presents F1 and kappa values estimated for the aggregate of all pixel-level predictions in all experiments. Here we concatenate the class predictions and truth labels for each pixel in each image for the listed experiment. This results in large arrays with
1365 n_{tot} size as indicated in the table (in excess of 2 billion). For computing reasons, we estimate a single F1 and Kappa value for these arrays by taking a random sample of 100 million predictions. Table 7 shows a slightly lower outcomes than tables 5 and 6. This is the most stringent test of our data, nevertheless, the observations of performance above the 90% level are not documented elsewhere in the literature. Interestingly, we note a
1370 sensitivity to the volume of training data. If we consider the outcomes for figure 1 rivers (in-sample data) we can see that quality improves as a function of training data volume: experiment 1 (12k tiles, 93%), experiment 3 (38k, tiles 96%) and experiment 2 (190k tiles, 97%). In the case of out-of-sample results from experiment 1 and 2 (the figure 2 rivers), the number of sample tiles has improved the performance from 65% F1 and a kappa of
1375 47% (considered as a 'fair' result) in the case of experiment 1 to 84% F1 and a kappa of 74% (considered a 'good' result).

1380

1385

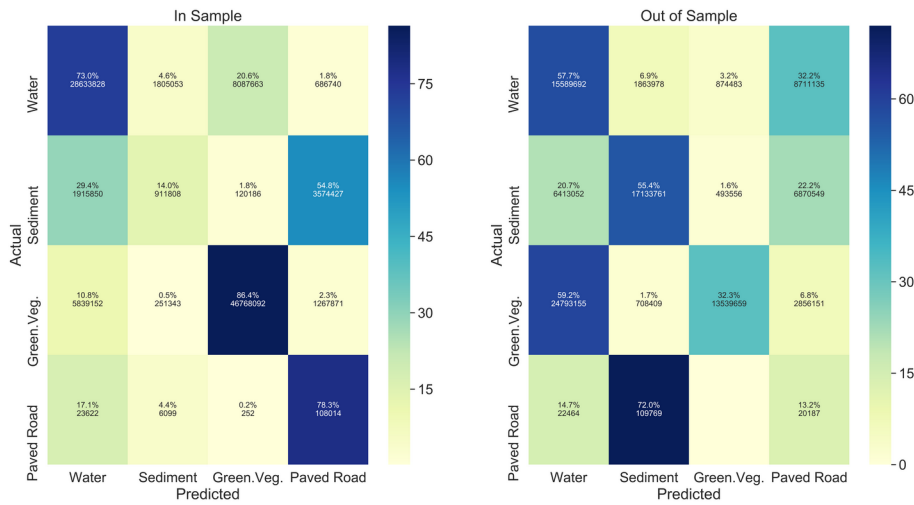
1390 **Table S3. large scale summary statistic estimates. The results for each listed experiment are aggregate for a total of n_tot pixel class predictions. F1 and kappa are estimated from a 100 million pixel sub-sample of this aggregate. Results are given as F1[%] with kappa as a 0-1 fraction. Data is only shown for the CSC method for both NASNet large and NASNet mobile when available. Figure 1 rivers are the in-sample data for experiment 2 and figure 2 rivers are the out-of-sample data for experiment 2. Experiment 3 was similarly dis-aggregated to allow for comparison, but all results from experiment 3 are from in-sample data.**

	NASNet Mobile		NASNet Large		n_tot
	CNN	CNN+MLP	CNN	CNN+MLP	
Experiment 1, Dartmouth, Ste-Marguerite	NA	NA	90(0.82)	93(0.88)	2.627E+09
Experiment 1, other 9 rivers	NA	NA	62(0.42)	71(0.53)	1.733E+09
Experiment 2, Fig. 1 rivers	95(0.91)	96(0.94)	96(0.92)	97(0.94)	1.615E+09
Experiment 2, Fig. 2 rivers	74(0.60)	81(0.70)	77(0.63)	84(0.74)	2.776E+09
Experiment 3, Fig. 1 rivers	NA	NA	95(0.91)	96(0.93)	1.615E+09
Experiment 3, Fig. 2 rivers	NA	NA	92(0.88)	93(0.89)	2.776E+09
Experiment 3, all 11 rivers	NA	NA	93(0.89)	94(0.91)	4.391E+09

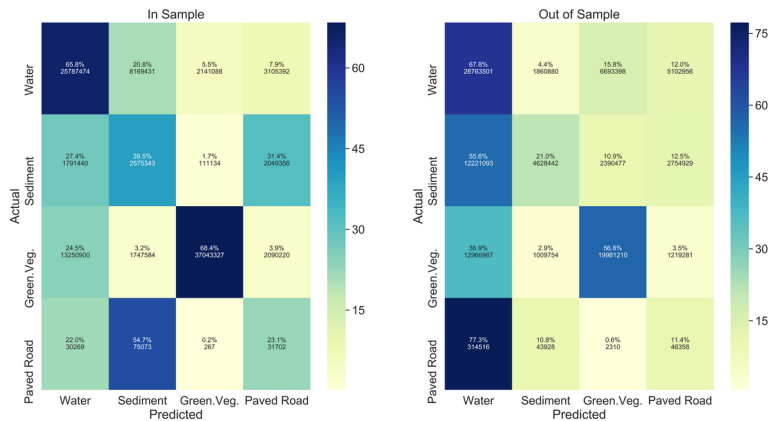
1405 3- Confusion Matrices.

Next, we present the confusion matrices associated with our results. In total, the results presented in the paper are based on 8620 evaluations of classification quality. Each of the 862 images was used 10 times throughout our various experiments; five times in
1410 experiment 1 (Maximum Likelihood, Random Forests, pixel-level MLP, CNN, and CNN+MLP), four times in the second experiment (NASnet Large (both CNN and CNN+MLP) and NASNet mobile (both CNN and CNN+MLP)) and twice in the third experiment. For a more synthetic view of confidence matrices, we have produced 30 confidence matrices in 15 figures. For each of the experiments, we cumulate and
1415 concatenate the entire set of pixel-level predictions and ground truth. This resulted in truth vs. predicted arrays with several billion rows. Then, to reduce the computational load to within our available resources, we randomly sample 100 million rows from these lists. Finally, we produce a total of 30 confidence matrices for each of the methods in each experiment. In the case of experiments 1 and 2, the figures are separated for in-sample
1420 and out-of-sample validation. In the case of experiment 3, we separate the figures according to the phase of the CSC process (CNN or CNN+MLP).

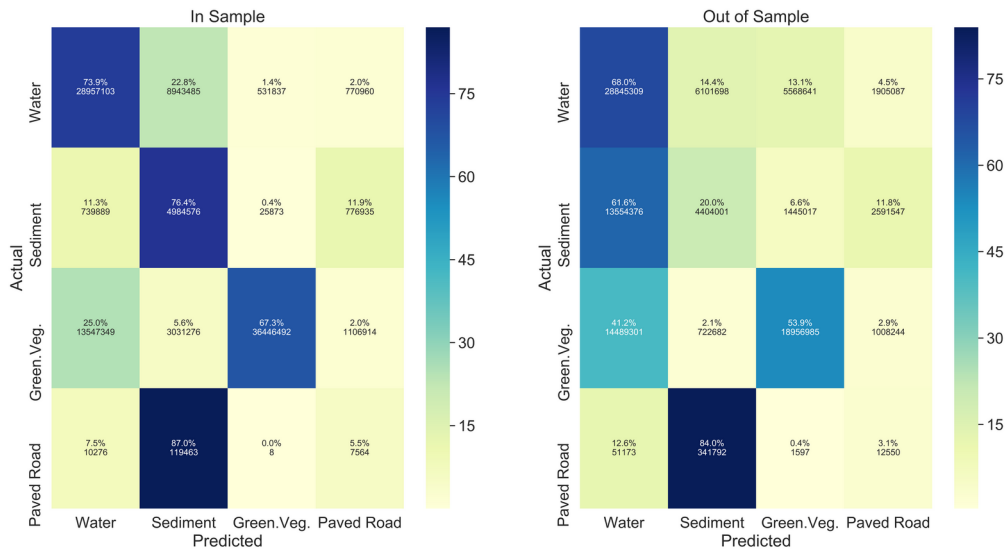
1425



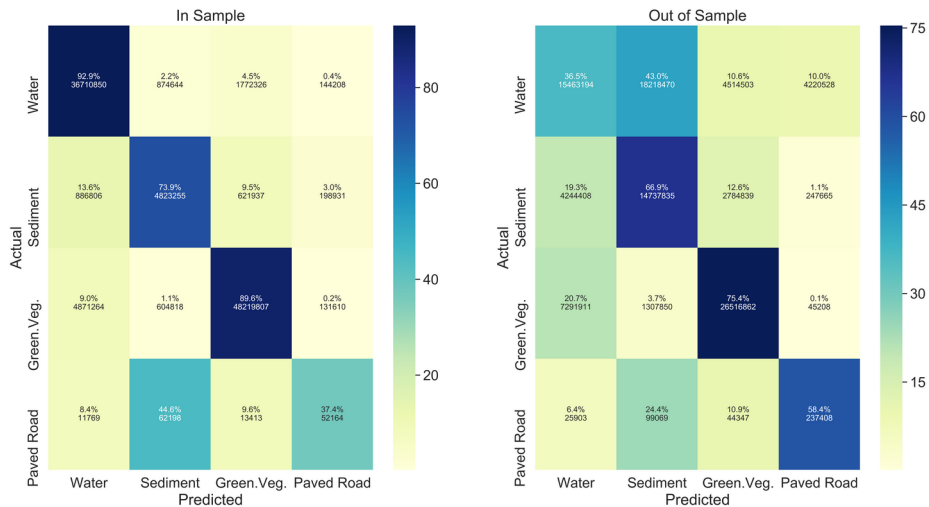
1430 *Figure S2. Confusion matrices for Maximum Likelihood outcomes in the first experiment. Left) In-sample data here from the rivers Dartmouth and Ste-Marguerite. Right) Out-of-*
 1435 *sample data here from the rivers Ouelle, Kananaskis, Pacuare, Kingie, Eamont, Sesia, Dora di Veny, Kurobe, and Kinu. For each location in the matrix, we give the number of*
samples and the percentage of the class. The color bars also indicate the percentage of
 1440 *each class in a given cell.*



1440 *Figure S3. Confusion matrices for Random Forest outcomes in the first experiment. Left) In-sample data here from the rivers Dartmouth and Ste-Marguerite. Right) Out-of-*
 1445 *sample data here from the rivers Ouelle, Kananaskis, Pacuare, Kingie, Eamont, Sesia, Dora di Veny, Kurobe, and Kinu.*



1455 *Figure S4. Confusion matrices for pixel-based Multilayer Perceptron (MLP) outcomes in the first experiment. Left) In-sample data here from the rivers Dartmouth and Ste-Marguerite. Right) Out-of-sample data here from the rivers Ouelle, Kananaskis, Pacuare, Kingie, Eamont, Sesia, Dora di Veny, Kurobe, and Kinu.*



1460 *Figure S5. Confusion matrices for CNN outcomes in the first experiment. Left) In-sample data here from the rivers Dartmouth and Ste-Marguerite. Right) Out-of-sample data here from the rivers Ouelle, Kananaskis, Pacuare, Kingie, Eamont, Sesia, Dora di Veny, Kurobe, and Kinu.*

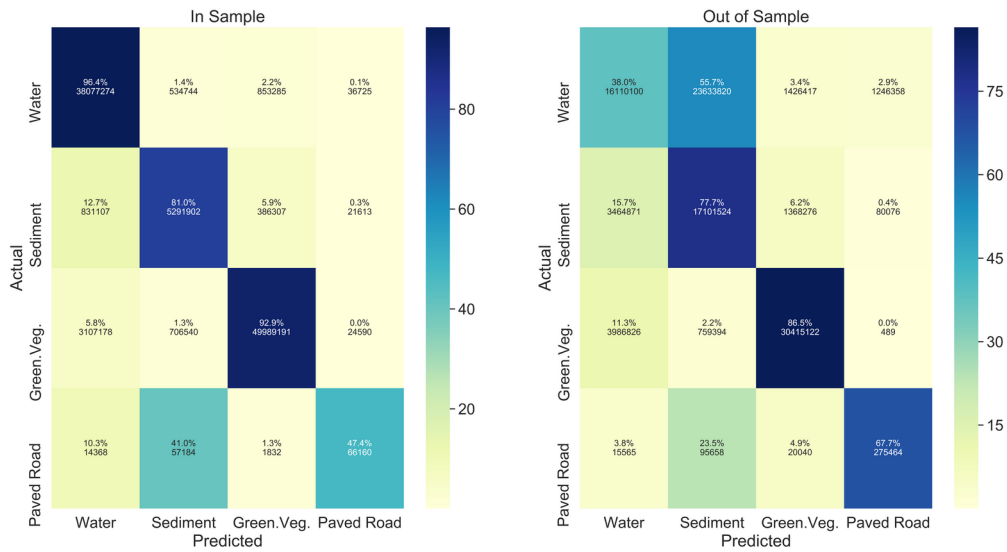


Figure S6. Confusion matrices for CNN+MLP (2 phase CSC process) outcomes in the first experiment. Left) In-sample data here from the rivers Dartmouth and Ste-Marguerite. Right) Out-of-sample data here from the rivers Ouelle, Kananaskis, Pacuare, Kingie, Eamont, Sesia, Dora di Veny, Kurobe, and Kinu.

1465

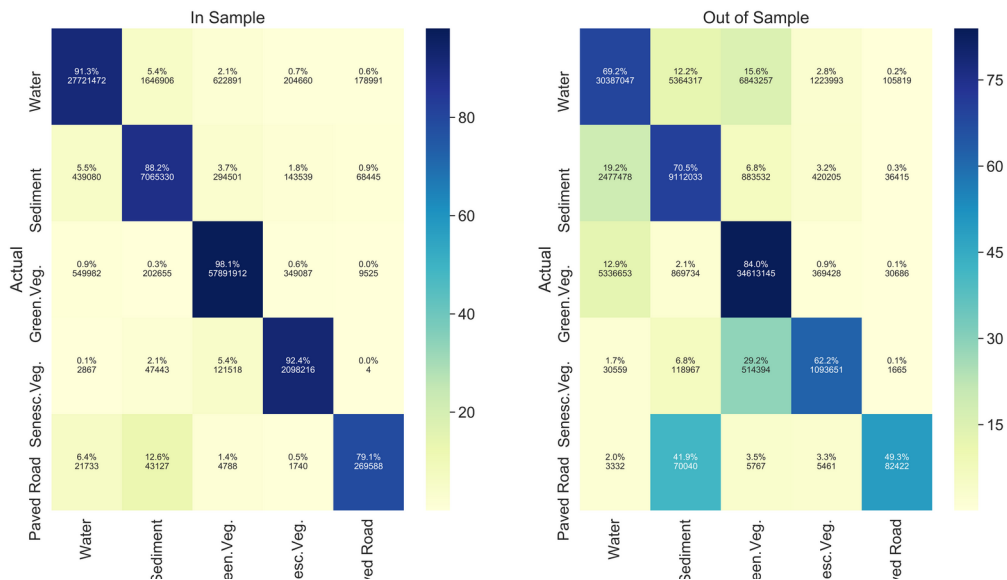
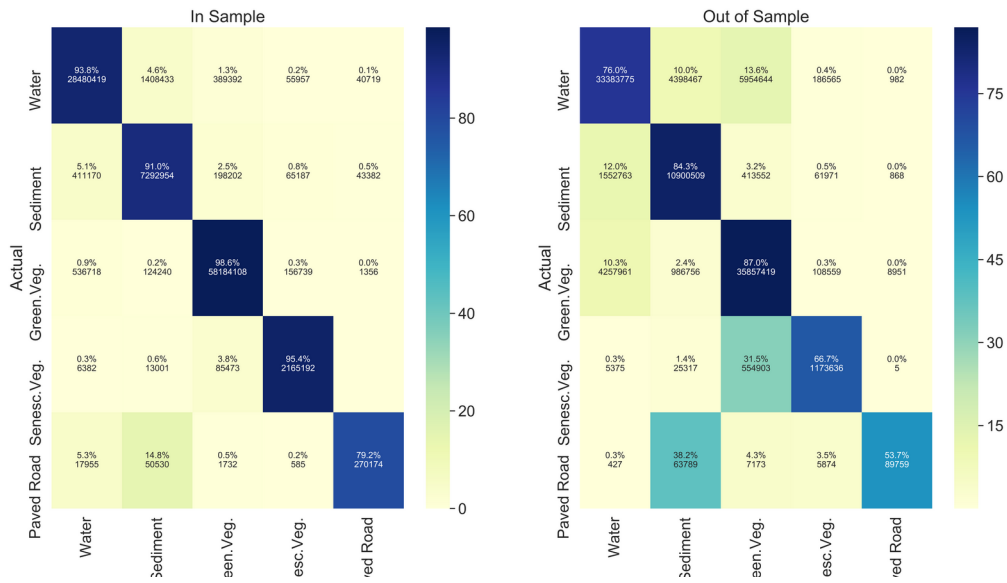
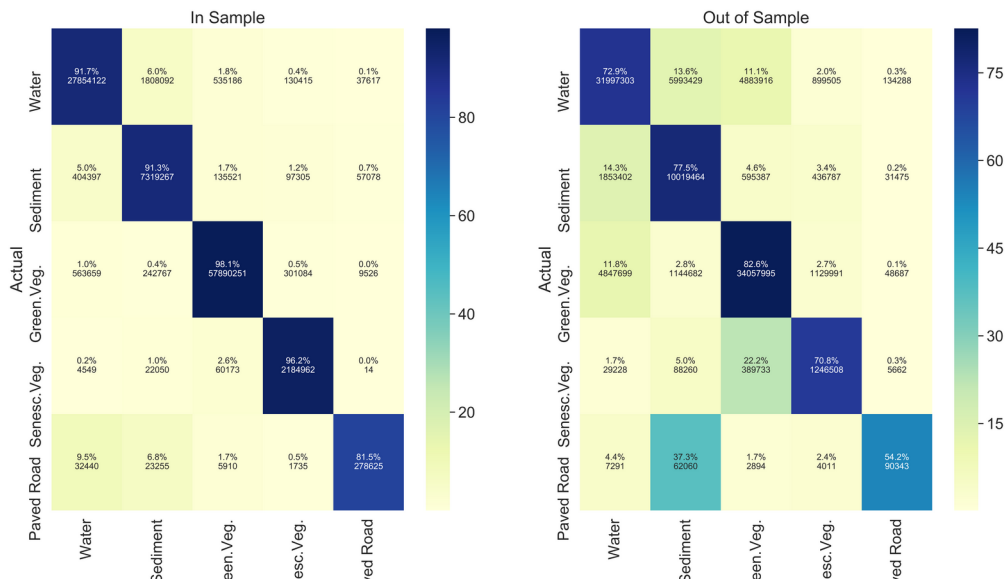


Figure S7. Confusion matrix for the second experiment, NASNet Mobile CNN, results. Left) In-sample data drawn from the rivers Ste-Marguerite, Kananaskis, Kingie, Sesia, and Kinu. Right) Out-of-sample data drawn from the rivers Dartmouth, Ouelle, Pacuare, Dora di Veny, Eamont, and Kurobe.

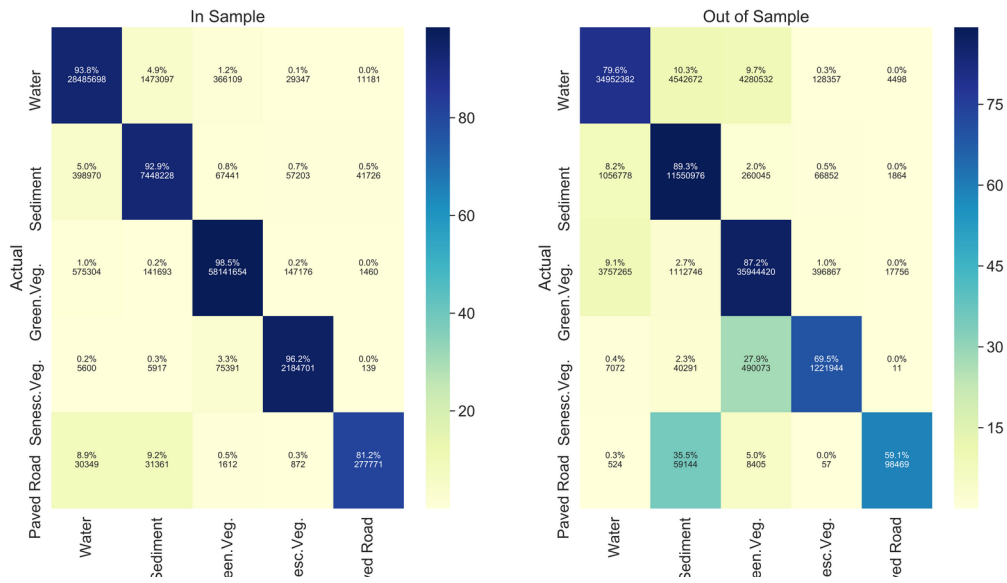
1470



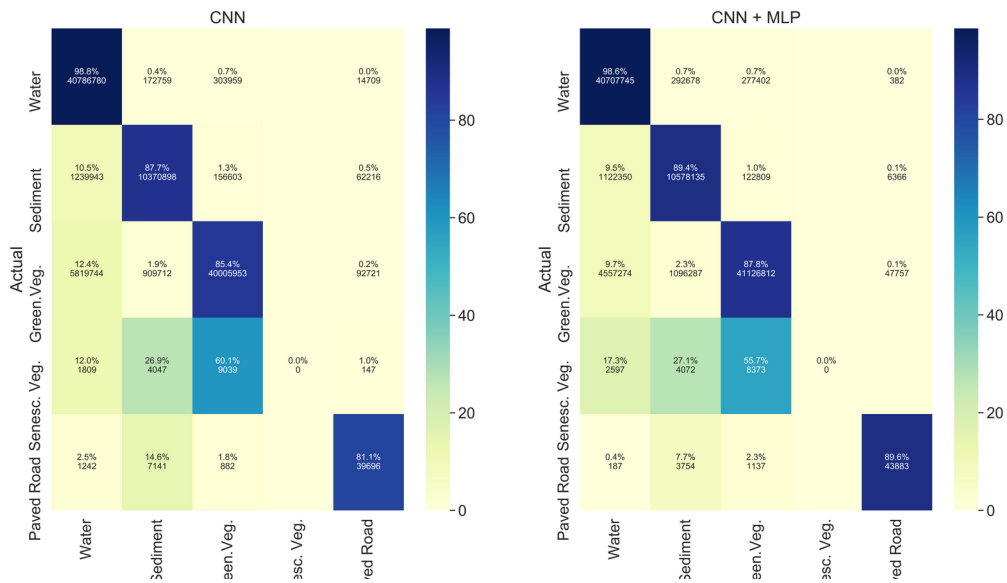
1475 *Figure S8. Confusion matrix for the second experiment, NASNet Mobile CNN+MLP, results. Left) In-sample data drawn from the rivers Ste-Marguerite, Kananaskis, Kingie, Sesia, and Kinu. Right) Out-of-sample data drawn from the rivers Dartmouth, Ouelle, Pacuare, Dora di Veny, Eamont, and Kurobe.*



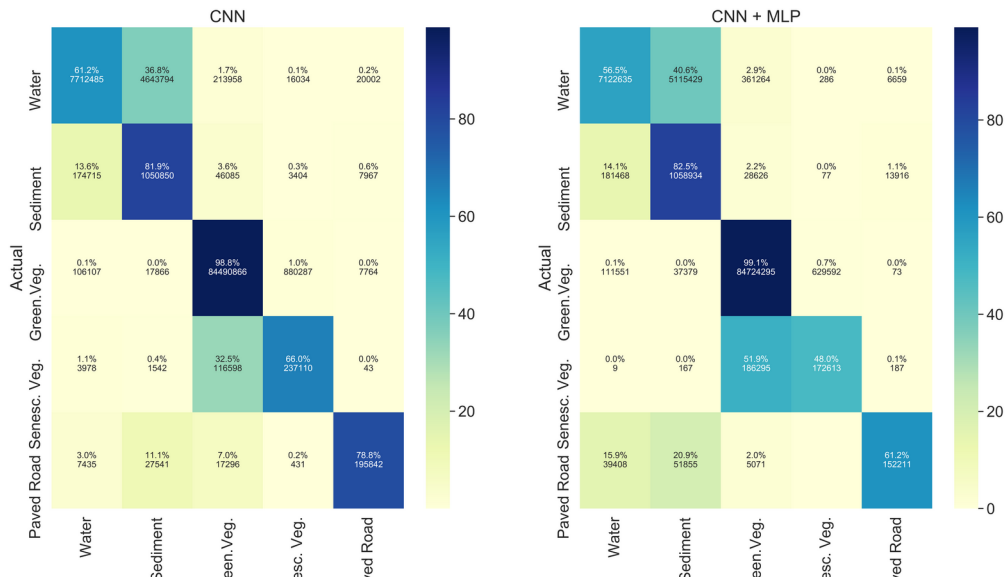
1480 *Figure S9. Confusion matrix for the second experiment, NASNet Large CNN, results. Left) In-sample data drawn from the rivers Ste-Marguerite, Kananaskis, Kingie, Sesia, and Kinu. Right) Out-of-sample data drawn from the rivers Dartmouth, Ouelle, Pacuare, Dora di Veny, Eamont, and Kurobe.*



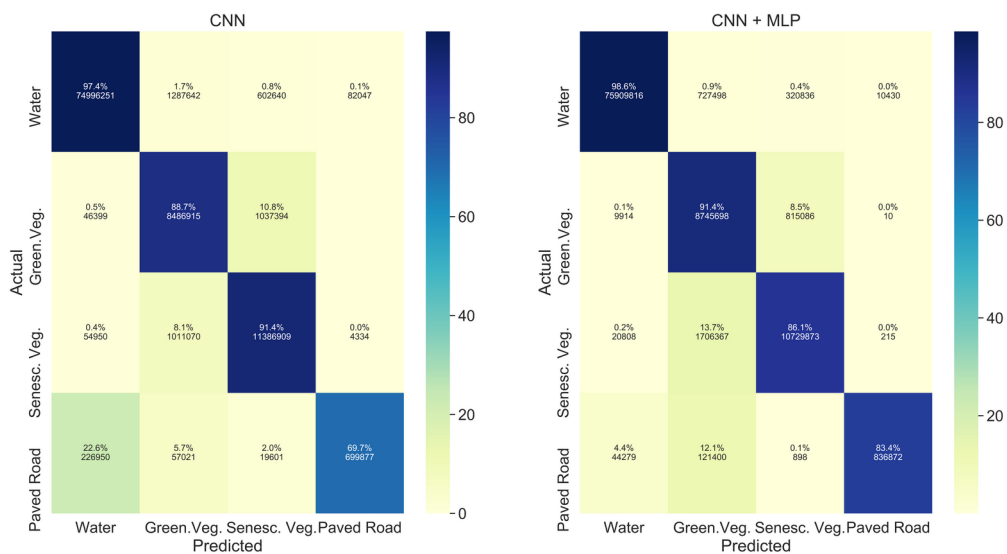
1485 **Figure S10.** Confusion matrix for the second experiment, NASNet Large CNN+MLP, results. Left) In-sample data drawn from the rivers Ste-Marguerite, Kananaskis, Kingie, Sesia, and Kinu. Right) Out-of-sample data drawn from the rivers Dartmouth, Ouelle, Pacuare, Dora di Veny, Eamont, and Kurobe.



1490 **Figure S11.** Confusion for the third experiment, Dartmouth river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1495 Figure S12. Confusion for the third experiment, Kananaskis river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1500 Figure S13. Confusion for the third experiment, Ouelle river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.

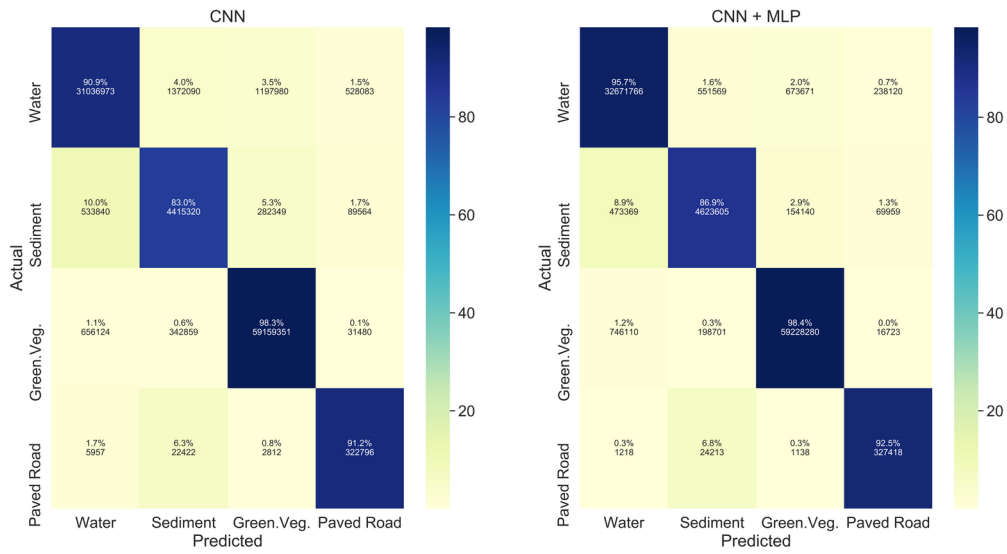
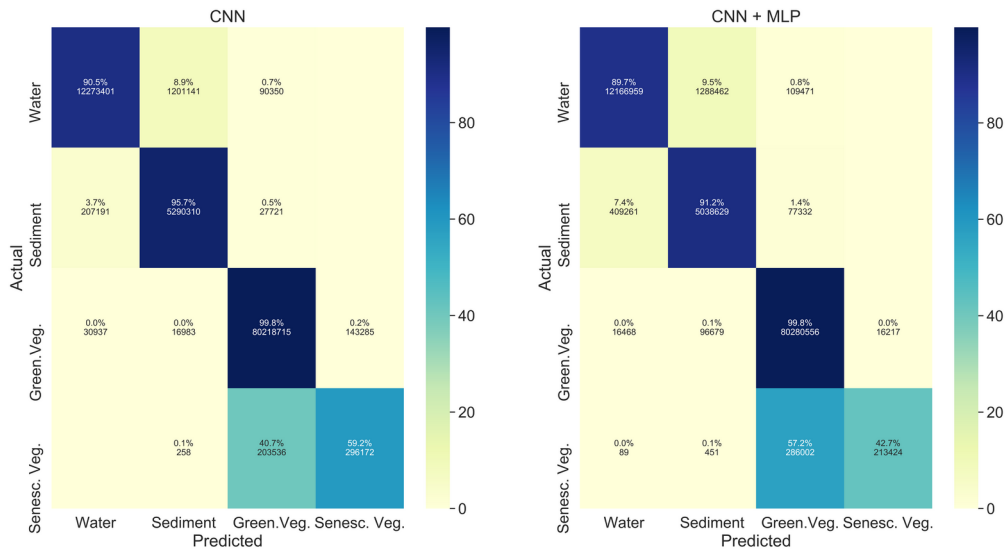


Figure S14. Confusion for the third experiment, Ste-Marguerite river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1510 Figure S15. Confusion for the third experiment, Pacuare river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.

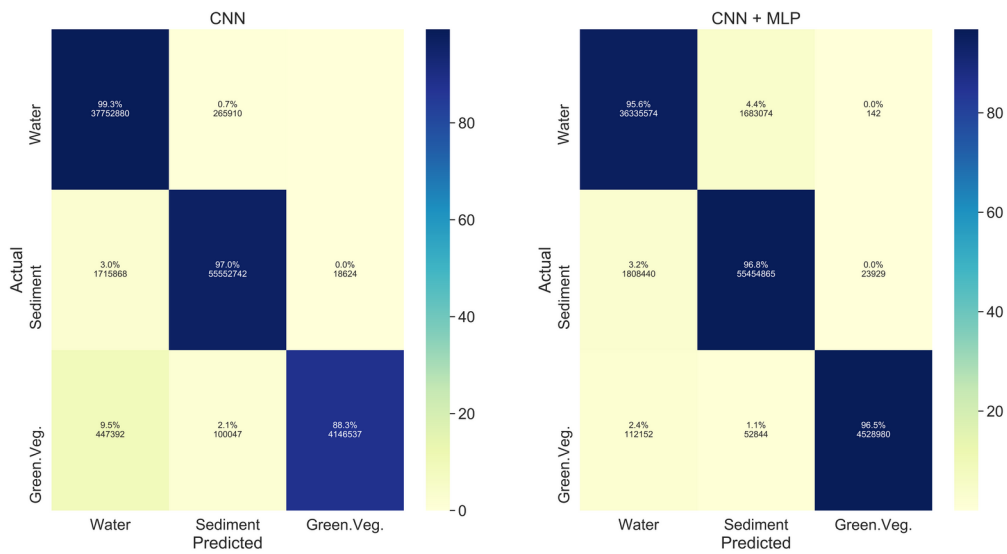
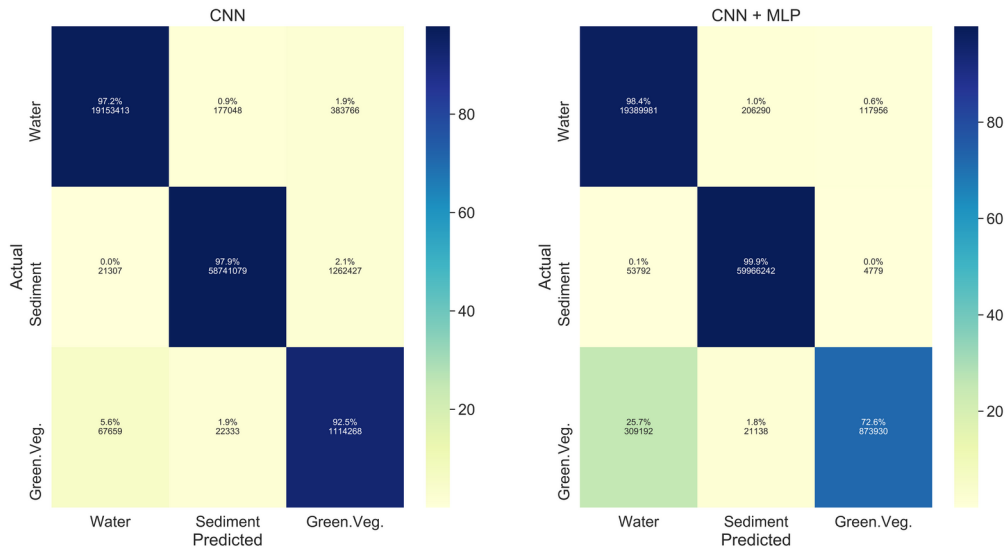
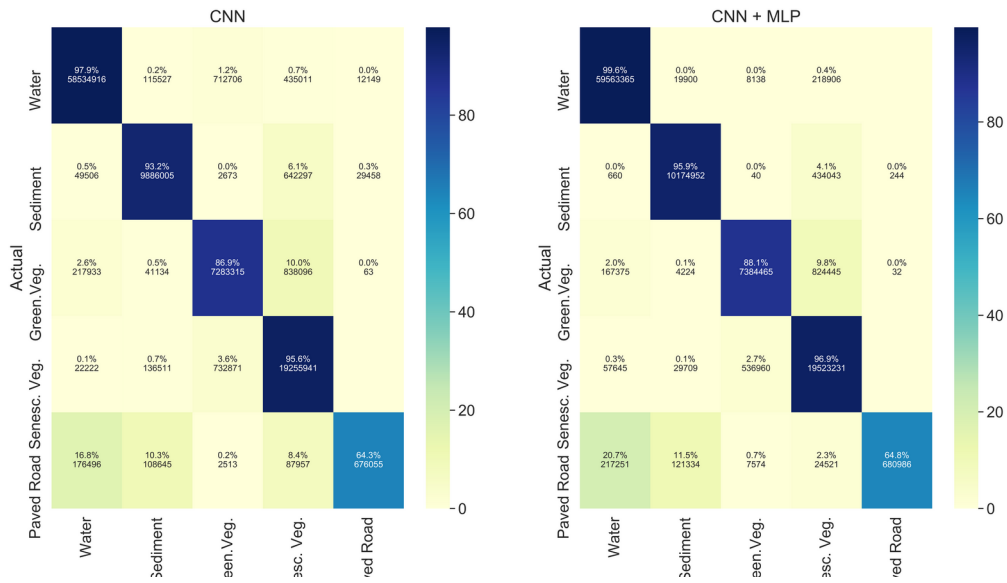


Figure S16. Confusion for the third experiment, Dora di Veny river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1520 Figure S17. Confusion for the third experiment, Sesia river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1525 Figure S18. Confusion for the third experiment, Kinu river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.

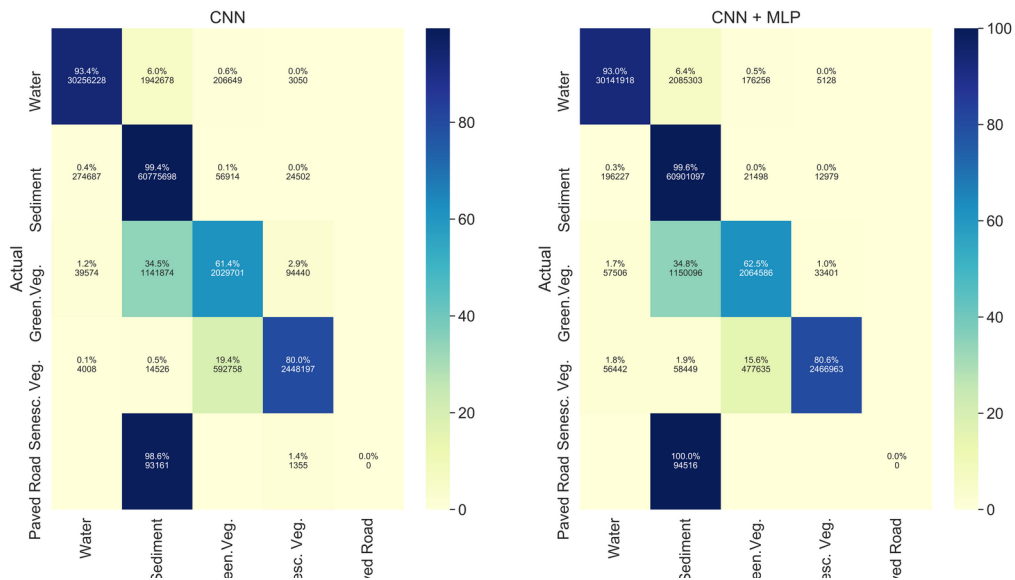
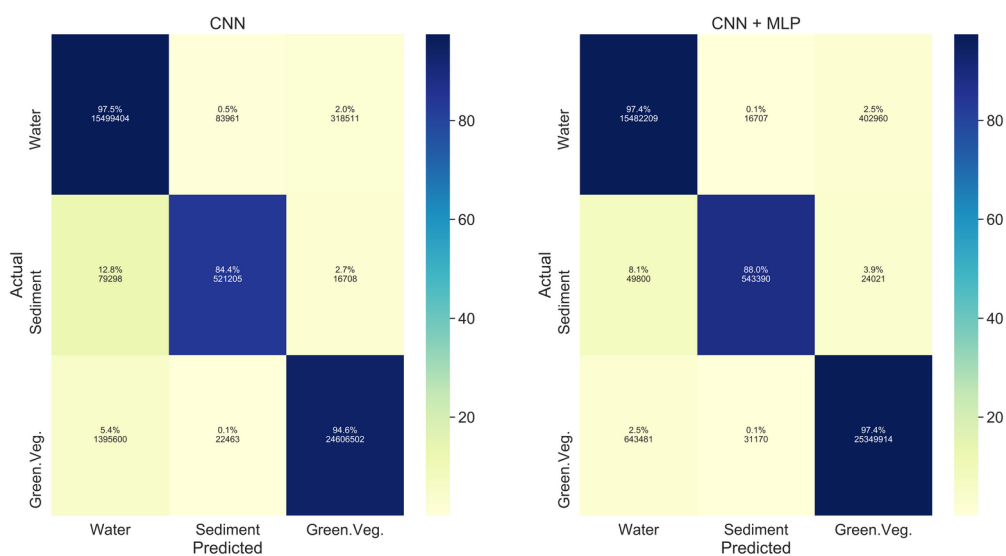
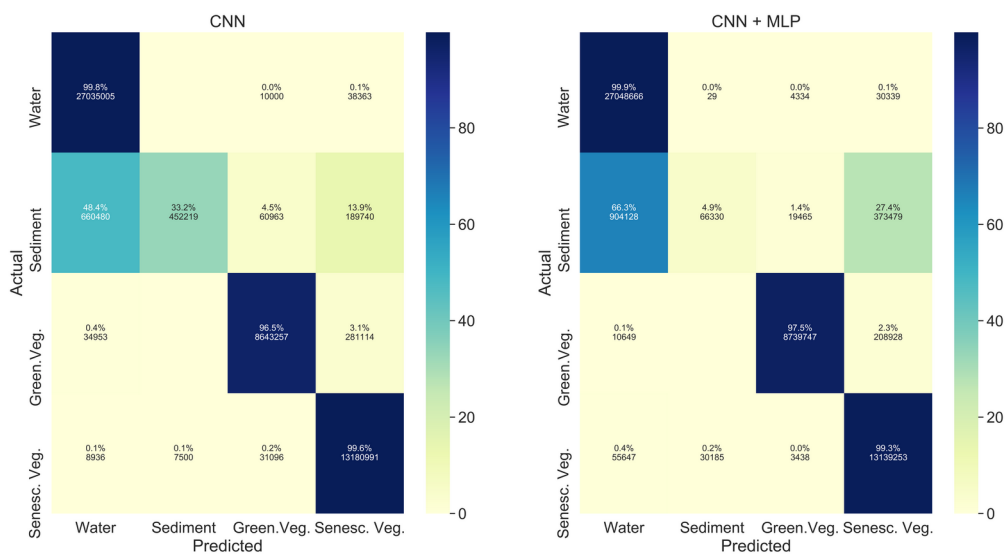


Figure S19. Confusion for the third experiment, Kurobe river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1535 Figure S20. Confusion for the third experiment, Eamont river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.



1540 Figure S21. Confusion for the third experiment, Kingie river. Left) Outcome of the CNN classification of the first CSC phase. Right) Outcome of the second CSC phase CNN+MLP.