

# **Earth and Space Science**



# RESEARCH ARTICLE

10.1029/2021EA002085

#### **Key Points:**

- Methodology for image segmentation based on agreement between a labeler and a machine learning model
- Faster and more accurate segmentation of interpretable imagery compared to traditional labeling
- Large multi-labeler consensus facilitates reproducible scientific inference from Earth surface images

#### **Supporting Information:**

Supporting Information may be found in the online version of this article.

#### Correspondence to:

D. Buscombe, dbuscombe@contractor.usgs.gov

#### Citation:

Buscombe, D., Goldstein, E. B., Sherwood, C. R., Bodine, C., Brown, J. A., Favela, J., et al. (2022). Humanin-the-loop segmentation of Earth surface imagery. *Earth and Space Science*, 9, e2021EA002085. https://doi. org/10.1029/2021EA002085

Received 15 OCT 2021 Accepted 25 JAN 2022

#### **Author Contributions:**

Conceptualization: D. Buscombe, E. B. Goldstein

Data curation: D. Buscombe, E. B.

Goldstein, C. Bodine, J. A. Brown, J. Favela, S. Fitzpatrick, C. J. Kranenburg, J. R. Over, A. C. Ritchie

Formal analysis: D. Buscombe

Funding acquisition: D. Buscombe, C. R. Sherwood, J. A. Warrick

Investigation: D. Buscombe, C. R.

Sherwood, C. Bodine, J. A. Brown **Methodology:** D. Buscombe, E. B. Goldstein, J. A. Brown

**Project Administration:** C. R. Sherwood, J. A. Warrick

© 2022 The Authors. This article has been contributed to by U.S. Government employees and their work is in the public domain in the USA.

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

# **Human-in-the-Loop Segmentation of Earth Surface Imagery**

D. Buscombe<sup>1</sup>, E. B. Goldstein<sup>2</sup>, C. R. Sherwood<sup>3</sup>, C. Bodine<sup>4</sup>, J. A. Brown<sup>5</sup>, J. Favela<sup>6</sup>, S. Fitzpatrick<sup>7</sup>, C. J. Kranenburg<sup>8</sup>, J. R. Over<sup>3</sup>, A. C. Ritchie<sup>9</sup>, J. A. Warrick<sup>9</sup>, and P. Wernette<sup>9</sup>

<sup>1</sup>Marda Science, LLC, Contracted to USGS Pacific Coastal and Marine Science Center, Santa Cruz, CA, USA, <sup>2</sup>Department of Geography, Environment, and Sustainability, University of North Carolina at Greensboro, Greensboro, NC, USA, <sup>3</sup>USGS Woods Hole Coastal and Marine Science Center, Woods Hole, MA, USA, <sup>4</sup>School of Informatics, Computing and Cyber Systems, Northern Arizona University, Flagstaff, AZ, USA, <sup>5</sup>USGS MD-DE-DC Water Science Center, Dover, DE, USA, <sup>6</sup>Department of Earth and Planetary Sciences, University of California Santa Cruz, Santa Cruz, CA, USA, <sup>7</sup>Department of Computer Science, California State University, Sacramento, CA, USA, <sup>8</sup>USGS St. Petersburg Coastal and Marine Science Center, St. Petersburg, FL, USA, <sup>9</sup>USGS Pacific Coastal and Marine Science Center, Santa Cruz, CA, USA

**Abstract** Segmentation, or the classification of pixels (grid cells) in imagery, is ubiquitously applied in the natural sciences. Manual methods are often prohibitively time-consuming, especially those images consisting of small objects and/or significant spatial heterogeneity of colors or textures. Labeling complicated regions of transition that in Earth surface imagery are represented by collections of mixed-pixels, -textures, and -spectral signatures, can be especially error-prone because it is difficult to reliably unmix, identify and delineate consistently. However, the success of supervised machine learning (ML) approaches is entirely dependent on good label data. We describe a fast, semi-automated, method for interactive segmentation of N-dimensional (x, y, N) images into two-dimensional (x, y) label images. It uses human-in-the-loop ML to achieve consensus between the labeler and a model in an iterative workflow. The technique is reproducible; the sequence of decisions made by human labeler and ML algorithms can be encoded to file, so the entire process can be played back and new outputs generated with alternative decisions and/or algorithms. We illustrate the scientific potential of segmentation of imagery of diverse settings and image types using six case studies from river, estuarine, and open coast environments. These photographic and non-photographic imagery consist of 1- and 3-bands on regular and irregular grids ranging from centimeters to tens of meters. We demonstrate high levels of agreement in label images generated by several labelers on the same imagery, and make suggestions to achieve consensus and measure uncertainty, ideal for widespread application in training supervised ML for image segmentation.

Plain Language Summary Labeling pixels in scientific images by hand is time-consuming and error-prone, so we would like to train computers to do that for us. We can use automated techniques from Artificial Intelligence or AI, like one called Deep Learning, but it needs a lot of example images and corresponding labels that have been made by hand. So, we still need to label quite a lot of images at the pixel level—called image segmentation. We made a computer program called Doodler that speeds up the process; you label some pixels, and it labels the rest. It is the fastest method we know of for image segmentation because it is semi-automated. We also show that it produces accurate and precise labeling, as we demonstrated by having multiple people use this method to label the same images. Because it is so fast and accurate, it allows us to get enough data to train Deep Learning models to do segmentation on all the images we have, from the past and in the future. Doodler therefore enables geoscientists to use Artificial Intelligence to extract much more information from their imagery, in service of geoscience in general.

#### 1. Introduction

#### 1.1. The Need for Data Labeling Tools for Earth Surface Processes Research

Automation of data-intensive tasks is increasingly important in Earth surface-processes research. Due to the availability of data at greater spatial and temporal coverages and resolutions (Farr et al., 2007; Gorelick et al., 2017; Wulder et al., 2019), and open-source geo-analytics tools (Richardson et al., 2018; Schwanghart & Scherler, 2014), it is increasingly possible to automate the discovery of patterns in processes operating over complex landscapes (Larsen et al., 2021; Walker et al., 2017). Scoping feasible applications of analytical tools

BUSCOMBE ET AL. 1 of 31



Resources: D. Buscombe
Software: D. Buscombe, E. B. Goldstein
Supervision: J. A. Warrick
Validation: D. Buscombe, C. Bodine, J.
Favela, S. Fitzpatrick, C. J. Kranenburg, J.
R. Over, A. C. Ritchie, P. Wernette
Visualization: D. Buscombe
Writing – original draft: D. Buscombe,
E. B. Goldstein
Writing – review & editing: D.
Buscombe, E. B. Goldstein, C. R.
Sherwood, C. Bodine, J. A. Brown, J. A.
Warrick

such as machine learning (ML) in the geosciences has become a useful way to rapidly explore and prototype ideas with data (Goldstein et al., 2019; Reichstein et al., 2019).

Given the wealth of available ML algorithms in open-access software, geomorphologists have an unprecedented set of available tools for data exploration and hypothesis testing. Machine learning allows us to teach a computer to learn by example, usefully approximating quantities from readily obtainable data that are otherwise hard to sense (Buscombe et al., 2017), parameterize (Beuzen et al., 2019; Ni et al., 2021; Tinoco et al., 2015), flag for quality control (Sugiura & Hosoda, 2020), or to visualize or make automated inference on high-dimensional datasets that a human could not (Chmiel et al., 2021; Plant & Stockdon, 2012), especially for phenomena without well-developed theory (Fox et al., 2015; Goldstein & Coco, 2015). However, the generation of the right type of examples for the machine to learn, or enough of sufficient quality, is a challenge that requires the development of specialist data labeling tools. These tools would allow Earth surface processes researchers to generate their own data representations for training ML to automate cleaning, distillation or classification of content, and make inference, on large geospatial data sets. An example is the segmentation of imagery.

### 1.2. The Need for Better Tools for Image Segmentation

What we hereafter call imagery is considered in the broadest sense as any data set on a regular grid that may or not have a regular spatial footprint, which is collected for scientific applications in the Earth and environmental sciences and in related scientific fields. This definition includes geospatial data sets or rasters, photographic imagery, imagery from satellites, sonar, radar, and other geophysical sensors, and any other gridded data that is visually interpretable (by a subject matter expert or otherwise). Such Earth surface imagery comes in a range of types, from single-band or greyscale commonly created by sensors used in geophysical applications that consist of interpretable textures and edges, to hyperspectral imagery where up to hundreds of coincident bands sense a different narrow portion of the electromagnetic spectrum. We use the term pixels to mean either pixels or voxels, depending on whether the imagery is two- or three-dimensional.

The increasing availability of imagery and increasing acceptance (Olhede & Wolfe, 2018), accessibility (Gil et al., 2016), and sophistication of human-supervised computerized analyses and classification workflows (Cheng et al., 2001; Hossain & Chen, 2019; Mi & Chen, 2020), mean that accurate image segmentation workflows — involving the classification of all pixels in an image —are ubiquitous in need and application in the geosciences (Carleer et al., 2005; Kotaridis & Lazaridou, 2021). Probabilistic segmentation of imagery using ML has various uses in Earth surface processes research (Lang et al., 2019) involving environmental monitoring (Anders et al., 2011; Bayr & Puschmann, 2019; Gaddes et al., 2019; Su et al., 2020). Detection of change in geomorphic studies has traditionally involved differencing of elevation surfaces (James et al., 2012). Segmentation of coincident imagery allows for additional insight, for example, the classification/attribution of the change, evaluation of the agent of change (Grams et al., 2019), the nature and persistence of change, and determination of implications (Barlow et al., 2006; Drăguţ & Eisank, 2012). Understanding these insights is key to habitat monitoring (Chilson et al., 2019; Gray et al., 2019; Ridge et al., 2019) and land use or cover (change) mapping (Buscombe & Ritchie, 2018; Carbonneau et al., 2020; Lefsky, 2010; Pandey et al., 2021) among many other examples (Chaudhary et al., 2019; Ching et al., 2018; Quinn et al., 2018; Weinstein, 2018).

State-of-the-art ML-based image segmentation requires at least some level of human supervision (Kotaridis & Lazaridou, 2021; Sultana et al., 2020). Often the greatest challenge to developing an automated workflow can be the creation of model training data that is internally consistent (Serre, 2019). In the case of image segmentation, training data consists of label imagery where each pixel is categorized into any number of pre-determined discrete nominal or ordinal classes. Many applications of segmentation of Earth surface imagery by definition are concerned with surfaces, therefore the focus of many labeling workflows, and also the present contribution, is the generation of 2D label images by segmenting visually interpretable imagery, that is, up to three coincident bands.

Such label imagery is typically acquired by either hand-digitizing vector polygons that are subsequently rasterized (Kotaridis & Lazaridou, 2021), or raster editing, which is hand classification of pixels directly. Creating label images through digitization of hand-drawn polygons is time-consuming; raster-editing can offer a quicker alternative, and most commercial and non-commercial image-editing software also have built-in tools that can select entire regions via similar colors or edge-detection techniques. These tools are typically (a) not reproducible because the outputs are generated by a sequence of clicks that are not recorded in a file (a fact that precludes

BUSCOMBE ET AL. 2 of 31



many of the analyses of multi-labeler agreement we present here), (b) proprietary or restrictively licensed, and/ or c still require significant amounts of time and effort to achieve good results. The largest error is at boundaries between classes, and arises due to two factors: (a) indistinct areas of transition where it is not always possible to make an objective decision about the class, and (b) it is almost never feasible to click the shape of a polygon outline at the pixel level.

Labeling Earth surface imagery using these traditional methods is especially time-consuming if images consist of small or unfamiliar objects and/or colors or textures exhibiting significant spatial heterogeneity and/or ambiguity, necessitating a high zoom level, or viewing at a range of scales. Moreover, labeling transition regions is difficult to do reliably because of mixed-pixels, -textures, and -spectral signatures, which can lead to significant amounts of error. Earth surface imagery is more likely to have these properties than much imagery used to develop image segmentation models, labeling tools, and benchmark datasets in ML research and applications (Everingham et al., 2010). In Earth surface imagery, and especially in transition areas, we argue that pixelwise classification needs a human for these transition regions and more complex textures, but could also be sped up by including techniques that aid the human labeler, such as ML models that are trained as a human annotates.

#### 1.3. Human-in-the-Loop Image Segmentation

Here, we describe and evaluate a so-called "human-in-the-loop" (Monarch, 2021) machine learning workflow for fast image segmentation, encoded in a computer program called Doodler, and we demonstrate its use for geophysical, photographic, and multispectral satellite images of natural environments. Doodler lies on the spectrum of what Monarch (2021) refers to as "assisted annotation," which is interaction with raw data, with ML assisting the data labeling process, and "predictive annotation," where ML generates outputs that can be edited. In fact, the program essentially does both, in a loop whose number of iterations for any given sample is dictated by the human labeler who acts to assure data quality.

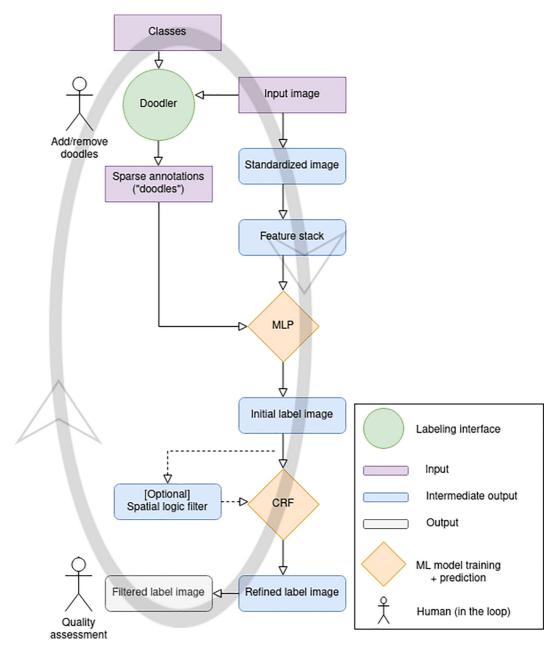
As supervised ML workflows gain popularity in the geosciences (Bergen et al., 2019; Zuo et al., 2019) and related fields (Crisci et al., 2012; Kashinath, Mustafa, et al., 2021), Doodler could be used in numerous contexts to reach a target ML model accuracy by training on large amounts of data acquired relatively quickly. It also serves as a case study in how to combine human and machine intelligence to label scientific data with increased efficiency and accuracy. In section two we introduce the human-in-the-loop labeling principles and graphical (in the sense of Koller and Friedman, 2009 of models consisting of nodes connected by vertices) model framework, followed by a description of the image feature-extraction methods, and the ML classifier. In Section 3 we describe six datasets that we use to demonstrate the approach. These are chosen to quantify and discuss variability among label images made by several independent labelers, and further to examine variability in image segmentation outputs due to image size and resolution.

Comparisons between images labeled by the same labeler at different scales, and multiple labelers of the same imagery are presented in Section 4. This section serves a few purposes. First, for subjective tasks involving interpretation of ambiguous data, or even objective tasks or relatively simple tasks where random human blunder may be a factor, no simple heuristics exist for deciding the correct label (Monarch, 2021) however some practical recommendations can be made using statistical metrics of multi-labeled datasets (Goldstein et al., 2021). Similarly, we offer some methods for identifying and quantifying uncertainty based on agreement over segmentations of the same imagery by multiple labelers. Second, this section serves to demonstrate that the methodology and implementation we present are reproducible between labelers, at different times, and using different computational infrastructure (computers, browsers, etc.), despite the fact that the label image is a model estimate from sparse annotations that would vary considerably from labeler to labeler. In Section 5 we make suggestions on how to achieve consensus and measure error, and recommendations over usage of the Doodler program, before drawing conclusions.

# 2. Human-in-the-Loop Labeling Using Machine Learning

The image labeling task (Figure 1) involves a human labeler providing sparse annotations (informally called "doodles") to inform and automate a process ("model") that estimates the label for all pixels in that image, then the same labeler refine the model predictions using a combination of adding/removing doodles and/or changing model hyperparameter values. A workable system necessitates a graphical user interface and a fast and accurate

BUSCOMBE ET AL. 3 of 31



**Figure 1.** Schematic of the approach encoded into the Doodler program. Most images-class set pairings trialed to date have been segmented successfully within one or two loops. Doodler also facilitates the user to modify the model hyperparameters that may be used iteratively the same way as adding or removing annotations ("doodles"). The human adjusts the hyperparameters and they feed into the Multilayer Perceptron and Conditional Random Field models.

image segmentation process. Each image is classified according to a set of pre-determined classes; we use the term label to refer to a single instance of an annotation of a specific class, such that each class present in every image is exemplified with numerous labels.

The images are segmented semi-interactively, one-by-one, so there is no need to specify an underlying prior statistical model, and we need not assume pixel values are conditionally independent of a given label. Therefore ML is ideally suited to the task; because it could learn how to map the features that may be readily extracted from imagery, to class labels, from a small proportion of labeled pixels. That model could then be used to estimate the class of the remaining pixels not labeled. More formally, we use a discriminative ML model, f, that has learned the conditional distribution  $P(y|\theta, x)$  directly, which reads as the probability of y, given  $\theta$  and x, where x are the

BUSCOMBE ET AL. 4 of 31



image features associated with annotated pixels y, and  $\theta$  are learned parameters. This approach is highly suited to task-specific prediction such as here; the models need not be portable among images, therefore no attempt is made to capture the distributions over x or model the correlations among x. The model then predicts the class  $\hat{y}$  of the unlabeled pixels  $\hat{x}$  by  $\hat{y} = f(\hat{x})$ , essentially by assuming  $P(\hat{y}|\theta,\hat{x}) \approx P(y|\theta,x)$ .

The system consists of (a) a human annotator providing sparse examples of each class of interest in a graphical user interface running in a web browser, (b) a Multi-Layer Perceptron (MLP) model (Bishop, 2006) for per-pixel class, based on a probabilistic model of how classes relate to a stack features extracted from standardized imagery based on intensity, texture, edges, and relative location, controlled by parameters learned during a discrete training period, and (c) a graphical model called a fully connected conditional random field (CRF; Kumar & Hebert, 2006) that refines estimates of the per-pixel class based on a probabilistic assessment of how classes relate to features extracted from imagery based on both color (if three or more dimensions) or intensity (if imagery is 2D) and relative location, controlled by hyperparameters set/tuned by the human labeler, who also acts to assess quality, and iterate as needed. We use the two ML models in conjunction with human annotations to classify each pixel of the scene and segment the image. At least one of the classes must exist in a given image, but otherwise there are no restrictions on the number of classes (other than practical considerations such as available time). Often models need the most detailed annotations or "doodles" near the boundaries where one class transitions to another.

The program facilitates human labeling, which also provides quality control. In effect, the labeler interacts with a machine to collectively decide on the most accurate and precise label image for any given image. Doodles are used to update the ML model iteratively, by adding/removing annotations, and also optionally changing hyperparameters for optimal segmentation on individual images, and retraining and implementing the model. The program relies on the labeler having the patience, dedication, and interest to do a good job, which may require a few iterations of the workflow (Figure 1). The design of the program would also be amenable to labeling in stages, with each stage perhaps employing people with different levels of expertise. We now describe the two ML models embedded within the doodler workflow, namely the Conditional Random Field (2.1) that uses a Multilayer Perceptron or MLP (2.3) as a sub-component. We conclude by describing our implementation of the Doodler workflow in 2.4.

The methodology not only facilitates much faster segmentation, which makes multiple labeler datasets more obtainable (affordable, and completed in a reasonable time), but also results in more accurate segmentations. That is because the labeler is asked only to provide true and unambiguous positive examples of each class. Errors at boundaries between classes that arise due to hand digitization, which can be significant because of mixed pixels or due to coarse digitization, are significantly reduced. That is because the program predicts at the pixel level much faster than a human could ever label at that scale, and also because our approach models the likelihood of uncertain regions. The latter is crucially important for class assignment in particularly difficult regions of imagery in a deterministic manner.

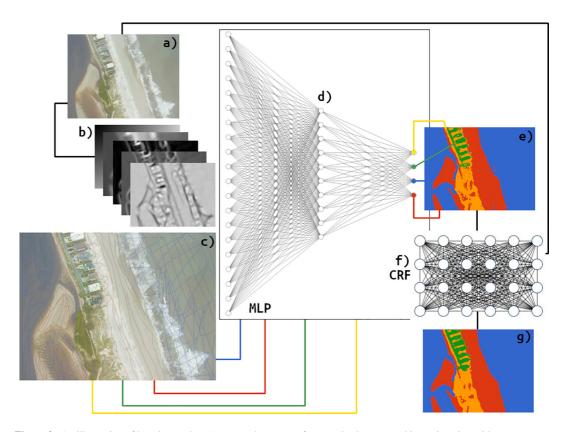
# 2.1. Conditional Random Field for Image Segmentation

We adopt a widely used approach to such task-specific probabilistic image segmentation, which is a Conditional Random Field or CRF model (Kumar & Hebert, 2006; Vosselman et al., 2017; Zhong et al., 2014) to estimate per-pixel class likelihoods (Figure 2). We use the similar CRF implementation of Krähenbühl and Koltun (2011) that was previously used by Buscombe and Ritchie (2018). Whereas Buscombe and Ritchie (2018) used a trained convolutional network to label regions of images that were used as unary potentials for a CRF model for pixel-level refinement, and Buscombe and Grams (2018) used sparse instrumental observations from the field in conjunction with geospatial imagery, here (Figure 2) labels of some regions of images are provided by humans, which are used to ascribe a probability of each class per pixel using a Multilayer Perceptron. Those outputs (per-pixel class likelihoods) are used as unary potentials for a CRF model for pixel-level refinement; the CRF model additionally models the joint likelihood of each pair of pixels, essentially checking for internal consistency of the MLP outputs.

The unary potentials define a log-likelihood over the label assignment y, and therefore represent the cost of assigning label  $y_i$  to grid node i. They are called "unary" potentials because they describe feature-class relations at every pixel, and to distinguish them from pairwise potentials, dependent on feature-class relations over pairs

BUSCOMBE ET AL. 5 of 31





**Figure 2.** An illustration of how image data (a) are used to extract features (b) that are used in conjunction with sparse annotations (c) to train an initial Multilayer Perceptron classifier (d) to extract unary potentials (e) that are refined by a Conditional Random Field (f) to create a refined label image (g).

of feature-class relations, which are also used in the CRF model and defined later. Here we use a Multilayer Perceptron (Bishop, 2006) as a classifier to generate unary potentials. In CRFs based on "local" connectivity, nodes connect adjacent pixels in x (Kumar & Hebert, 2006), whereas in the fully connected definition such as here (Figure 2f), each node is linked to every other (Krähenbühl & Koltun, 2011). Linking each node of the graph created from x to every other enables modeling of the long-range connections within the data by considering both proximal and distal pairs of nodes, resulting in refined labeling at boundaries and transitions between different classes. We use a global probability prior  $p_u$  of the unary potentials, that is, a prior probability that any random sample correctly labels the underlying image features. It is exposed to the user as a seldom-varied hyperparameter, defaults to 0.9, and generally has limited effect unless provided annotations are actually of poor quality, which we assume is rarely the case.

There are two non-dimensional hyperparameters exposed to labelers using the Doodler program. The first is  $\theta_{\beta}$  (default = 1) is used by the CRF feature extractor to extract color image features and map them to classes. These features are engineered, by convolving Gaussian kernels with the imagery (in much the same way as features are extracted as inputs to the MLP model –see Section 2.2). Hyperparameter  $\theta_{\beta}$  controls the degree of allowable similarity in image features among classes, therefore  $\theta_{\beta}=1$  only tolerates image features with small differences in intensity being assigned the same class label.

The second hyperparameter,  $\mu$ , is used within a Potts label "compatibility" function (Krähenbühl & Koltun, 2011) to define pairwise potentials used by the model to encourage adjacent pixels to be the same class label, defined as  $\Lambda(i,j) = \mu$  if i = j and 0 otherwise. By default, Doodler uses  $\mu = 1$ , meaning  $\Lambda$  is simply a  $k \times k$  identity matrix, whereby all classes are equally "compatible" (as likely as each other to be adjacent in either image or feature space). Values greater than 1 weight the pairwise potentials more than the unary potentials, which might be useful when the MLP prediction is poor, in which case the pairwise potentials count by a factor of  $\mu$  greater than the unary potentials.

BUSCOMBE ET AL. 6 of 31



By definition,  $\theta_{\beta}$  and  $\mu$  are task-specific, so their respective effects are hard to generalize, but it can be said that, in general, larger values of  $\mu$  tend to give the model greater independence, resulting in the reclassification of more pixels. The importance of pairwise potentials becomes much greater than unary potentials, and spatial inconsistencies in feature-label pairings have greater likelihood of being reclassified. In general,  $\theta_{\beta}$  has a more muted effect and generally controls the sharpness of the class boundaries in the label image. Note that neither effect necessarily improves the result. Please refer to Figure S1 in Supporting Information S1 for visualizations of the effects of varying  $\theta_{\beta}$  and  $\mu$  on sample imagery from the Sandwich data set, expressed in terms of where the labels of pixels are altered by the CRF compared to the MLP output. The reader is also referred to the Supporting Information S1 section entitled "Fully Connected Conditional Random Field for Image Segmentation" for more technical details about its implementation and interpretation of parameters.

By design, the CRF solution is not overly sensitive to hyperparameter values. First, imagery is standardized therefore the model does not need to use parameters for brightness (related to non-zero image mean) and contrast (related to non-unit image variance). Second, we use spatial logic to filter CRF inputs, which eliminates a major source of uncertainty for the CRF solution employing pairwise potentials, because the CRF model will be given more consistent spatial pairs of feature-class-pairings to make inference from. Finally, hyperparameter sensitivity increases if the sparse annotations are used alone (Buscombe & Grams, 2018), and/or if the unary potentials estimated by the MLP model are spatially sparse (Buscombe & Ritchie, 2018).

#### 2.2. Image Standardization and Feature Extraction

Each input image, I(i, j, d), where i and j describe 2D pixel locations and d indicates the number of coincident data layers, is standardized such that it has zero mean and unit variance (see Supporting Information S1 section entitled "Image Standardization and Feature Extraction"). This ensures the values are distributed within the range -1 and 1, which helps numerical stability and builds insensitivity to outliers, as well as removing any bias from any channel as a function of the mean image intensity.

Raw pixel values are not used as inputs to the MLP classification model described in Section 2.3. Instead, features are extracted in a prescribed way that is, the image features are extracted in the same way each time, known as feature engineering. Features relating to image intensity, edges, texture, and relative location are extracted, all at a range of scales. Then a stack of features are provided to the classifier. We use kernel convolution methods for feature extraction because they are already common in numerous geophysical applications concerning interpretation and quantification of spatially distributed imagery. Image intensity features  $I_f(i, j)$  are extracted from  $I_s(i, j, d)$  by convolving with filter bank  $\Sigma_s$ , or  $I_f = \Sigma_s * I_s$  where \* denotes convolution, and where  $\Sigma_s$  consists of s 2D Gaussian kernels.

Edge features are extracted using the Sobel operator, computing an approximation of the gradient magnitude of  $I_p$ ,  $\nabla_{I_f}(i,j)$ . Location is encoded as the kernel-convolved bank of 2D features given by  $L(i,j) = \Sigma_s * \sqrt{(i^2 + j^2)}$ . Finally, texture features are computed as the first and second eigenvalues of Hessian matrix of  $I_f(i,j)$ , or  $H_1(i,j)$  and  $H_2(i,j)$ . Eigenvalue analysis of the Hessian is commonly used in geophysical and medical image feature-extraction (Bishop, 2006) because of its formalized relationship to physical quantities, extracting the principal directions in which the local second order structure of the image, that is, its spatial covariance structure, can be decomposed. The eigenvectors and eigenvalues of the Hessian are known as principal directions and principal curvatures respectively (Koenderink & Van Doorn, 1992). The first two eigenvalues are the magnitudes of the maximum and minimum curvature, respectively.

# 2.3. Initial Segmentation Using a Multilayer Perceptron

The feature stack used for initial segmentation consists of a set of 3D (i, j, d) grids, each flattened to 1D (1, ijd), then stacked columnwise to create a model input vector. The feature stack is then subsampled row-wise by a factor defined by the user. For larger imagery, this subsampling factor may be as large as six, but typically it is one (i.e., no subsampling) to three, and depends on the available computer memory and processing time.

Our entire model framework implementation (see Supporting Information S1 section entitled "Multilayer Perceptron") consists of an input layer of *ijd* neurons, two hidden layers, the first consisting of 100 neurons and the

BUSCOMBE ET AL. 7 of 31



second of 60 neurons, each linked to each other (i.e., fully connected), and finally a classifying layer consisting of k neurons, where k is the set of classes with labels, that is, present in the scene, determined a priori for the scene. Through extensive experimentation, we are satisfied that model outputs are not overly sensitive to the specification of the number of neurons in each of the two hidden layers. However, hidden layers or neurons could be added for greater discriminative power at the expense of model parsimony and computational efficiency. MLPs have previously been successfully used for Earth surface image segmentation (Kurnaz et al., 2005; Villmann et al., 2003), as have other types of artificial neural networks (Buscombe & Ritchie, 2018; Kemker et al., 2018).

While any number of similar deterministic ML algorithms could have been used, MLPs are attractive due to their relative simplicity and longevity which has created a widespread use of them among many geoscience and related fields (Gardner & Dorling, 1998). Because this is task-specific prediction, 90% of the input feature data are used for training and only 10% for validation, for iteratively adjusting **w** and **b** through back-propagation and solved using stochastic gradient descent, with a maximum of 2,000 training epochs. We use the Adam stochastic gradient-based optimizer method proposed by Kingma and Ba (2014), using early stopping to terminate training when the validation score does not improve by 1e<sup>-4</sup> for at least 10 consecutive training epochs. Model outputs are not very sensitive to hyperparameters, that is, choices about percentage of data used for validation, number of training epochs, or criteria for terminating the training (see for example, Figure S1 in Supporting Information S1). For brevity, this sensitivity analysis is not shown here but the program documentation explains where these hyperparameters may be adjusted and their resulting outputs compared.

Whereas there is no drop-in replacement for the CRF, the MLP could be switched to a different ML framework. In fact, we have also extensively trialled a Random Forest model framework but decided that the MLP performed better; see Figure S2 in Supporting Information S1 for an example, based on data set A.

#### 2.4. Implementation: The Doodler Program

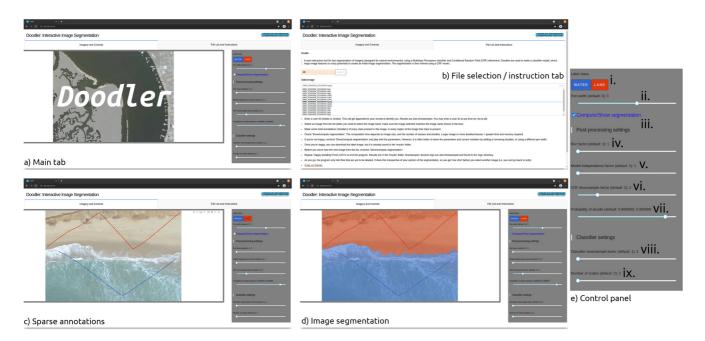
In a human-in-the-loop data labeling system, the design of the front-end annotation interface is as or more important than the back-end ML model framework. At a minimum, the user interface must allow for image annotation and a mechanism for launching the image segmentation process (Figure 3). Optionally, it can also expose controls to facilitate image curation and class (label) definition, mechanisms to adjust hyperparameters, and controls for re-segmentation. We have created several versions of the program, including some that store images locally, and others that retrieve imagery from a remote server. The latter case is useful for collaborative labeling projects, because the application can be hosted on the worldwide web and the results can be stored centrally.

The default version of the program that we have made publicly available allows the user to place images for classification in a local "assets" folder. The program tracks images that have been classified, therefore the list of files available for classification gets smaller during a labeling session. Users can also modify hyperparameters and redo segmentations as many times as desired, as well as the "pen" width (width in pixels to ascribe each annotation). These controls can optionally be hidden from the user in order to only collect the sparse annotations, and/or (pixelwise) label images with a fixed set of default hyperparameters.

Each MLP prediction is a matrix of dimension ijk encoding the probabilities of each pixel i, j and each class k. The discrete class is found as the maximum over i, j in the k dimension, or argmax, resulting in a label matrix of integer values, each integer corresponding to a unique class. Often there can be high-frequency noise in the resulting 2D discrete label image of pixelwise predictions, that is, small islands of misclassified pixels. Since classifier outputs are probabilistic, instead of using argmax we could choose to filter these islands based on logic or some other process operating on the probabilities themselves, or we could filter islands by operating in the spatial domain on the label image. Doodler implements the latter, using two complementary filtering procedures. Therefore we implement an additional, but optional, step is performed in which the label matrix output from the MLP model is spatially filtered. The filtered label is then used as input to the CRF model. The reader is referred to Supporting Information S1 section "Spatial Filtering of Initial Segmentation" for more details. An illustration of the full workflow described in Section 2.1 through to the present section, including the spatial filtering of the initial segmentation, is presented as Figure S3 in Supporting Information S1.

BUSCOMBE ET AL. 8 of 31





**Figure 3.** The graphical user interface of the program Doodler for a simple two-class (water and land) labeling task. (a) Initial view of the primary interface tab; (b) view of the second tab showing the (optional) input for name or identifier that gets appended to every output filename (to help identify labelers in multi-labeler trials such as presented here), and the drop-down list of image filenames yet to label (lists are cross-checked every second or some user-defined amount); (c) view of the primary tab with doodles; (d) view of the result of the initial segmentation; and (e) detailed view of the control panel. The control panel shows classes as different color buttons (i), pen width (ii), a checkbox for computing the image segmentation (iii) (exposed) hyperparameters that relate to the Conditional Random Field (CRF) (iv:  $\theta_p$ ; v:  $\mu$ ; vi: CRF downsample factor; vii: the prior global probability of the unary potentials,  $p_u$ ), and (exposed) hyperparameters that relate to the feature extraction (viii: feature downsample factor; ix: number of scales over which to extract features). The feature downsample factor downsamples the entire feature stack before being classified with an MLP, and the resulting outputs are upsampled using nearest-neighbor interpolation back to the original size. The CRF downsample factor downsamples the 3D one-hot encoded stack of unary potentials that result from the MLP, before being used as input to the CRF model, and again the resulting CRF outputs are upsampled using nearest-neighbor interpolation.

# 2.5. Comparison of Segmentations

In order to quantify inter-labeler differences, the canonical metric to evaluate the difference between two thematic maps or label images (Costa et al., 2018) is the mean Intersection over Union score (*IOU*, or Jaccard Index) averaged over *k* classes. For a collection of overlapping regular shapes, an IOU value of 0.5 would imply average overlapping by 50%, but in this context contributions are summed over fields, therefore when IOU reflects 50% average overlap between each contiguous region of labeled pixels. However, many label images are class-imbalanced, which is to say there tends to be a majority class and one or more minority classes.

The mean Dice score is relatively insensitive to the number of pixels total in each class, because the numerator is the number of correctly classified pixels, and the denominator is the total number of pixels in a class that is in both estimated and observed. It has therefore been suggested to be a more accurate metric for the overall agreement between two label images for class-imbalanced label images, whereas an IOU score is not as sensitive to contributions from the smaller class (Csurka et al., 2004). The reader is referred to Figure S4 in Supporting Information S1 for the functional relationship between mean Dice and mean Intersection over Union, and to Figure S5 in Supporting Information S1 for an illustration of the behavior of these metrics for a sample comparison using one data set, and to the in Supporting Information S1 section entitled "Comparison of Segmentations" for mathematical details about the two metrics.

For both IOU and Dice scores, where different numbers of unique classes exist, that is, two different candidates for k, we could choose to set k as the minimum number of the two respective class sets, or the maximum number. We chose the maximum, therefore scores are conservative in these situations. It might be surmised that Dice measures average accuracy, while IOU measures something closer to the worst-case accuracy. However, they vary nonlinearly and, due to averaging over classes, exhibit independently useful properties. We present both scores for each data set, and also use them to discuss ways to detect class imbalance, outlier labelers, and label

BUSCOMBE ET AL. 9 of 31



Table 1 Case Study Datasets					
Data set	Name	Type	Classes	Labelers	Source
A	Beach sedimentology	Orthomosaic	6	2	U.S. Geological Survey data release <sup>a</sup> .
					Sherwood et al. (2021)
В	Post-hurricane assessment	Oblique aerial	4	2	National Geodetic Survey emergency response imagery <sup>b</sup>
C	Shoreline identification	Nadir aerial	4	5	U.S Geological Survey data release <sup>c</sup>
					Kranenburg et al. (2020)
D	Riverbed structure	Sidescan	9	2	Used with permission from U.S. Fish and Wildlife Service
E	Barrier breach	False-color satellite	7	2	Sentinel-2 imagery courtesy of European Space Agency (ESA)
F	Coastal evolution	Visible-band satellite	4	1	Landsat-8 imagery courtesy of U.S. Geological Survey

 $^a https://doi.org/10.5066/P9BFD3YH.\ ^b https://storms.ngs.noaa.gov.\ ^c https://doi.org/10.5066/P9CA3D8P.$ 

images in multi-labeler contexts, as well as reporting mean agreement for multi-labeled datasets as an uncertainty and quality metric.

# 3. Datasets and Case Studies

We demonstrate our approach using several case studies from riverine, estuarine, and coastal environments of the United States, chosen to illustrate the scientific potential of image segmentation in diverse environments and image types, and more specifically to quantify inter-labeler-agreement under various contexts. The datasets (Table 1) consist of one- and three-band imagery on regular and irregular grids ranging from centimeters to tens of meters, including photographic and non-photographic imagery. Segmentation of this imagery can be used to answer a range of scientific questions concerning landscape change, which we exemplify for each data set below. In each case, the labelers were issued instructions only verbally, rather than demonstrating with examples. The task was discussed, then attempted once and not redone.

# 3.1. Sedimentary Mapping of a Mixed-Sand-Gravel Beach From Visible-Band Aerial Orthomosaic Imagery

Data set A (Sherwood et al., 2021) consists of one, three-band orthomosaic image (Figure 4a, Table 1), at 5-cm and also downsampled to a resolution of 25-cm, for mapping beach substrates of Sandwich Town Neck Beach on Cape Cod, Massachusetts. The orthomosaics are created from photographs collected from a low-altitude Uncrewed Aircraft System (UAS) on 21 September 2016, using a structure-from-motion workflow similar to that described by Over et al. (2021) for high-resolution elevation mapping of coasts from aerial imagery (Warrick et al., 2019). The 5-cm and 25-cm pixel imagery are divided into  $1,024 \times 1,024$  pixel, 3-band (RGB) tiles for annotation, which results in 99 and six tiles for the respective resolutions. The two data sets were labeled by different individuals. The reader is referred to Figures S1, S2, and S3 in Supporting Information S1 for more example imagery. The following categories are used; (a) water, (b) sand, (c) gravel, (d) cobble/boulder, (e) vegetated, (f) development.

The orthomosaics are used to evaluate the products resulting from labeling images at two resolutions. They are also used to illustrate how to determine optimal image and pixel size for annotation. Such imagery is used for tracking changes to the beach morphology and sedimentology, such as tracking the position of the shoreline, berm, and scarp to indicate the nature of morphological change, as well as individual sediment fractions such as gravel patches that may have a morphodynamic role or could be sensitive coastal state indicators. Segmentation is also useful for determining which parts of the scene are useable data for subsequent analyses. In some situations when working with large imagery, it is difficult to know a priori what image size to use when annotating using the methods described here; while the program facilitates zooming and panning (see Section 2.4 for details on our

BUSCOMBE ET AL. 10 of 31



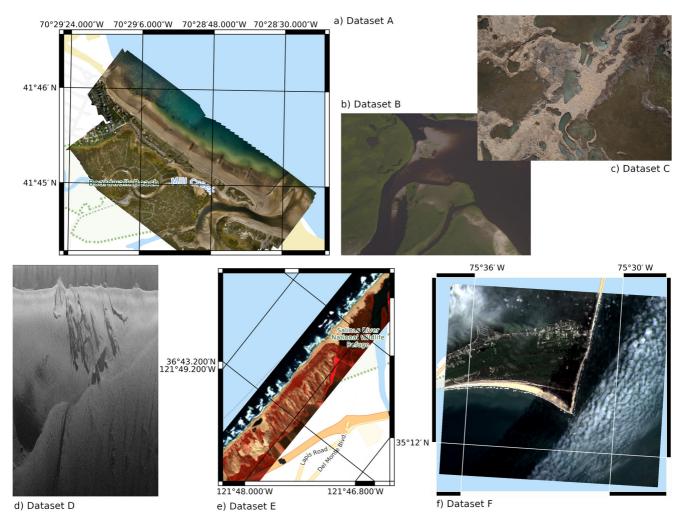


Figure 4. One example image from each of the six data sets used in this study, from left to right; (a) a portion of an orthomosaic image of a beach, (b) an aerial image of a marsh environment, (c) an aerial image of a backbarrier coastal dune environment, (d) a portion of a sidescan echogram from a coastal plain river, and (e) a false-color multispectral satellite image of a coastal lagoon and vicinity. Geospatial imagery on regular grids are shown with latitude and longitude grids and labels.

program implementation), sometimes it is more efficient to use smaller image tiles. In other situations, there is a choice over what grid size to use when making the imagery, such as when converting from ungridded to gridded data. The orthomosaics are created from color-attributed 3D point clouds (Over et al., 2021), therefore we use data set A (Table 1) to discuss a workflow designed to experimentally determine optimal grid size and image size ahead of a large labeling task.

# 3.2. Flood Detection in Post-Hurricane Aerial Photographic Imagery

Data set B (Figure 4b and Table 1) consists of a non-continuous spatial series of 80, three-band image tiles  $(1,000 \times 750 \times 3 \text{ pixels})$ , which are from Emergency Response Imagery collected by the National Geodetic Survey Remote Sensing Division of the US National Oceanographic and Atmospheric Administration, NOAA (NOAA, 2021), that have been each divided into four tiles. The imagery is from North and South Carolina taken after Hurricane Florence (2018). Post storm imagery can be used to monitor the effects of hurricanes on coastal communities (Chen et al., 2018) and ecosystems (Barnard et al., 2021) and coastal change (Goldstein et al., 2020). The images are labeled using the following classes: (a) water, (b) sand, (c) vegetated surface, and (d) development. We compare the segmentations from two labelers labeling the same complex imagery that is readily interpretable without specialist knowledge, but nevertheless difficult to interpret all classes consistently. The reader is referred to Figure S6 in Supporting Information S1 for more example imagery.

BUSCOMBE ET AL. 11 of 31



#### 3.3. Delineating Land From Water in Intertidal Areas of Aerial Photographic Imagery

Data set C (Figure 4c and Table 1) consists of a series of 10, three-band arbitrary images of shoreline environments such as could be collected from a low-altitude aircraft in numerous locations, each labeled by five people using the following four classes; (a) deep water, (b) whitewater, (c) intertidal area (including all visibly shallow water where the surface below the water is visible, swash regions, and wet sand), and (d) dry land. The reader is referred to Figure S7 in Supporting Information S1 that depicts all 10 images. Such imagery is useful for basic monitoring and photogrammetric reconstruction of shoreline environments.

Five labelers examined the same complex imagery that is readily interpretable without specialist knowledge but like data set B, is not necessarily straightforward to consistently interpret. It is a complex labeling task involving identification and lumping of intertidal areas of what are in fact two distinct classes, namely wet sand and shallow water, into a single "shallow" class. The task is made even more complex by asking the labelers to distinguish between that shallow class and "water", a subjective choice requiring identification of water that is deep enough so as not to be confused with shallow water through which the underlying surface is visible. On this occasion, the labeling team of five people discussed the challenges of reliably distinguishing among these four classes beforehand, and this labeling exercise was to determine the utility of the class set before a larger labeling exercise was conducted.

#### 3.4. Benthic Physical Habitat Mapping in Sidescan Sonar Data

Data set D (Figure 4d and Table 1) consists of a non-continuous spatial series of 51, one-band (greyscale) image tiles, each a short section of port or starboard scan consisting of 1,024 consecutive sonar pings stacked as image columns. The length of each ping varied due to sonar range, resulting in the number of image rows varying between 1,300 and 2,000 pixels. The scans are collected using a Humminbird® Solix sidescan sonar emitting a frequency modulated sound pulse with a nominal carrier frequency of 1.2 MHz, from sections of the Pearl River and its tributary the Bogue Chitto, and from the Chickasawhay, Buoy and Leaf tributaries of the Pascagoula River, in Spring 2021, for mapping in-stream physical habitats in coastal plain rivers of Louisiana and Mississippi. Data set D (Table 1) consists of 10 example scans from the Bogue Chitto River, four from the Buoy River, two from the Chickasawhay River, 12 from the Leaf River, and the remaining 23 from the main stem Pearl River. The samples are selected for a variety of substrate types, water depths, and turbidities. Data are decoded and processed following Buscombe (2017). The reader is referred to Figure S8 in Supporting Information S1 for more example images.

In these so-called "waterfall" images, the pixels represent acoustic backscatter intensity (brighter = higher intensity) of the 80-ms pulse, mapped in a non-linear coordinate system representing two-way travel time on the y-axis, and pulse number on the x-axis. Because the transducer moves, pulse number corresponds to along-track distance, but the scale varies with boat and current speed. The top portion of the y-axis records backscatter from the water column and represents a nearly vertical domain between the transducer and the river bed. The lower portion records backscatter from the river bed at increasing distances from nadir. As the distance increases (lower in the images), the sound-path angle of incidence increases, changing the distance scale. The pixels representing the water column are oriented perpendicular to the bed, and the remaining pixels representing the riverbed and shadows in the lee of the bed and other objects. The water column pixels are therefore 2D (x, z) and the remaining pixels are 2D (x, z) representations of the 3D (x, z) bed relief; objects on the bed cast shadows in their lee, the length of which depends on the geometry of the object with respect to the sonar (Buscombe et al., 2016). The length of each ping is variable, depending on the characteristics of the sound pulse that collectively determine range, and the fact that the amount of useable data also varies strongly across-track (the vertical image dimension) due to attenuation of sound by water and the bed (Buscombe, 2017).

For these data, we found it preferable to annotate the waterfall imagery rather than the georectified mosaic images, which would account for heading, slant-range, and time-varying gain. Visual identification of features is easier and more consistent when applied to the unrectified scans rather than the georectified mosaics, because they are collected in fast, shallow water, boat speed is variable, course is difficult to maintain exactly, and navigating meanders can obscure features. Time-varying gain is minimal because the water depths are much less than 10 m. Therefore we have found it preferable, in this case, to segment the unrectified imagery, and the subsequent label images would then be georectified.

BUSCOMBE ET AL. 12 of 31



Like many scientific images, there are unusable portions of the imagery that would need to be removed through classification and removal by an automated process; in this case, they are the bank shadows and water classes, because the others are mappable in 2D space. There are many low-signal-to-noise (dark, grainy) textures that at small scale are not distinguishable without some spatial context - such as water, and shadows cast by variously sized objects. The full class list is as follows: (a) water; (b) shadow/riverbank; (c) shadows cast by instream objects and morphologies; (d) submerged wood; (e) fine sediment bedforms; (f) flat, fine sediment; (g) coarse sediment (gravel through boulders), bedrock, and vegetation; (h) anthropogenic (human-made objects); and (i) unknown (rare blank regions where the sonar recording cut out). Of the above, all but "anthropogenic" are present in the data set used for this study.

Such imagery is used to compare the products resulting from two labelers annotating the same complex imagery requiring specialist interpretation. Such imagery is used for mapping riverbed sediments (Buscombe, 2017; Buscombe et al., 2016) to provide basic information for benthic habitat mapping, and morphodynamic and sediment transport studies in rivers. It is also an example of a geophysical data set with features in common with other Earth surface imagery, such as slices from 3D tomography data, Synthetic Aperture Radar (SAR), multibeam sonar backscatter, seismic reflection and refraction, to name but a few. The sidescan data set requires the most training and expertise to interpret. It is the only data set used here that is actively sensed (using an emitted sound wave and recording the echo).

Other than the false-color satellite imagery, this sidescan imagery is the only data set that requires specialist knowledge to even sensibly interpret. Those data are therefore labeled by two experts with extensive prior experience in visual/manual interpretation of fluvial morphosedimentary forms. The other datasets (aerial and orthomosaic imagery) are passively sensed (photographic) and readily interpretable in the visible color spectrum (Table 1), requiring no special training however, that does not necessarily mean the labeling task is less difficult.

# 3.5. Coastal Lagoon and Barrier Beach Dynamics in False-Color Satellite Imagery

Data set E (Figure 4e and Table 1) consists of a time-series of 40, three-band false-color 10-m (122 × 342 × 3 pixels) cloud-free Sentinel-2 satellite images of coastal lagoon environments in Salinas Rivermouth Natural Preserve and National Wildlife Refuge in Monterey, California, collected between 31 December 2018 and 19 May 2021. The false color images consist of near infrared (band eight), red (band four), and green (band three). This three-band combination is commonly used for visual landscape classification where vegetation is present (Vuolo et al., 2016) because plant-covered land appears deep red, and denser plant growth is darker red. Water appears blue/black. The spatio-temporal time-series depicts various changes on the landscape, including the dynamics of the Salinas River mouth into the coastal ocean, surfzone and riverplume characteristics, changes to marsh and dune vegetation, and agricultural crop rotation. Therefore we defined the following classes: (a) water, (b) whitewater, (c) bare sand, (d) marsh veg, (e) dune veg, (f) crop/woody, (g) soil. The reader is referred to Figure S9 in Supporting Information S1 for more example imagery.

This imagery is further used to study the dynamics of beach breaching by a coastal river, and to compare the variability in geomorphic interpretation resulting from automated analysis of labels from three labelers labeling the same relatively complex imagery. Such imagery could be useful for opportunistic monitoring of coastal change from, among many potential uses, shoreline detection and characterization to assess trends in erosion and deposition, to assessments of habitat loss, flooding, surf zone hydrodynamics, agricultural development, bluff and sand dune dynamics. The frequency of important change at the coast is often greater than the frequency of available aerial platforms to provide imagery, especially in remote locations at short notice, and this makes the vertical and time-varying components of these landscapes especially difficult to unravel from opportunistic surveying/sampling. Satellite imagery with its regular timestamp therefore has a crucial role to play in linking time and spatial scales at coasts (McCarthy et al., 2017), and will play an increasingly important role in facilitating coastal science as imagery becomes higher resolution and better quality, and new sensors provide capabilities to sense new quantities (Vos et al., 2020).

#### 3.6. Coastal Evolution in Satellite Imagery

Data set F (Figure 4f and Table 1) consists of a time-series of 43, three-band visible-band pan-sharpened 15-m Landsat-8 satellite images ( $768 \times 768 \times 3$  pixels) of Cape Hatteras, Cape Hatteras National Seashore, North

BUSCOMBE ET AL. 13 of 31



Carolina, collected between 15 February 2015 and 27 September 2021. Data set F differs from data set E in three important respects; (a) imagery represent a larger area of over 10 km in each horizontal dimension; (b) imagery is visible-band; and (c) the dynamics captured, consisting of changing sandbars, sandwaves, beaches and wave breaking patterns, manifest over a larger timescale (79 months compared to 18 months of data set E). We labeled the following classes: (a) water, (b) whitewater (surf), (c) sand, (d) land (all dry land that is not sand). There are also some small clouds and shadows of clouds in the scene, all occurring above water, therefore they are labeled "water". However, separate classes for clouds and shadows might also be a valid strategy. The reader is referred to Figure S10 in Supporting Information S1 for more example imagery. This larger-scale (multi-km) imagery is used to demonstrate the utility in segmenting natural features at relatively large scales, and is also used to compare hand-digitization workflows with the methodology presented here.

# 4. Case Study Results

# 4.1. Image Size and Resolution

A comparison of label images at the two different grid sizes helps us understand at what grid size, and perhaps more importantly image size, we should ideally use for a given scene. A region of the 5-cm and 25-cm pixel imagery in data set A (Sherwood et al., 2021) are divided into  $1,024 \times 1,024 \times 3$  pixel tiles for annotation, which resulted in 99 and six tiles for the respective resolutions. It is more difficult to accurately label the larger, coarser resolution imagery for two reasons: the 25-cm imagery covers a much greater spatial extent than the 5-cm imagery, so features are smaller, and the imagery is less well resolved, therefore features are less distinct. However, images can be over-resolved for the task, and the time it takes to label a set scales approximately proportionally, at best, with the number of images in the set.

Each of the image tiles are labeled, then merged back into large label orthomosaics on the same spatial grids as the original orthomosaic images (Figure 5). In this case, errors are more readily observed when image tiles are merged, and assessed visually. We found this for both the 25-cm imagery and the 5-cm imagery; in Figure 5, those regions appear as abrupt changes in label values and are indicated by white boxes in Figures 5e through 5h. This artifact is more common for the coarser-resolution 25-cm imagery. The purpose of tiling of large imagery is to make the labeling tasks more manageable, and it also typically makes labeling faster. The disadvantage is that many of the errors in the higher resolution imagery occur or become apparent at tile boundaries. These errors are generally either caused by (a) a relatively low spatial density of annotations compared to the higher-resolution imagery, or (b) by annotations omitted by the labeler due to the larger size of the imagery. The majority of such errors occur at label boundaries and could be ameliorated through use of a spatial low-pass filter.

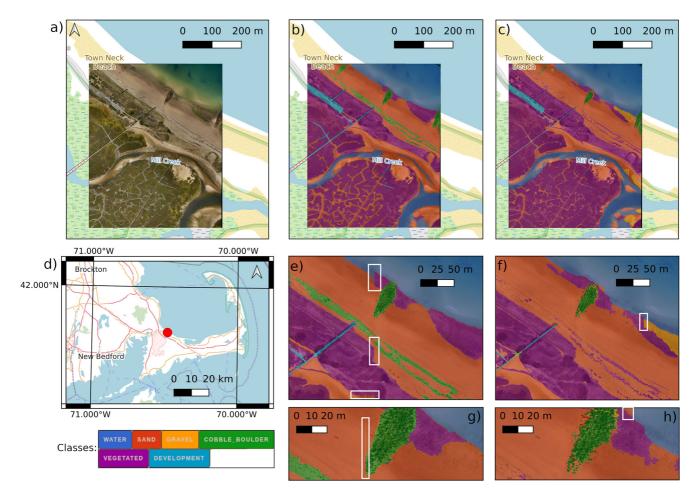
Other errors are due to misidentifications due to the lower resolution of the imagery; note how in Figure 5e the wrack line is labeled green (cobble/boulder), whereas in Figure 5f it is labeled "vegetated." The latter is perhaps more correct, because it is composed of dead vegetation. The task became ambiguous, because wrack is rough like cobbles but composed of organic matter. In addition, the wrack is much better resolved and identifiable in the 5-cm imagery. For this class set, we would use moderately low resolution imagery for this segmentation task, but small image tiles. However, the decision is dependent on the processes of interest. In this example, spatially less extensive, higher-resolution image tiles would be useful for delineating subtle differences in sedimentary grade or texture that only manifests at that scale, such as the difference between fine and coarse sand. Coarser resolution imagery may be sufficient for delineating the more obvious sedimentary transitions, such as gravel to boulders. Before embarking on segmentation tasks where image grid size can be varied it is recommended to use an exercise similar to this to determine a grid resolution and image size that is a good compromise for available time, required spatial density of annotations, and ideal image size where the smallest important features are visible (e.g., higher resolution may be needed for identifying animals or distinguishing between subtle sediment or vegetation types).

#### 4.2. Inter-Labeler Differences

Data set B is used to compare the products resulting from two labelers labeling the same complex data set. The mean agreement is high (Figure 6), as evidenced by a median of mean Dice scores of 0.76, and Dice scores are generally only marginally higher than equivalent IOU scores, suggesting class imbalance is not too much of a

BUSCOMBE ET AL. 14 of 31





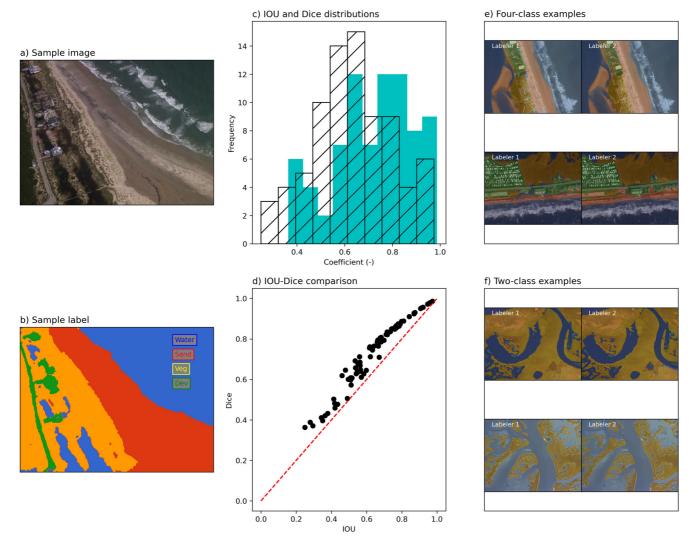
**Figure 5.** (a) A region of orthomosaic of Sandwich Town Beach (data set A); (b) 25-cm label imagery as a semi-transparent color overlay; (c) 5-cm label imagery as a semi-transparent color overlay; (d) geographic location of the site; (e) closer detail of (b); (f) closer detail of (c); (g) yet closer detail of (e); and (h) yet closer detail of (f). In (b), (c), (e) and (f), label imagery consists of small  $1024 \times 1024$  pixel label tiles that have been combined into a raster of full extent in a GIS. Classes are also depicts as colorful buttons (the same buttons used in the program Doodler when used to make the label tiles). White boxes highlight regions discussed in the text.

factor for this data set. There are many more examples of where Dice  $\gg$  IOU (i.e., IOU-Dice residual in Figure 6c is greater than, for example, 0.075), than where Dice and IOU are close.

# 4.3. Class Selection

An analysis of the labels generated from data set C presents an opportunity to discuss labeler agreement when a classification task is somewhat subjective, and how to achieve consensus by identifying which classes to lump together, and which to keep separate. IOU and Dice scores are surprisingly good (Dices scores range from 0.87 to 0.93) when evaluated over the full set of 4 classes (Figure 7a) and show greatest improvement (Dice scores range from 0.94 to 0.97) when the whitewater class is included with the deep water class and shallow is lumped with the dry land class, to create a binary or two-class set (Figure 7b). Any remaining low scores are partially the result of confusion over whether to include swash foam as whitewater. All Dice and IOU scores increased when evaluated over two classes instead of four, although not uniformly (Figure 7), suggesting class imbalance is variable. Analysis of a set of labels in this way from multiple labelers could also be used to identify any outlier labelers whose interpretations are different from the rest of the group. As in evident in Figure 7, there are no individuals among the five labelers who have a noticeably lower agreement.

BUSCOMBE ET AL. 15 of 31



**Figure 6.** (a) Sample image from the data set; (b) Label image associated with (b); (c) Histograms of Intersection over Union (IOU) and Dice scores for the 80 pairs of labeled aerial images; (d) IOU-Dice comparison; (e) Examples where mean Dice >0.075 than mean IOU; (f) Examples where mean Dice and mean IOU are within 0.075.

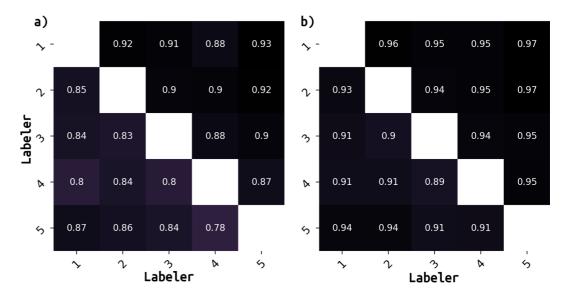
# 4.4. Specialized Labeling

Data set D used to compare the products resulting from two labelers labeling the same complex imagery requiring specialist interpretation. In this case, the mean agreement is lower than for the NOAA aerial imagery, as evidenced by a comparatively low median of mean Dice score of 0.43 compared to 0.76 for the NOAA aerial imagery (compare Figures 8a and 6a). This is possibly due to the task being more difficult, meaning large areas can be legitimately called two different classes (examples are shown in Figures 8d and 8e), and because there are more classes (eight instead of four), meaning the class-averaged IOU or Dice is affected by outlier classes.

Another major reason for the generally lower scores is that having more classes presents greater opportunity for a mismatch in the number of respective classes in each of a pair of label images. Recall that where different numbers of unique classes exist, that is, two different candidates for k, we choose k as the maximum length of the two respective class sets. The sidescan label set has, among those used in the present study, a greater percentage of images like this where there are unequal numbers of labels per image, therefore a greater percentage of conservative scores, which further decreases the class-averaged score.

Set-averaged Dice and IOU scores (i.e., the scalar mean of a distribution of mean scores) are close (Figure 8a), suggesting any class imbalance is not affecting the comparison between labels. Class imbalance may not be

BUSCOMBE ET AL. 16 of 31



**Figure 7.** Matrices quantifying agreement among five labelers numbered one through five. The upper-right half of each matrix shows Dice scores, and the lower-left have shows Intersection over Union (IOU) scores. Two labeling experiments are shown: left (a) used four classes (deep, white, shallow, and dry); right (b) used two classes, combining "deep" and "whitewater" as one class, and "shallow" and "dry' as the other.

avoidable if specific classes must be used for the scientific purpose the labeled imagery serves, however the effects of class imbalance can be reduced by merging appropriate classes, that is, a minority class into a majority class, where possible. If a class is infrequent, but deemed too important to miss, imagery could be cropped so the class imbalance issue is ameliorated, or the algorithms could be modified to use class weights.

The two examples shown in Figure 8e with relatively poor agreement do so for different reasons; in the upper example the two labelers have disagreed over the two shadow classes, and in the lower example the two labelers have disagreed where one identifies a region as coarse whereas the other identifies it as wood. In these examples, consensus could be achieved through some rules-based process, or by redoing the labels with lower-than-average IOU and/or Dice scores in order to achieve greater label precision through consensus (Goldstein et al., 2021; Monarch, 2021).

## 4.5. Multi-Labeler Comparison of Quantifying a Geomorphic Process

Data set E is used to compare the products resulting from three labelers labeling the same complex imagery of a geomorphic process. The overall agreement between Labelers 2 and 3 is very high, as evidenced by a mean Dice of 0.9 (Figure 9a). Additionally, the distribution of scores between Labeler 1 and Labelers 2 and 3 are almost identical.

In this case, mean Dice scores always exceed mean IOU scores (Figures 9b and 9c), suggesting class imbalance does affect the comparison between labels (water is by far the dominant class in every image). The two largest discrepancies between mean Dice and IOU scores are shown in Figure 9d; in each case, the white arrow highlights the major error, which in both cases is the mislabeling of water, which, as the dominant class, has a disproportionately negative affect on mean IOU compared to mean Dice. A comparison between IOU and Dice can also be used to detect outliers. The highlighted outlier in Figure 9e corresponds to the pair of labels shown in Figure 9f, in which the one from labeler 3 is missing one category, whitewater, which the program has called sand and which would have to be relabeled.

As for the geomorphic event we wished to describe using the segmentation data, namely the barrier breaching and "resealing" event that happened between 25 January 2019 and 10 April 2019, captured by seven cloud-free images, Figure 10 depicts the breach vicinity in each of the seven images, with the contoured outline of the sand category of the image segmentation created by each of the three labelers overlain. In all but one case, shown by the white rectangle in Figure 10g, all three labelers captured the outline of the barrier correctly, in the vicinity of

BUSCOMBE ET AL. 17 of 31



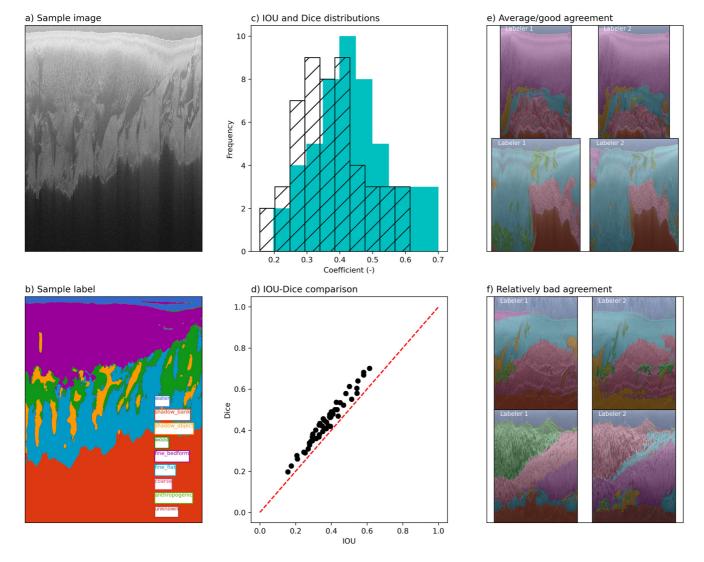


Figure 8. (a) Sample image from the data set; (b) Label image associated with (b); (c) Histograms of mean (class-averaged) Intersection over Union (IOU) and Dice scores for the 51 pairs of labeled sidescan images; (d) sample mean IOU –mean Dice comparison; (e) two examples of average/good agreement; and (f) two examples of relatively bad agreement.

the breach, plus the back barrier and shoreline areas. There are two additional images showing more temporary breaching events (on 24 April 2020 and 28 February 2021) in which all three labelers captured the outline of the barrier correctly (not shown). The average horizontal variability between outlines for the three respective labelers is within two pixels (20-m horizontal ground distance).

Aside from specific cases like those described above, a potential more generic downside of using highly discriminative models optimized for specific tasks is that they do not necessarily transfer well to out-of-distribution data. This is why Doodler works well to generate training data for other types of models that carry out segmentation on datasets at scale (i.e., with much more variety than a single image). To demonstrate how the MLP model framework does not transfer well to unseen data, and hence why for fully automated segmentation of unseen sample imagery requires a more powerful approach such as a deep neural network trained on thousands of examples, we use data set E once again. For each of the 40 images, we used the MLP model built on the small annotated scene to apply to a scene with an extent twice as large, extending down coast. The MLP model trained on each half image is able to extrapolate the broad categories that are significant at the boundary of the extent of annotations well, that is, at the bottom edge of the top half of the image (Figure 11) such as water, dune, and crops. However, it tends to under-predict the less dominant classes whitewater (surf), soil and sand, and predictions get worse the farther away from the boundary. The CRF model cannot fix all the errors in these under-predicted classes

BUSCOMBE ET AL. 18 of 31



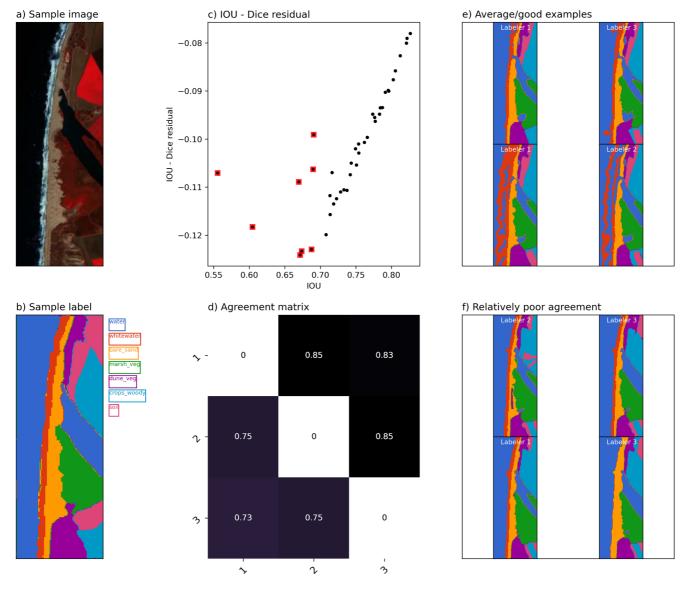


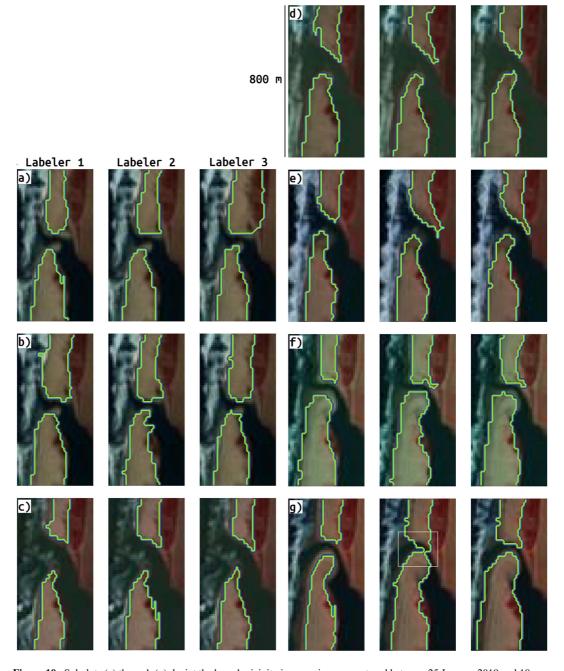
Figure 9. (a) Sample image from the data set (e); (b) Label image associated with (b); (c) mean Intersection over Union (IOU) vs. mean IOU—mean Dice residual for the 80 pairs of labeled multispectral satellite images, highlighting outlier labels; (d) IOU (bottom left matrix elements) and Dice (top right matrix elements) scores among all 3 labelers; (e) two examples of average/good agreement; and (f) two examples of relatively bad agreement.

however, the Doodler program itself results in annotations that could be used within alternative ML frameworks and it is likely that annotations with sufficient density for a good MLP solution would easily be sufficient for a more sophisticated model (perhaps at greater computational expense) because MLPs are relatively simple ML architectures. The fact that the annotations have been optimized through guided iteration towards a solution for a particular ML algorithm, does not mean they cannot be repurposed for, after all, they are simply example pixels of each class. And, as we mentioned above, Doodler is designed for both one-time data set segmentation and for generation of label imagery for training ML models such as deep learning models for fully automated image feature-extraction and class segmentation at scale, for application to Earth surface imagery.

# 4.6. Comparison With Manual Digitization

A single scene collected in 15 February 2015 (the first image in the collection) was annotated in a traditional way using hand digitization of polygons, then again using Doodler. This was conducted by the same individual on the

BUSCOMBE ET AL. 19 of 31



**Figure 10.** Subplots (a) through (g) depict the breach vicinity in seven images captured between 25 January 2019 and 10 April 2019, with the contoured outline of the sand category of the image segmentation created by each of the three labelers overlain. The white rectangle in (g) shows the only case where the sand polygon would suggest the barrier is still sealed, albeit by a single connecting pixel. Otherwise, the agreement is very close, within two pixels typically with a maximum discrepancy of four.

same day. It took 7.5 min to carefully label the scene and compute the segmentation using Doodler. We used an open-source annotation software (Skalski, 2019) to efficiently hand-digitize polygons for the entire scene. This program has similar zoom and pan tools to Doodler, which enables careful labeling of small features such as the relatively narrow sand beach and the surf zone (multiple lines of breaking waves). Additional imagery showing the stages of digitization is provided as Figure S11 in Supporting Information S1. The manual digitization took 25 min, or more than three times as long. Whereas we could have conducted this comparison using any of the datasets presented in Table 1, we chose this data set because the imagery is sufficiently large, and some classes

BUSCOMBE ET AL. 20 of 31

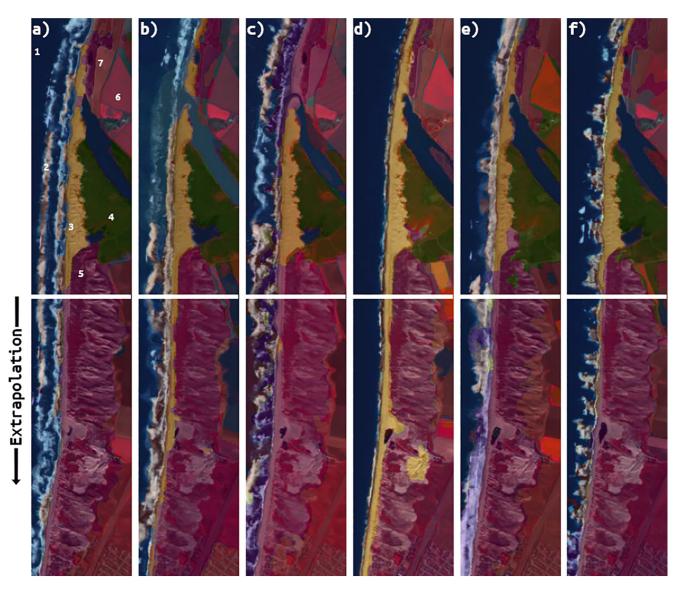


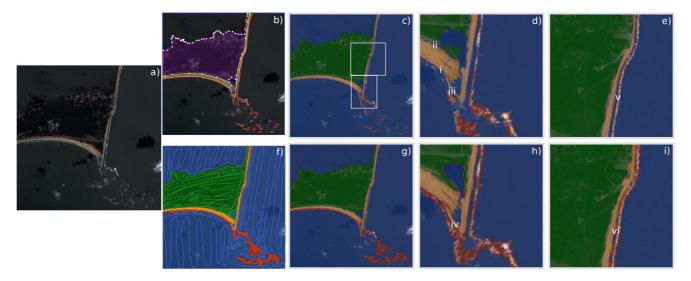
Figure 11. Output label images from a Multi-Layer Perceptron model built on the small annotated scene above the white horizontal white line in the center of the scene, then applied to the entire scene with an extent twice as large. In the extrapolated region, water, dunes, and crops are reasonably well predicted, but sand, whitewater (surf), and soil are not as well predicted.

sufficiently spatially limited, to warrant zooming and panning in order to accurately label. We note that the degree of zoom and pan is somewhat comparable between the two annotation programs, however the extent of annotation is much less with Doodler, and each annotation is much quicker to complete.

The digitized polygons were converted into a label image for direct comparison with the label image obtained using Doodler. A comparison of the inputs and results is presented in Figure 12. The mean IOU and Dice scores that quantify the agreement between the two label images are 0.48 and 0.5, respectively. This is low because the mean agreement for the two minority classes "surf" and "sand: are only approximately 0.015, whereas the agreement over "water" and "land" are approximately 0.97 each. Owing to the large class imbalance in this scene, quantitative comparison is limited. Qualitatively, we observe that the two label images differ in three important ways. First, there are a few small gaps in the label image where the labeler did not ensure matchup (or overlap) between adjacent polygons. This is a common limitation of hand-digitization, and here manifests most significantly as gaps between sand polygons, as indicated in Figure 12d by numeral i, and between the marsh and the beach, as indicated by numeral ii. Second, extremely small/thin objects are more difficult to hand digitize,

BUSCOMBE ET AL. 21 of 31





**Figure 12.** A comparison of hand-digitization versus human-in-the-loop segmentation workflows. The image (a) is the first in data set F, captured by Landsat 8 on 15 February 2015. The hand-drawn polygons (b) are rasterized to create a label image (c). Subplots (d) and (e) show details from the two regions identified in (c). The same image is segmented using sparse annotations (f), resulting in label image (g). The same regions highlighted in (d) and (e) are shown in (h) and (i), respectively, for the image segmented using Doodler. Numerals i through vi are discussed in the text.

resulting in the omission of the very thin sand bar, indicated by numeral iii in Figure 12d. The presence of this bar is marginally visible but also indicated by the adjacent breaking waves. Doodler was able to capture this feature properly (Figure 12h, numeral iv) with a few annotated pixels in this region. Third, in complex regions of transition where adjacent classes are indistinct at the level of zoom at which the labeler has chosen to label, such as near shore where waves are breaking on the sand beach, hand annotation generally results in overly coarse digitization compared to Doodler. Doodler is able to predict at the pixel level, whereas it is overly time consuming for hand digitization of polygons at the same scale. However, there are also advantages to relatively coarse hand digitation if it preserves actual boundaries better than a model prediction instance. An example is indicated by numerals v and vi in Figures 12e and 12i, respectively; hand digitization has labeled the ocean side of the beach better than Doodler, however Doodler has better labeled the pixel-level detail in the lagoon side of the beach.

# 5. Discussion

# 5.1. Obtaining High Levels of Agreement

The results suggest that given knowledgeable labelers, the Doodler program produces consistent label images (segmentations), even for complex scenes with numerous classes, indicating that multiple labelers can be used to label a data set and the results will be consistent and cohesive. The majority of errors in the labels are not necessarily due to the model but are consistent among labelers. The data sets shown here (Table 1) are a few among numerous datasets we have already successfully used the program with, from millimeter-scale grid sizes in close-range photography to multi-decimeter-scale pixels in satellite imagery, using between two and many tens of classes. We also tried several previous software implementations for the basic idea, and have arrived at a user interface by testing hundreds to thousands of individual samples by dozens of individual labelers. By combining unary potentials from a discriminative MLP model that encodes the conditional likelihood of a class given an image feature, with pairwise potentials that encode the joint likelihood of image features and classes together, the CRF technique exploits the benefits of both discriminative and generative ML model frameworks, and almost always results in an as or more accurate image segmentation than using the discriminative MLP model alone as determined visually on thousands of label samples; the program can generate a side-by-side comparison of the MLP output and CRF output for any sample image.

An advantage of using a so-called "cascade" of ML models whereby the outputs of the first is the inputs to the next (Figure 2), is that the second model can and often does revise the predictions of the first if they are

BUSCOMBE ET AL. 22 of 31



inconsistent with the second. This situation can often arise because the confidence of discriminative ML models, such as MLPs, are as much a reflection of the model feature-extraction and classification processes (summarized by learned parameters,  $\theta$ ) as the input data. That is why we say the model output is  $P(y|\theta, \mathbf{x})$  rather than simply  $P(y|\mathbf{x})$ , to acknowledge the joint importance of model parameters  $\theta$  with the specific image features  $\mathbf{x}$  used during training.

Outputs are further improved by having a human in the loop, that is, to immediately visually inspect segmentations for quality, and to add/remove annotations where necessary in places the model has mispredicted, and/or to adjust model hyperparameters (on an image-by-image basis if necessary). The percentage of imagery where such correction is necessary varies considerably by task (and to a certain degree the diligence of the individual labeler); on datasets tested to date, we estimate that approximately half or more of images require the addition of annotations beyond the initial sparse set, and approximately a tenth or less require the removal of annotations or the adjustment of hyperparameters. It is generally considered a good thing that the CRF solution is not overly sensitive to hyperparameter values, and that happens for several reasons by design (see Section 2.4), because that allows the instructions given to labelers to focus on how to annotate well.

Based on comparitive exercises between hand-digitization using polygons and our alternative workflow, we conclude that our methodology encoded into the Doodler program is always faster; approximately 3 times faster for the imagery used in Figure 12, and up to 10 times faster for other imagery we tested that does not require as much (or any) zooming and panning. Faster labeling makes multiple labeler data sets easier to obtain, and multi-labeler contexts have been shown to provide reliable label uncertainty metrics.

We also conclude that Doodler generally results in a segmentation that is as-or-more accurate than slower hand digitization workflows. First, Doodler ensures every pixel is labeled, whereas ensuring no gaps in the label raster that is the result of a hand-digitization workflow is difficult and often not managed. Additionally, Doodler picks up on pixel-level features that are too time-consuming to label or invisible at a reasonable zoom level, especially in complicated regions of transition. As a result, labels are finer-scale and more accurate at the pixel level because errors at boundaries between classes that arise due to hand digitization, which can be significant because of mixed pixels or due to coarse digitization, are significantly reduced. Modeling the likelihood of uncertain regions is crucially important for class assignment in particularly difficult regions of imagery in a deterministic manner.

# 5.2. Measuring Agreement

In general, it may be qualitatively observed that any IOU score above 0.5 is a very high level of agreement at the whole-image level, especially for high-resolution imagery. One of the really useful aspects of both IOU and Dice as metrics is that they both penalize pixel-level noise, and scores are therefore an accurate reflection of high-frequency label noise, which tends to increase with higher resolution imagery. A comparison of aggregated IOU scores between pairs of labels in whole datasets also meaningfully reflects the difficulty of the task; sidescan scores are typically lower than aerial and satellite imagery due to relative difficulty in interpretation.

However, due to averaging over classes and uneven numbers of classes among samples and datasets, both IOU and Dice scores are best treated as comparatives within datasets. In fact, when evaluating agreement (uncertainty) on individual datasets, computing and comparing both Dice and IOU scores can be useful for various reasons. We have shown it is possible to use them to discuss ways to detect class imbalance, outlier labelers, and label images in multi-labeler contexts, as well as reporting mean agreement for multi-labeled data sets as an uncertainty and quality metric, among other potential uses. IOU is always the more conservative metric than Dice, and that can sometimes be useful when deciding on the subsequent uses of the data. While it is very sensitive to class imbalance, there are potentially a lot of advantages to measuring total error rate, the sum all different pixels (i.e., all false positives and false negatives) divided by the number of pixels in the image. The per-class IOU and/or Dice scores can show problematic classes where there is lack of agreement (Figure 13). For example, in the sidescan data set (data set D), the distribution of per-class scores has the largest range; shadow and wood classes achieve relatively little consensus (Figure 13b). The two shadow classes would likely have to be merged for consistency, and better agreement over wood and all the other categories might be possible if a manual documenting examples is prepared (Goldstein et al., 2021). In the post-hurricane data set (data set B), sand is often difficult to distinguish from water for the same reasons as described for data set.

BUSCOMBE ET AL. 23 of 31



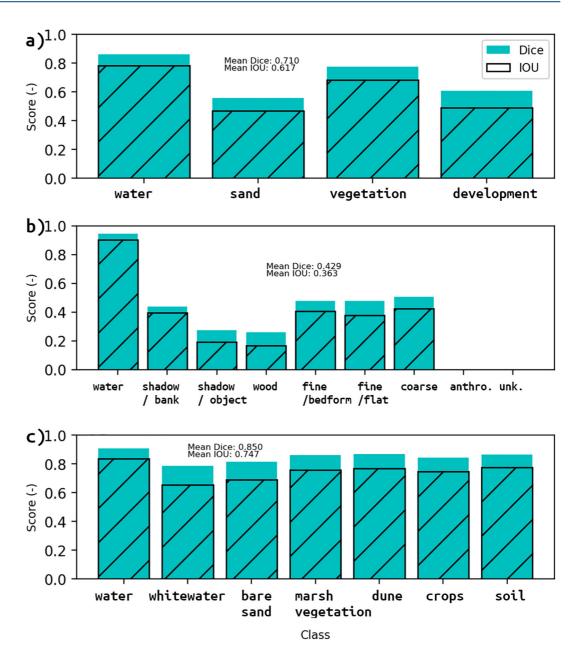


Figure 13. Per-class Dice and Intersection over Union (IOU; hatching) scores for (a) post-hurricane aerial imagery, (b) sidescan imagery, and (c) satellite imagery.

# **5.3.** The Value of Sparse Annotations

The sparse annotations provided by the human labeler are more valuable than the specific realisation of the fully labeled image. There are several reasons for this assertion. First, we tested alternative discriminative algorithms to the MLP that evaluate  $P(y|\theta, x)$  on features x that have already been extracted in a prescribed way. Among the alternative algorithms tested included the Random Forest and Support Vector Machine, both of which are used extensively in Earth surface processes research (Perry & Dickson, 2018; Provost et al., 2017; Yao et al., 2008) and worked well here too (see Figure S2 in Supporting Information S1 for representative comparison between MLP and Random Forest outputs). We chose the MLP because it is as or more accurate, with fewer model parameters, generally less overfitting, and had faster computation times. The key insight here is that the sparse annotations could be used with similar effect using a range of ML algorithms. This means that the Doodler program provides

BUSCOMBE ET AL. 24 of 31



a means to acquire sparse labels that are optimal for a many ML frameworks to carry out segmentation, not just the specific ML framework (MLP and CRF) that we have presented.

Second, as labels, annotations are more valuable than the pixelwise label imagery because there may be better ML model frameworks to predict pixelwise class from the sparse annotations in the near future, but it may be much longer before computers are able to label complex Earth surface imagery unaided with human-level accuracy. In fact there may already be viable ways to use the sparse annotations directly to train deep learning models for image segmentation, for example, by exploiting the variable spatial autocorrelation of each class (Hua et al., 2021) or by classifying image features as nearest neighbors in embedding space (Ke et al., 2021), however these techniques are currently much more computationally demanding, and would need large sparsely labeled datasets to achieve training convergence.

Third, the sparse annotations themselves encode the pixels chosen to represent the class in that region of the image, thus they are likely much better than a random selection of pixels from each scene and class at representing that class, perhaps efficiently encoding the line of greatest spatial transition (i.e., class boundaries). The CRF may on occasion (and by design) override the human label, and this may be quantified by locating (and/or counting) the pixels that differ in class between human input and CRF output. An analysis may reveal the degree to which and conditions under which the CRF over-rides the decisions of the human labeler. The Doodler program provides tools for extracting not only the sparse annotations and final projects, but also interim products, for any type of post facto analyses and evaluations.

Finally, the annotations themselves may be a proxy for other interesting properties of the data. For example, the spatial density of annotations may reveal areas of the scene that are more important for classification than others, or less ambiguous, or where the difficult transition areas are that the model is expected to predict. It is an interesting and as-yet under-explored supposition that there is some minimum sparsity of annotation necessary for a given target accuracy, but that would be complicated by the fact that multiple sets of annotations might give rise to identical outputs.

Other potentially informative derived attributes that relate to spatial autocorrelation and other spatial properties of the labeled regions include the spatial extent of each prediction, the shape the outline of that contiguously labeled region makes, and the spatial density with which annotations need to be made to properly segment the image. We find that the percentage of scene that is labeled for a satisfactory outcome varied with image size. It is between 10% and 20% for the sidescan imagery  $(1,024 \times 1,300-2,000 \text{ pixels})$  and between 10% and 30% for the NOAA imagery  $(1,000 \times 750 \times 3 \text{ pixels})$ .

In both cases, there is no systematic tendency for one labeler to spend more time on labeling overall, although there can be significant differences over individual images. However for the  $122 \times 342 \times 3$  pixel satellite imagery, the percentage is between 40% and 65%. The percentages may be an overestimate of the labels actually necessary for a good image segmentation, because the default "pen" (cursor) width is three pixels. That value is rarely changed by the majority of labelers in this study, although individual labelers tend to adopt that practice more readily than others, typically varying between 2 and 5 pixels depending on the scene. That is to say, it is possible that 1- or 2-pixel width annotations would have resulted in an equally good segmentation. That could be tested by using a morphological erosion operator on the sparse annotations then using the eroded doodles as inputs to the MLP and CRF estimation pipeline, and finally comparing outputs from full and thinned pen strokes. In some imagery used here, some labelers used thicker pens for the dominant classes, but others realized may have not done so because of the extra time it takes to change pen width. The number or spatial density of doodles, rather than thickness of pen, is generally a better local indication of scene complexity.

We found no significant correlation between either IOU or Dice score and percentage of the image annotated, either for individual images or for scores averaged over sets of labeled images. However, that is likely due to the fact that all labelers here are attentive and generally labeled a large percentage of the scene (between 10% and 65% of the scene, depending on image size) and in all areas of the image. Additionally labelers likely did so until the segmentation created from their sparse labels is satisfactory, that is, it seemed to accurately represent the underlying scene. Annotations are somewhat different, and individual labelers were even sometimes identifiable by their unique style. However, in this study agreement among labels was not identifiably related to a labeler's individual labeling sytle.

BUSCOMBE ET AL. 25 of 31



The program outputs also provide the means to analyze the annotations (like quantify their spatial density) and compare them. It is generally a more effective and efficient strategy to add and remove annotations than use model hyperparameters to modify CRF model predictions, although of course both are sometimes necessary of the most difficult imagery. Other useful metrics to track include the total percentage of the image labeled, although in this study that is not correlated with any qualitative accuracy metric or quantitative agreement metric because all labelers were careful and attentive and not more detailed with one class than with another. However, total percentage labeled would reveal situations where a labeler consistently annotated too much or too little of the image, both of which can be a problem due to either model underfitting or overfitting the data.

#### 5.4. Future Work

The most difficult imagery for Doodler would arguably be regarded as the most difficult for any image labeling program, namely degraded or poor quality imagery, and especially imagery where features and objects are small and hard to resolve because of low spatial resolution. Additionally, Doodler is not particularly well suited to labeling especially thin and short objects consisting of only tens to hundreds of pixels. For example, in large-format aerial imagery that represent large areas of ground, such hard-to-label objects would include individual pieces of driftwood, short and narrow paths and roads, vehicles, small buildings like cabins, people and other animals, among other common things. The common solution is to (a) exhaustively label almost every occurrence of the small, thin classes, and (b) to use a lot of zoom and panning, or smaller images, in which the labeler can better resolve the class and position the pen more accurately and precisely. However, because the CRF has agency it can override the human labels, and unfortunately tends to do so disproportionately for the more infrequent classes, which is almost always the classes associated with the small, thin objects. However, there are often trade-offs between available time and target accuracy with any labeling task. Therefore, on occasions when it is not efficient to use smaller images or spend time zooming and panning, especially if the main classification target is spatially extensive and/or continuous, the recommendation we would make is to classify the scene without employing the small, thin class(es); polygonal labels of those classes could be added later, rasterized, and merged with the label images of the other classes.

In Section 5.3 we stated that annotations are more valuable than the pixelwise label imagery because there may already be viable ways to use the sparse annotations directly to train deep learning models for image segmentation. The recent semi-supervised method of Ke et al. (2021) is particularly representative of current trends in this scope, utilizing a concept known as contrastive learning (Wei & Ji, 2021) that learn the similarity between labeled and unlabeled data and base classifications on that similarity. The similarity is learned from the data, and the regions considered to be adjacent to each require some form of abstraction such as defining superpixels (contiguous segments of image based on location and color obtained by clustering algorithms) or perhaps another trainable model component. It is therefore a more complex solution. Whereas, Doodler uses labeled pixels to assign classes to unlabeled pixels within each image, emerging ML techniques like Ke et al. (2021) also use those labels to assign classes within and across images. Such advances are possible by utilizing learned embedding representations of class-image pairings over larger datasets. Tools like Doodler would still be necessary to both collect the sparse annotations, and to generate independent data to evaluate the outputs of an automated technique for collections of images.

Although that was not carried out here in order to measure agreement over class sets and imagery among several labelers based on verbal instructions alone, upon inspection of the results we now recommend discussing and practicing candidate class sets with a small sample of imagery, and then having small a group of labelers trial, no matter how trivial the task may seem beforehand (Geiger et al., 2021). Regardless of hypothesized degree of ambiguity in a given labeling task, individual labelers vary a little in terms of diligence and skill, and with a lot of Earth surface imagery there is an expectation for different labels in ambiguous regions of imagery, for the reasons discussed in Section 1. Therefore achieving consensus is (a) part design, by using a modeling framework that is designed to objectively arrive at consensus in labels across the scene based on class-feature pixel pairings and (b) part analysis, by analyzing agreement in segmentations of the same imagery by multiple labelers. Analysis of labels for the purposes of deciding on optimal class sets, and achieving consensus, is only possible when multiple labelers are used, although analysis of labels made by the same labeler on separate occasions might also have some value.

BUSCOMBE ET AL. 26 of 31



More sophisticated labeling workflows would include those that modeled the likelihood (confidence) of the sparse annotations themselves, or provided ways for the labelers themselves to provide that assessment (Monarch, 2021) and that may be the subject of future work. There is also much more work that needs to be done concurrently into strategies for selecting images to be labeled, such as active learning (Goldstein et al., 2020), automatically labeling data using embeddings (Ding et al., 2020) and other data representations that have been found by application learning and transfer-learning algorithms (Cunha et al., 2020), or discovered using synthetic data (Wu et al., 2019).

# 5.5. Human-in-the-Loop Image Segmentation

Scoping feasible applications of Deep Learning in the geosciences benefits from rapid prototyping of ideas, model frameworks, and trained models used in a transfer-learning workflows that are often inherited from other disciplines (Buscombe & Carini, 2019; Buscombe et al., 2020; Cunha et al., 2020; Goldstein et al., 2020; Yang et al., 2020). The challenge is to evaluate their utility using domain-specific labeled datasets, perhaps against baseline methods that may already exist in that domain. The availability of labeled data, and especially the availability in analysis ready formats that might be readily ingested into a model training workflow, is the major impediment to uptake of advanced data analytics such as Deep Learning among the community of Earth surface scientists. While semi- or un-supervised classification methods are gaining more attention in many research contexts (Le et al., 2019) and are a staple method in landcover classification of mostly relatively coarse-resolution imagery (Deng & Clausi, 2005; Smits & Dellepiane, 1997), human annotators will continue to be vital for the success of many tasks that can be automated using ML. Despite the fact that the development of unsupervised methods require labeled data for development and, especially, evaluation, supervised methods at the time of writing are still state-of-the-art, and considered necessary to model imagery with high intra-class variance, such as a lot of Earth surface imagery. Supervised ML will therefore continue to be popular, and powerful, if facilitated by open-source tools that make data labeling more efficient, and analyses of uncertainty that add vital context to its use. Doodler, as what Monarch (2021) refers to as a "smart interface for semantic segmentation," is one of many specific software tools or interfaces (Bueno et al., 2020; Goldstein et al., 2021; Zhao et al., 2020) for the generation of large labeled data sets (Kashinath, Mudigonda, et al., 2021; Sumbul et al., 2019) that can be used for teaching and self-exploration of Deep Learning techniques, for use in transfer learning, and for new model development. Doodler is an open-source program that runs in a web browser, and may be one of many similar future implementations that might use human-in-the-loop ML for efficient labeling of other scientifically relevant label data such as those generated from time-series signals or social media content (Cai et al., 2017). These methods also have great potential for segmenting digital elevation model imagery, which may have great utility in a great variety of geomorphic applications.

The use of an ML model cascade, whereby the outputs of one classifier (MLP) is checked for consistency by another independent classifier (CRF), is crucial to the success of the approach for a wider variety of imagery and class sets. Image standardization, image feature engineering, spatial filtering, and the use of an ML model cascade all help reduce sensitivity of model outputs to user hyperparameters. These allow the human labeler to concentrate on annotating well, rather than spend time adjusting hyperparameters. We show that the proportion of the image pixels that require annotation for accurate pixelwise label image is relatively low around 10% of pixels for images of a size that is typically suitable for the program without excessive use of zoom and pan tools, which is imagery typically 3,000 pixels in either horizontal dimension or less. Discrepancies in agreement are unavoidable with multiple labelers and represent a source of irreducible uncertainty in all image segmentation workflows. Doodler provides the means to rapidly label images, therefore multi-labeler label data sets are more readily acquired and the irreducible error can be quantified. Further, we show how combining agreement metrics can be used to flag inconsistent label images and annotation styles, and identify the effects of class imbalance. Dice and IOU scores are shown to be useful metrics for reporting agreement between segmentations of the same data by more than one labeler, and we recommend reporting mean agreement for multi-labeled datasets as an uncertainty and quality metric, per image, per class, or aggregated over images and/or classes. We also show how the metrics can be used to detect class imbalance, outlier labelers, and label images in multi-labeler contexts. Even though segmentations vary from person to person, that does not introduce unreasonable variance in label images created by different people, at different times, or using different computational infrastructure.

BUSCOMBE ET AL. 27 of 31



# 6. Conclusions

We describe a human-in-the-loop machine learning system involving a graphical user interface for fast, interactive segmentation of N-dimensional (x, y, N) images into two-dimensional (x, y) label images. It is designed to meet two objectives: (a) segmentation of relatively small datasets for specific geoscientific inquiries, and (b) segmentation of small to large amounts of imagery for subsequent training of other types of ML models for fully automated segmentation of large datasets. The program is designed to work with any type of Earth surface imagery. We demonstrate the approach using five case study datasets from river, estuarine, and open coast environments of the United States; (a) segmentation of beach sediments in visible-band aerial orthomosaic imagery to document change to beaches of Cape Cod, Massachusetts; (b) segmentation of post-hurricane aerial imagery from North and South Carolina, for assessment of storm impacts; (c) segmentation of aerial imagery for delineating complex shoreline environments; (d) segmentation of sidescan sonar imagery for mapping in-stream physical habitats in coastal plain rivers of Mississippi; (e) segmentation of false-color Sentinel-2 satellite imagery of coastal lagoon environments in Monterey, California, to study the dynamics of river breaching of beaches; and (f) segmentation of larger visible-band Landsat-8 satellite imagery of Cape Hatteras, North Carolina, to study coastal landform evolution at a regional scale. The datasets consist of irregular grids (each pixel does not represent the same spatial footprint), as well as regular grids. Based on comparative exercises between hand-digitization using polygons and our alternative workflow, we conclude that our methodology encoded into the Doodler program is always faster, and also generally results in a segmentation that is as-or-more accurate than slower hand digitization workflows. We thereby demonstrate the effectiveness of the approach using geophysical, photographic, and multispectral imagery, as well as regular and irregular grids, and several different class sets and pixel sizes. The technique is reproducible in the sense that all decisions made by human labeler and ML algorithms (and their specific sequence) can be encoded to file, therefore the entire process can be played back and new outputs generated with alternative decisions and/or algorithms. We therefore expect our human-in-the-loop labeling workflow to have widespread applicability in Earth and Space scientific applications.

# **Data Availability Statement**

363(6433). https://doi.org/10.1126/science.aau0323

All data used for this study are available on Dryad via Buscombe et al. (2022). The code for the labeling interface is available on Zenodo and GitHub via Buscombe (2022).

# Acknowledgments References

Thanks to Tanja Williamson, Meg Palm-

sten, Chris Magirl, and two anonymous

use of trade, firm, or product names is for

reviewers for helpful suggestions. Any

descriptive purposes only and does not

imply endorsement by the U.S. Govern-

by the U.S. Geological Survey Coastal/

Marine Hazards and Resources Program

Appropriations for Disaster Relief Act of

ment. This work has been supported

and by Congressional appropriations

through the Additional Supplemental

2019 (H.R. 2157), EBG acknowledges

support from USGS (G20AC00403).

USGS Community for Data Integra-

tion-funded Coast Train project.

C. Bodine acknowledges support from

USFWS (F19AC00836), J. Favela and S.

Fitzpatrick acknowledge support from the

Anders, N. S., Seijmonsbergen, A. C., & Bouten, W. (2011). Segmentation optimization and stratified object-based analysis for semi-automated geomorphological mapping. Remote Sensing of Environment, 115(12), 2976–2985. https://doi.org/10.1016/j.rse.2011.05.007

Barlow, J., Franklin, S., & Martin, Y. (2006). High spatial resolution satellite imagery, DEM derivatives, and image segmentation for the detection of mass wasting processes. *Photogrammetric Engineering & Remote Sensing*, 72(6), 687–692. https://doi.org/10.14358/pers.72.6.687

Barnard, P. L., Dugan, J. E., Page, H. M., Wood, N. J., Hart, J. A. F., & Cayan, D. R. (2021). Multiple climate change-driven tipping points for coastal systems. *Scientific Reports*, 11(1), 1–13. https://doi.org/10.1038/s41598-021-94942-7

coastal systems. Scientific Reports, 11(1), 1–13. https://doi.org/10.1038/s41598-021-94942-7
Bayr, U., & Puschmann, O. (2019). Automatic detection of woody vegetation in repeat landscape photographs using a convolutional neural

network. Ecological Informatics, 50, 220–233. https://doi.org/10.1016/j.ecoinf.2019.01.012

Bergen, K. J., Johnson, P. A., Maarten, V., & Beroza, G. C. (2019). Machine learning for data-driven discovery in solid Earth geoscience. Science,

Beuzen, T., Goldstein, E. B., & Splinter, K. D. (2019). Ensemble models from machine learning: An example of wave runup and coastal dune erosion. *Natural Hazards and Earth System Sciences*, 19(10), 2295–2309. https://doi.org/10.5194/nhess-19-2295-2019

Bishop, C. (2006). Pattern recognition and machine learning. Springer. Retrieved from https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning/

pattern-recognition-machine-learning/
Bueno, A., Zuccarello, L., Díaz-Moreno, A., Woollam, J., Titos, M., Benítez, C., et al. (2020). PICOSS: Python interface for the classification of

seismic signals. Computers & Geosciences, 142, 104531. https://doi.org/10.1016/j.cageo.2020.104531

Buscombe, D. (2017). Shallow water benthic imaging and substrate characterization using recreational-grade sidescan-sonar. Environmental

Modelling & Software, 89, 1–18. https://doi.org/10.1016/j.envsoft.2016.12.003

Buscombe, D. (2022). Dash-Doodler, https://doi.org/10.5281/zenodo.5847379

Buscombe, D., & Carini, R. J. (2019). A data-driven approach to classifying wave breaking in infrared imagery. *Remote Sensing*, 11(7), 859. https://doi.org/10.3390/rs11070859

Buscombe, D., Carini, R. J., Harrison, S. R., Chickadel, C. C., & Warrick, J. A. (2020). Optical wave gauging using deep neural networks. *Coastal Engineering*, 155, 103593. https://doi.org/10.1016/j.coastaleng.2019.103593

Buscombe, D., Goldstein, E. G., Sherwood, C. R., Bodine, C., Favela, J., Fitzpatrick, S., et al. (2022). Dataset accompanying Buscombe et al. Human-in-the-loop segmentation of Earth surface imagery. https://doi.org/10.5061/dryad.2fqz612ps

Buscombe, D., & Grams, P. E. (2018). Probabilistic substrate classification with multispectral acoustic backscatter: A comparison of discriminative and generative models. *Geosciences*, 8(11), 395. https://doi.org/10.3390/geosciences8110395

Buscombe, D., Grams, P. E., & Kaplinski, M. A. (2017). Compositional signatures in acoustic backscatter over vegetated and unvegetated mixed sand-gravel riverbeds. *Journal of Geophysical Research: Earth Surface*, 122(10), 1771–1793. https://doi.org/10.1002/2017jf004302

BUSCOMBE ET AL. 28 of 31



- Buscombe, D., Grams, P. E., & Smith, S. M. (2016). Automated riverbed sediment classification using low-cost sidescan sonar. *Journal of Hydraulic Engineering*, 142(2), 06015019. https://doi.org/10.1061/(asce)hy.1943-7900.0001079
- Buscombe, D., & Ritchie, A. C. (2018). Landscape classification with deep neural networks. *Geosciences*, 8(7), 244. https://doi.org/10.3390/geosciences8070244
- Cai, J., Huang, B., & Song, Y. (2017). Using multi-source geospatial big data to identify the structure of polycentric cities. Remote Sensing of Environment, 202, 210–221. https://doi.org/10.1016/j.rse.2017.06.039
- Carbonneau, P. E., Dugdale, S. J., Breckon, T. P., Dietrich, J. T., Fonstad, M. A., Miyamoto, H., & Woodget, A. S. (2020). Adopting deep learning methods for airborne RGB fluvial scene classification. *Remote Sensing of Environment*, 251, 112107. https://doi.org/10.1016/j.rse.2020.112107
- Carleer, A., Debeir, O., & Wolff, E. (2005). Assessment of very high spatial resolution satellite image segmentations. *Photogrammetric Engineering & Remote Sensing*, 71(11), 1285–1294. https://doi.org/10.14358/pers.71.11.1285
- Chaudhary, P., D'Aronco, S., Moy de Vitry, M., Leitão, J. P., & Wegner, J. D. (2019). Flood-water level estimation from social media images. ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences, 4(2/W5), 5–12. https://doi.org/10.5194/ isprs-annals-iv-2-w5-5-2019
- Chen, S. A., Escay, A., Haberland, C., Schneider, T., Staneva, V., & Choe, Y. (2018). Benchmark dataset for automatic damaged building detection from post-hurricane remotely sensed imagery. Retrieved from https://arxiv.org/abs/1812.05581
- Cheng, H.-D., Jiang, X. H., Sun, Y., & Wang, J. (2001). Color image segmentation: Advances and prospects. Pattern Recognition, 34(12), 2259–2281. https://doi.org/10.1016/s0031-3203(00)00149-7
- Chilson, C., Avery, K., McGovern, A., Bridge, E., Sheldon, D., & Kelly, J. (2019). Automated detection of bird roosts using NEXRAD radar data and Convolutional Neural Networks. Remote Sensing in Ecology and Conservation, 5(1), 20–32. https://doi.org/10.1002/rse2.92
- Ching, T., Himmelstein, D. S., Beaulieu-Jones, B. K., Kalinin, A. A., Do, B. T., Way, G. P., et al. (2018). Opportunities and obstacles for deep learning in Biology and medicine. *Journal of the Royal Society Interface*, 15(141), 20170387. https://doi.org/10.1098/rsif.2017.0387
- Chmiel, M., Walter, F., Wenner, M., Zhang, Z., McArdell, B. W., & Hibert, C. (2021). Machine learning improves debris flow warning. Geophysical Research Letters, 48(3), e2020GL090874. https://doi.org/10.1029/2020gl090874
- Costa, H., Foody, G. M., & Boyd, D. S. (2018). Supervised methods of image segmentation accuracy assessment in land cover mapping. Remote Sensing of Environment, 205, 338–351. https://doi.org/10.1016/j.rse.2017.11.024
- Sensing of Environment, 203, 338–351. https://doi.org/10.1016/j.rse.2017.11.024

  Crisci, C., Ghattas, B., & Perera, G. (2012). A review of supervised machine learning algorithms and their applications to ecological data. *Ecolog-*
- ical Modelling, 240, 113–122. https://doi.org/10.1016/j.ecolmodel.2012.03.001 Csurka, G., Larlus, D., Perronnin, F., & Meylan, F. (2004). What is a good evaluation measure for semantic segmentation. *IEEE PAMI*, 26(1).
- Cunha, A., Pochet, A., Lopes, H., & Gattass, M. (2020). Seismic fault detection in real data using transfer learning from a convolutional neural network pre-trained with synthetic seismic data. Computers & Geosciences, 135, 104344. https://doi.org/10.1016/j.cageo.2019.104344
- Deng, H., & Clausi, D. A. (2005). Unsupervised segmentation of synthetic aperture radar sea ice imagery using a novel Markov random field model. IEEE Transactions on Geoscience and Remote Sensing, 43(3), 528–538. https://doi.org/10.1109/tgrs.2004.839589
- Ding, L., Tang, H., & Bruzzone, L. (2020). LANET: Local attention embedding to improve the semantic segmentation of remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 59(1), 426–435.
- Drăguţ, L., & Eisank, C. (2012). Automated object-based classification of topography from STRM data. Geomorphology, 141, 21-33.
- Everingham, M., Van Gool, L., Williams, C. K., Winn, J., & Zisserman, A. (2010). The PASCAL visual object classes (VOC) challenge. International Journal of Computer Vision, 88(2), 303–338. https://doi.org/10.1007/s11263-009-0275-4
- Farr, T. G., Rosen, P. A., Caro, E., Crippen, R., Duren, R., Hensley, S., & Roth, L. (2007). The shuttle radar topography mission. Reviews of Geophysics, 45(2), RG2004. https://doi.org/10.1029/2005rg000183
- Fox, M., Bodin, T., & Shuster, D. L. (2015). Abrupt changes in the rate of Andean Plateau uplift from reversible jump Markov chain Monte Carlo inversion of river profiles. *Geomorphology*, 238, 1–14. https://doi.org/10.1016/j.geomorph.2015.02.022
- Gaddes, M., Hooper, A., & Bagnardi, M. (2019). Using machine learning to automatically detect volcanic unrest in a time series of interferograms. *Journal of Geophysical Research: Solid Earth*, 124(11), 12304–12322. https://doi.org/10.1029/2019jb017519
- Gardner, M. W., & Dorling, S. (1998). Artificial neural networks (the multilayer perceptron) a review of applications in the atmospheric sciences. *Atmospheric Environment*, 32(14–15), 2627–2636. https://doi.org/10.1016/s1352-2310(97)00447-0
- Geiger, R. S., Cope, D., Ip, J., Lotosh, M., Shah, A., Weng, J., & Tang, R. (2021). "Garbage in, garbage out" revisited: What do machine learning application papers report about human-labeled training data? *Quantitative Science Studies*, 1–32.
- Gil, Y., David, C. H., Demir, I., Essawy, B. T., Fulweiler, R. W., Goodall, J. L., et al. (2016). Toward the Geoscience paper of the future: Best practices for documenting and sharing research from data to software to provenance. *Earth and Space Science*, 3(10), 388–415. https://doi.org/10.1002/2015ea000136
- Goldstein, E. B., Buscombe, D., Lazarus, E., Mohanty, S. D., Rafique, S. R., Anarde, K. A., et al. (2021). Labeling post-storm coastal imagery for machine learning: Measurement of inter-rater agreement. Earth and Space Sciences, 8. https://doi.org/10.1029/2021EA001896
- Goldstein, E. B., & Coco, G. (2015). Machine learning components in deterministic models: Hybrid synergy in the age of data. Frontiers in Environmental Science, 3, 33. https://doi.org/10.3389/fenvs.2015.00033
- Goldstein, E. B., Coco, G., & Plant, N. G. (2019). A review of machine learning applications to coastal sediment transport and morphodynamics. Earth-Science Reviews, 194, 97–108. https://doi.org/10.1016/j.earscirev.2019.04.022
- Goldstein, E. B., Mohanty, S. D., Rafique, S. N., & Valentine, J. (2020). An active learning pipeline to detect hurricane washover in post-storm aerial images. *EarthArXiv*.
- Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., & Moore, R. (2017). Google Earth engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment*, 202, 18–27. https://doi.org/10.1016/j.rse.2017.06.031
- Grams, P. E., Buscombe, D., Topping, D. J., Kaplinski, M., & Hazel, J. E. (2019). How many measurements are required to construct an accurate sand budget in a large river? Insights from analyses of signal and noise. *Earth Surface Processes and Landforms*, 44(1), 160–178. https://doi.org/10.1002/esp.4489
- Gray, P. C., Fleishman, A. B., Klein, D. J., McKown, M. W., Bézy, V. S., Lohmann, K. J., & Johnston, D. W. (2019). A Convolutional Neural Network for detecting sea turtles in drone imagery. *Methods in Ecology and Evolution*, 10(3), 345–355. https://doi.org/10.1111/2041-210x.13132
- Hossain, M. D., & Chen, D. (2019). Segmentation for object-based image analysis (OBIA): A review of algorithms and challenges from remote sensing perspective. ISPRS Journal of Photogrammetry and Remote Sensing, 150, 115–134. https://doi.org/10.1016/j.isprsjprs.2019.02.009
- Hua, Y., Marcos, D., Mou, L., Zhu, X. X., & Tuia, D. (2021). Semantic segmentation of remote sensing images with sparse annotations. IEEE Geoscience and Remote Sensing Letters, 19.
- James, L. A., Hodgson, M. E., Ghoshal, S., & Latiolais, M. M. (2012). Geomorphic change detection using historic maps and DEM differencing: The temporal dimension of geospatial analysis. Geomorphology, 137(1), 181–198. https://doi.org/10.1016/j.geomorph.2010.10.039

BUSCOMBE ET AL. 29 of 31

s00267-017-0880-x



- Kashinath, K., Mudigonda, M., Kim, S., Kapp-Schwoerer, L., Graubner, A., Karaismailoglu, E., & Mahesh, A. (2021). ClimateNet: An expert-labeled open dataset and deep learning architecture for enabling high-precision analyses of extreme weather. *Geoscientific Model Development*, 14(1), 107–124.
- Kashinath, K., Mustafa, M., Albert, A., Wu, J., Jiang, C., Esmaeilzadeh, S., et al. (2021). Physics-informed Machine Learning: Case studies for weather and climate modelling. *Philosophical Transactions of the Royal Society A*, 379, 20200093. https://doi.org/10.1098/rsta.2020.0093
- Ke, T.-W., Hwang, J.-J., & Yu, S. X. (2021). Universal weakly supervised segmentation by pixel-to-segment contrastive learning. arXiv preprint arXiv:2105.00957.
- Kemker, R., Salvaggio, C., & Kanan, C. (2018). Algorithms for semantic segmentation of multispectral remote sensing imagery using deep learning. ISPRS Journal of Photogrammetry and Remote Sensing, 145, 60–77. https://doi.org/10.1016/j.isprsjprs.2018.04.014
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. Proceedings of the 3rd International Conference on learning representations (ICLR). arXiv:1412.6980.
- Koenderink, J. J., & Van Doorn, A. J. (1992). Surface shape and curvature scales. Image and Vision Computing, 10(8), 557–564. https://doi.org/10.1016/0262-8856(92)90076-f
- Koller, D., & Friedman, N. (2009). Probabilistic graphical models: Principles and techniques. MIT press.
- Kotaridis, I., & Lazaridou, M. (2021). Remote sensing image segmentation advances: A meta-analysis. *ISPRS Journal of Photogrammetry and Remote Sensing*, 173, 309–322. https://doi.org/10.1016/j.isprsjprs.2021.01.020
- Krähenbühl, P., & Koltun, V. (2011). Efficient inference in fully connected CRFs with Gaussian edge potentials. Advances in Neural Information Processing Systems, 24, 109–117.
- Kranenburg, C., Ritchie, A., Brown, J., Over, J., Buscombe, D., Sherwood, C., & Wernette, P. (2020). Post-hurricane Florence aerial imagery: Cape fear to Duck. U.S. Geological Survey data release. https://doi.org/10.5066/P91KB9SF
- Kumar, S., & Hebert, M. (2006). Discriminative random fields. International Journal of Computer Vision, 68(2), 179–201. https://doi.org/10.1007/s11263-006-7007-9
- Kurnaz, M. N., Dokur, Z., & Ölmez, T. (2005). Segmentation of remote-sensing images by incremental neural network. Pattern Recognition Letters, 26(8), 1096–1104. https://doi.org/10.1016/j.patrec.2004.10.004
- Lang, S., Hay, G. J., Baraldi, A., Tiede, D., & Blaschke, T. (2019). GEOBIA achievements and spatial opportunities in the era of Big Earth Observation Data. ISPRS International Journal of Geo-Information, 8(11), 474. https://doi.org/10.3390/ijgi8110474
- Larsen, A., Nardin, W., Van de Lageweg, W., & Bätz, N. (2021). Biogeomorphology, quo vadis? On processes, time, and space in biogeomorphology. Earth Surface Processes and Landforms, 46(1), 12–23. https://doi.org/10.1002/esp.5016
- Le, H. M., Goncalves, B., Samaras, D., & Lynch, H. (2019). Weakly labeling the antarctic: The penguin colony case. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops* (pp. 18–25).
- Lefsky, M. A. (2010). A global forest canopy height map from the moderate resolution imaging spectroradiometer and the geoscience laser altim-
- eter system. Geophysical Research Letters, 37(15), L15401. https://doi.org/10.1029/2010gl043622

  McCarthy, M. J., Colna, K. E., El-Mezayen, M. M., Laureano-Rosario, A. E., Méndez-Lázaro, P., Otis, D. B., et al. (2017). Satellite remote sensing for coastal management: A review of successful applications. Environmental Management, 60(2), 323–339. https://doi.org/10.1007/
- Mi, L., & Chen, Z. (2020). Superpixel-enhanced deep neural forest for remote sensing image semantic segmentation. ISPRS Journal of Photogrammetry and Remote Sensing, 159, 140–152. https://doi.org/10.1016/j.isprsjprs.2019.11.006
- Monarch, R. (2021). Human-in-the-Loop machine learning: Active learning and annotation for human-centered AI. Manning Publications.
- Ni, J., Wu, T., Zhu, X., Hu, G., Zou, D., Wu, X., & Pang, Q. (2021). Simulation of the present and future projection of permafrost on the Qinghai-Tibet Plateau with statistical and machine learning models. *Journal of Geophysical Research: Atmospheres*, 126(2), e2020JD033402. https://doi.org/10.1029/2020jd033402
- NOAA. (2021). National geodetic Survey emergency response imagery. Retrieved from https://storms.ngs.noaa.gov/
- Olhede, S. C., & Wolfe, P. J. (2018). The growing ubiquity of algorithms in society: Implications, impacts and innovations. *Philosophical Transactions of the Royal Society A: Mathematical, Physical & Engineering Sciences*, 376(2128), 20170364. https://doi.org/10.1098/rsta.2017.0364
- Over, J.-S. R., Ritchie, A. C., Kranenburg, C. J., Brown, J. A., Buscombe, D. D., Noble, T., & Wernette, P. A. (2021). Processing coastal imagery with agisoft metashape professional edition, version 1.6—structure from motion workflow documentation (Tech. Rep.). US Geological Survey.
- Pandey, P. C., Koutsias, N., Petropoulos, G. P., Srivastava, P. K., & Ben Dor, E. (2021). Land use/land cover in view of Earth observation: Data sources, input dimensions, and classifiers—A review of the state of the art. *Geocarto International*, 36(9), 957–988. https://doi.org/10.1080/10106049.2019.1629647
- Perry, G. L., & Dickson, M. E. (2018). Using machine learning to predict geomorphic disturbance: The effects of sample size, sample prevalence, and sampling strategy. *Journal of Geophysical Research: Earth Surface*, 123(11), 2954–2970. https://doi.org/10.1029/2018jf004640
- Plant, N. G., & Stockdon, H. F. (2012). Probabilistic prediction of barrier-island response to hurricanes. *Journal of Geophysical Research*, 117(F3), F03015. https://doi.org/10.1029/2011jf002326
- Provost, F., Hibert, C., & Malet, J.-P. (2017). Automatic classification of endogenous landslide seismicity using the random forest supervised classifier. *Geophysical Research Letters*, 44(1), 113–120. https://doi.org/10.1002/2016gl070709
- Quinn, J. A., Nyhan, M. M., Navarro, C., Coluccia, D., Bromley, L., & Luengo-Oroz, M. (2018). Humanitarian applications of Machine Learning with remote-sensing data: Review and case study in refugee settlement mapping. *Philosophical Transactions of the Royal Society A: Mathematical, Physical & Engineering Sciences*, 376(2128), 20170363. https://doi.org/10.1098/rsta.2017.0363
- Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., & Carvalhais, N. (2019). Deep learning and process understanding for data-driven Earth system science. *Nature*, 566(7743), 195–204. https://doi.org/10.1038/s41586-019-0912-1
- Richardson, A. D., Hufkens, K., Milliman, T., & Frolking, S. (2018). Intercomparison of phenological transition dates derived from the Pheno-Cam Dataset V1. 0 and MODIS satellite remote sensing. Scientific Reports, 8(1), 1–12. https://doi.org/10.1038/s41598-018-23804-6
- Ridge, J. T., Gray, P. C., Windle, A. E., & Johnston, D. W. (2019). Deep learning for coastal resource conservation: Automating detection of shellfish reefs. Remote Sensing in Ecology and Conservation, 6.
- Schwanghart, W., & Scherler, D. (2014). TopoToolbox 2–MATLAB-based software for topographic analysis and modeling in Earth surface sciences. Earth Surface Dynamics, 2(1), 1–7. https://doi.org/10.5194/esurf-2-1-2014
- Serre, T. (2019). Deep learning: The good, the bad, and the ugly. Annual Review of Vision Science, 5, 399–426. https://doi.org/10.1146/annurev-vision-091718-014951
- Sherwood, C., Over, J., & Soenen, K. (2021). Structure from motion products associated with uas flights in sandwich, Massachusetts. U.S. Geological Survey data release. https://doi.org/10.5066/P9BFD3YH
- Skalski, P. (2019). Make sense. Retrieved from https://github.com/SkalskiP/make-sense/

BUSCOMBE ET AL. 30 of 31



- Smits, P. C., & Dellepiane, S. G. (1997). Synthetic aperture radar image segmentation by a detail preserving Markov random field approach. IEEE Transactions on Geoscience and Remote Sensing, 35(4), 844–857. https://doi.org/10.1109/36.602527
- Su, H., Wu, L., Jiang, J. H., Pai, R., Liu, A., Zhai, A. J., & DeMaria, M. (2020). Applying satellite observations of tropical cyclone internal structures to rapid intensification forecast with machine learning. *Geophysical Research Letters*, 47(17), e2020GL089102. https://doi.org/10.1029/2020gl089102
- Sugiura, N., & Hosoda, S. (2020). Machine learning technique using the signature method for automated quality control of argo profiles. *Earth and Space Science*, 7(9), e2019EA001019. https://doi.org/10.1029/2019ea001019
- Sultana, F., Sufian, A., & Dutta, P. (2020). Evolution of image segmentation using deep convolutional neural network: A survey. *Knowledge-Based Systems*, 201, 106062. https://doi.org/10.1016/j.knosys.2020.106062
- Sumbul, G., Charfuelan, M., Demir, B., & Markl, V. (2019). BigEarthNet: A large-scale benchmark archive for remote sensing image understanding. In *The 2019 IEEE International geoscience and remote sensing Symposium* (pp. 5901–5904). https://doi.org/10.1109/igarss.2019.8900532
- $Tinoco, R., Goldstein, E., \&\ Coco, G.\ (2015).\ A\ data-driven\ approach\ to\ develop\ physically\ sound\ predictors:\ Application\ to\ depth-averaged\ velocities\ on\ flows\ through\ submerged\ arrays\ of\ rigid\ cylinders.\ \textit{Water}\ Research, 51(2),\ 1247-1263.\ https://doi.org/10.1002/2014wr016380$
- Villmann, T., Merényi, E., & Hammer, B. (2003). Neural maps in remote sensing image analysis. Neural Networks, 16(3-4), 389-403. https://doi.org/10.1016/s0893-6080(03)00021-2
- Vos, K., Harley, M. D., Splinter, K. D., Walker, A., & Turner, I. L. (2020). Beach slopes from satellite-derived shorelines. Geophysical Research Letters, 47(14), e2020GL088365. https://doi.org/10.1029/2020gl088365
- Vosselman, G., Coenen, M., & Rottensteiner, F. (2017). Contextual segment-based classification of airborne laser scanner data. ISPRS Journal of Photogrammetry and Remote Sensing, 128, 354–371. https://doi.org/10.1016/j.isprsjprs.2017.03.010
- Vuolo, F., Zółtak, M., Pipitone, C., Zappa, L., Wenng, H., Immitzer, M., et al. (2016). Data service platform for Sentinel-2 surface reflectance and value-added products: System use and examples. Remote Sensing, 8(11), 938. https://doi.org/10.3390/rs8110938
- Walker, I. J., Davidson-Arnott, R. G., Bauer, B. O., Hesp, P. A., Delgado-Fernandez, I., Ollerhead, J., & Smyth, T. A. (2017). Scale-dependent perspectives on the geomorphology and evolution of beach-dune systems. *Earth-Science Reviews*, 171, 220–253. https://doi.org/10.1016/j. earscirey.2017.04.011
- Warrick, J. A., Ritchie, A. C., Schmidt, K. M., Reid, M. E., & Logan, J. (2019). Characterizing the catastrophic 2017 Mud Creek landslide, California, using repeat structure-from-motion (SfM) photogrammetry. Landslides, 16(6), 1201–1219. https://doi.org/10.1007/s10346-019-01160-4
- Wei, Y., & Ji, S. (2021). Scribble-based weakly supervised deep learning for road surface extraction from remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*.
- Weinstein, B. G. (2018). A computer vision for animal ecology. *Journal of Animal Ecology*, 87(3), 533–545. https://doi.org/10.1111/1365-2656.12780 Wu, X., Liang, L., Shi, Y., & Fomel, S. (2019). FaultSeg3D: Using synthetic data sets to train an end-to-end convolutional neural network for 3D seismic fault segmentation. *Geophysics*, 84(3), IM35–IM45. https://doi.org/10.1190/geo2018-0646.1
- Wulder, M. A., Loveland, T. R., Roy, D. P., Crawford, C. J., Masek, J. G., Woodcock, C. E., et al. (2019). Current status of Landsat program, science, and applications. *Remote Sensing of Environment*, 225, 127–147. https://doi.org/10.1016/j.rse.2019.02.015
- Yang, C., Zhao, H., Bruzzone, L., Benediktsson, J. A., Liang, Y., Liu, B., & Ouyang, Z. (2020). Lunar impact crater identification and age estimation with Chang'E data by deep and transfer learning. *Nature Communications*, 11(1), 1–15. https://doi.org/10.1038/s41467-020-20215-y
- Yao, X., Tham, L., & Dai, F. (2008). Landslide susceptibility mapping based on support vector machine: A case study on natural slopes of Hong Kong, China. *Geomorphology*, 101(4), 572–582. https://doi.org/10.1016/j.geomorph.2008.02.011
- Zhao, J., Wang, X., & Zhou, Y. (2020). A crowdsourcing-based platform for labelling remote sensing images. In The 2020 IEEE International geoscience and remote sensing Symposium (pp. 3227–3230). https://doi.org/10.1109/igarss39084.2020.9323820
- Zhong, Y., Zhao, J., & Zhang, L. (2014). A hybrid object-oriented conditional random field classification framework for high spatial resolution remote sensing imagery. IEEE Transactions on Geoscience and Remote Sensing, 52(11), 7023–7037. https://doi.org/10.1109/tgrs.2014.2306692
- Zuo, R., Xiong, Y., Wang, J., & Carranza, E. J. M. (2019). Deep learning and its application in geochemical mapping. *Earth-Science Reviews*, 192, 1–14. https://doi.org/10.1016/j.earscirev.2019.02.023

#### **References From the Supporting Information**

Dash. (2021). A productive python framework for building web analytic applications. Plotly. Retrieved from https://dash.plotly.com/introduction Grinberg, M. (2018). Flask web development: Developing web applications with python. O'Reilly Media, Inc.

Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., & Oliphant, T. E. (2020). Array programming with NumPv. Nature, 585(7825), 357–362. https://doi.org/10.1038/s41586-020-2649-210.1038/s41586-020-2649-2

Holoviz. (2021). High-level tools to simplify visualization in python. Anaconda, Inc. Retrieved from https://holoviz.org/index.html

Merkel, D. (2014). Docker: Lightweight linux containers for consistent development and deployment. Linux Journal, 2014(239), 2.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., & Grisel, O. (2011). others, Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12, 2825–2830.

Plotly. (2015). Collaborative data science. Plotly Technologies Inc. Retrieved from https://plot.ly

React. (2021). A javascript library for building user interfaces. Facebook Inc. Retrieved from https://reactjs.org/

BUSCOMBE ET AL. 31 of 31