

Guiding Labelling Effort for Efficient Learning With Georeferenced Images

Takaki Yamada¹, Miquel Massot-Campos², Adam Prügel-Bennett³, Oscar Pizarro⁴, *Member, IEEE*, Stefan B. Williams⁵, *Senior Member, IEEE*, and Blair Thornton⁶, *Member, IEEE*

Abstract—We describe a novel semi-supervised learning method that reduces the labelling effort needed to train convolutional neural networks (CNNs) when processing georeferenced imagery. This allows deep learning CNNs to be trained on a per-dataset basis, which is useful in domains where there is limited learning transferability across datasets. The method identifies representative subsets of images from an unlabelled dataset based on the latent representation of a location guided autoencoder. We assess the method's sensitivities to design options using four different ground-truthed datasets of georeferenced environmental monitoring images, where these include various scenes in aerial and seafloor imagery. Efficiency gains are achieved for all the aerial and seafloor image datasets analysed in our experiments, demonstrating the benefit of the method across application domains. Compared to CNNs of the same architecture trained using conventional transfer and active learning, the method achieves equivalent accuracy with an order of magnitude fewer annotations, and 85 % of the accuracy of CNNs trained conventionally with approximately 10,000 human annotations using just 40 prioritised annotations. The biggest gains in efficiency are seen in datasets with unbalanced class distributions and rare classes that have a relatively small number of observations.

Index Terms—Semi-supervised learning, convolutional neural network, autoencoder, georeferenced imagery, pseudo-labelling

1 INTRODUCTION

GEOREFERENCED visual images taken by aircraft, satellites and submersibles are widely used in environmental monitoring. Modern robotic surveys using aerial drones and Autonomous Underwater Vehicles (AUVs) can collect thousands to tens of thousands of georeferenced images in a single mission [1], [2], [3]. As the influx of images gathered by these platforms increases, the need for domain expertise to generate appropriate annotations becomes a bottleneck in our ability to efficiently interpret the data. Supervised machine learning techniques are potentially useful for automated interpretation. However, environmental studies have reported limited transferability of learning from generic training datasets [4], [5], citing the need for application-specific expert-annotated training examples. This is limiting since comprehensive

training datasets do not yet exist for many environmental monitoring applications. The main reasons for this are the high sensitivity of image appearance to environmental conditions (e.g., lighting, atmosphere/water turbidity), observation variables (e.g., range to target, spatial resolution, observation footprint), the large variability in the appearance of unstructured scenes and the complexity of the annotation schemes used in environmental monitoring applications [6], [7]. These factors combined with the large number and different specification of the imaging platforms used (e.g., wavelength sensitivity, dynamic range, illumination source for underwater applications) limit crossover between datasets. Although unsupervised methods can efficiently process large volumes of imagery without relying on human annotations, their outputs typically do not align with the class boundaries of interest to experts, which limits their value for environmental monitoring and infrastructure inspection [8], [9].

This paper develops a novel semi-supervised method that improves learning efficiency when using georeferenced imagery, and reduces the human effort needed to train classifiers for environmental monitoring applications. The method is designed for whole image classification of natural scenes in downward looking imagery and consists of the following parts:

- Unsupervised learning - extracting latent representations of an unlabelled image dataset
- Prioritised labelling - identifying a subset of representative images for human annotation, and assigning predictive pseudo-labels to the remaining data.
- Supervised learning - use of prioritised annotations and pseudo-labels to train CNNs

For unsupervised learning, we investigate the impact on downstream accuracy when two different types of autoencoder are used to learn latent representations. The

- Takaki Yamada, Miquel Massot-Campos, and Adam Prügel-Bennett are with the Faculty of Engineering and Physical Science, University of Southampton, SO16 7QF Southampton, U.K. E-mail: {T.Yamada, miquel.massot-campos}@soton.ac.uk, apb@ecs.soton.ac.uk.
- Oscar Pizarro and Stefan B. Williams are with the Australian Centre for Field Robotics, The University of Sydney, Sydney, NSW 2006, Australia. E-mail: {o.pizarro, stefan.williams}@sydney.edu.au.
- Blair Thornton is with the Faculty of Engineering and Physical Science, University of Southampton, SO16 7QF Southampton, U.K., and also with the Institute of Industrial Science, The University of Tokyo, Tokyo 153-8505, Japan. E-mail: B.Thornton@soton.ac.uk.

Manuscript received 26 Nov. 2020; revised 18 Oct. 2021; accepted 20 Dec. 2021. Date of publication 4 Jan. 2022; date of current version 5 Dec. 2022.

This work was supported in part by U.K. Natural Environment Research Council's Oceanids Biocam under Grant NE/P020887/1 and in part by Australian Research Council's Automated Benthic Understanding Discovery Project under Grant DP190103914.

(Corresponding author: Blair Thornton.)

Recommended for acceptance by Z. Akata.

Digital Object Identifier no. 10.1109/TPAMI.2021.3140060

first uses only the information in images and the second is a location guided autoencoder (LGA) that also uses georeference information to regularise learning [9]. For prioritised labelling, we investigate the impact of using different methods to automatically identify a small subset of images for prioritised annotation and estimate class decision boundaries when assigning predictive pseudo-labels in unannotated images. The prioritised annotations and pseudo-labels can be used to train different CNNs. We analyse the success and sensitivity of the proposed method using four different real-world datasets consisting of tens of thousands of georeferenced environmental monitoring image patches that have expert human labels for training and validation. The gains in learning efficiency are assessed based on the achieved accuracy and number of human annotations used in comparison to CNNs trained using well established transfer and active learning methods.

The advantages of the semi-supervised method for downstream classification tasks are:

- Unlike unsupervised methods, classifier outputs are aligned with class boundaries of interest to humans
- Accurate results can be achieved with significantly reduced human annotation effort compared to conventional supervised methods, and significantly reduced human and computational effort compared to iterative training approaches (e.g., active learning)

The reduction in human effort needed to achieve an equivalent accuracy to current state of the art approaches means that end-to-end training can be achieved on a per-dataset basis, making our approach suitable for use in domains where there is limited transferability of learning between datasets. The rest of this paper is structured as follows; section 2 reviews relevant machine learning literature and section 3 describes the semi-supervised training method. Experimental results for georeferenced seafloor and aerial image datasets are presented in section 4.

2 BACKGROUND

2.1 Machine Learning for Environmental Monitoring

Determining the distribution of land cover, land use, habitats, substrates and infrastructures are tasks that lie at the core of environmental monitoring. One way of achieving these tasks is to interpret imagery using established classification schemes [10], [11], where often a small subset of images are selected for human annotation from which aggregate statistics can be derived. For more comprehensive analysis, many groups have reported automated interpretation of imagery using machine learning, with representative literature described in the following subsections.

2.1.1 Supervised Learning

A large proportion of automated classifiers have used a combination of hand-picked features chosen based on expert knowledge of the application domain or through a reward-based selection process [12], [13]. In [12] the authors apply a Support Vector Machine (SVM) to texture- and colour-based features designed to classify seafloor images into different substrates types for reef ecology surveys. In [14] hand-picked geometric features are combined with SVM for

classification of satellite images. In [13] a similar approach is applied for seafloor mineral prospecting. Spatial invariant features such as Local Binary Patterns (LBP) [15] and Spatial Pyramid Matching (SPM) [16] have also been effectively applied to classification problems for land [17], [18] and seafloor imagery [19], [20]. However, these types of features require manual tuning of parameters, or feature engineering, to efficiently describe each independent dataset. Furthermore, a separate classification process is needed, which typically requires further parameter tuning. As such these methods often require expert knowledge of both the data and application domain, and have limited versatility when applied to multiple datasets.

A key advantage of deep learning techniques is that both the latent representation of data and classification can be simultaneously optimised in a single end-to-end training process. This avoids the need for costly and potentially subjective feature engineering and reduces the need for parameter tuning, making deep learning techniques a compelling choice for image classification tasks. Deep learning techniques are widely used for interpreting aerial and satellite imagery [21]. In [22] the ResNet [23] deep learning CNN is used to classify images of coral into nine separate classes, achieving higher classification resolution than prior studies and demonstrating the ability of deep learning to effectively model class boundaries used in scientific taxonomy. However, to work effectively, deep learning classification techniques typically require a large number of annotated examples of each class. Although labelling platforms tailored to aerial imagery [24] and seafloor imagery exist [2], [5], the sensitivity of images to environmental and acquisition conditions, complexity of annotation schemes and comparatively small size of each environmental monitoring community means that large-scale label repositories such as those in terrestrial imaging [25] and autonomous driving [26] do not yet exist. Several annotated datasets exist for satellite imagery [27]. However, most of these target built environments and artificial objects, and the annotations are not suitable for monitoring and conservation of the natural environment, where standardised but complex hierarchical annotation schemes that consist of hundreds to several thousands of terms are used [6], [7]. Furthermore, for sub-sea imaging, most groups gather images using custom built imaging hardware, where in [28] the authors reported that even small differences in sub-sea imaging hardware limits learning transferability and distorts deep learning classifier outputs. In [29] a pipeline to make training datasets transferable for inference on images from other datasets is proposed for segmentation of marine organism. The work proposes how to reduce scale variance across multiple datasets, which is highlighted as an important consideration for seafloor imagery. A detailed description of this and other domain specific distortions (e.g., blur, haziness, and colour distortion) that affect seafloor imagery can be found in Appendix A, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2021.3140060>.

For datasets where these disturbances are non-negligible, training on a per-dataset basis as is common in unsupervised learning can be considered a potentially effective solution.

Under these constraints, a reasonable approach for effective use of deep learning techniques is to train models on the target dataset itself. However, the implied requirement to annotate large numbers of images every time a new

dataset is obtained is unlikely to be justified for most applications, forming a barrier to wide-spread adoption of deep learning for image interpretation in environmental monitoring applications. This motivates research into techniques for effort reduction.

2.1.2 Unsupervised Learning

Unsupervised learning techniques have great potential for image interpretation in environmental monitoring because they do not require annotations, and so can be efficiently trained and applied on a per-dataset basis. As with any automated image analysis, feature engineering is crucial for effective interpretation. In [8] LBP [15] features derived from greyscale images, 3D rugosity and colour are applied to seafloor image clustering. The authors later applied SPM [16] as a more generic approach to describe seafloor images [30]. These scale invariant features are also used for clustering of aerial and satellite imagery [31]. In [32], the non-parametric Bayesian clustering technique used in [8] and [30] is extended to incorporate annotations made during active learning [33] for seafloor imagery. In [34] the accumulated histogram of oriented gradients from keypoints are used to describe each image, and this is applied to clustering and anomaly detection. More recently, Shields *et al.* [35] used unsupervised clustering results generated from visual images as labels for supervised learning of terrain elevation datasets. To avoid the demanding trial and error process of feature engineering, we developed an unsupervised deep learning LGA in our previous work [9]. The proposed LGA learns latent representations without the need for feature engineering. The georeference information attached to each image is used to regularise learning, allowing CNN architectures to leverage this information and describe patterns that occur on spatial scales larger than a single image frame in a single end-to-end process. Since the LGA does not require any human annotations, it can be efficiently trained and applied on a per-dataset basis, and this has been shown to be effective for clustering and content-based query of seafloor images. Tile2Vec [36] is a method proposed for representation learning of aerial and satellite imagery, where a similar approach based on the physical distances between cropped image patches are leveraged during training.

However, a disadvantage of unsupervised approaches is that the resulting clusters do not attempt to align with the class boundaries of interest to humans, and when latent representations are optimised on a per-dataset basis, it is not possible to make direct comparisons between clusters or perform content-based queries across multiple processed datasets.

2.2 Methods to Reduce Annotation Effort

The shortage of annotations is a common problem when supervised learning is applied to real-world problems, and a number of concepts have emerged to address this issue.

2.2.1 Transfer Learning

Transfer learning allows supervised learning models to be trained using a relatively small number of annotations in the target dataset by making use of much larger annotated

datasets from a different domain. Several frameworks have been proposed to implement this concept [37]. Network-based transfer learning has been applied in many application domains including medical [38], satellite [39], and seafloor imaging [40]. This approach works by reusing networks that have been pre-trained using large, generic datasets (e.g., ImageNet [41], COCO [42], Pascal VOC [43]) that consist of hundreds of thousands to more than ten million labels as an initial model. Though the number of dataset specific annotations needed depends on the domain, number of classes and data augmentation methods used, previous studies on satellite [44] and medical imagery [38] have required several hundreds of domain specific labels for effective use.

2.2.2 Prioritised Labelling

Images in a dataset do not have equal value for CNN training. In [45] the authors demonstrate that training data selection can have a significant impact on learning, where CNNs trained on a well selected subset of annotations can outperform CNNs trained using a larger number of annotations. In [46] annotation efforts are prioritised using k means clustering to estimate the entropy of each sample, showing significant gains in performance compared to random selection.

In active learning [33], the learner interacts with human annotators by iteratively proposing data samples that it considers will most efficiently improve performance. Several strategies have been proposed to achieve this. Most approaches prioritise unlabelled samples that have the highest estimated uncertainty, or are predicted to have the biggest impact on the model. However, the heuristics used to suggest samples can only be calculated after the initial subset has been analysed by the algorithm. Although the initial subset can impact subsequent learning performance, its selection falls outside of the scope of most active learning techniques [33], [47].

In [40] an autoencoder is used to locate objects of interest in an unsupervised manner. The method highlights these regions to human experts in order to facilitate efficient use of time for manual segmentation. The approach leverages the assumption that interesting objects are relatively rare in the original image datasets they are applied to. Regions with a high autoencoder reconstruction loss value are considered likely to include targets of potential interest, and these regions are flagged for prioritised annotation by humans. Active learning is also applied for seafloor image interpretation in [32], [35], where the authors implemented this with SPM as the feature descriptor.

2.2.3 Group Labelling and Label Extrapolation

Group-based labelling [48], [49] is a technique that assigns annotations to subgroups of clustered data in order to reduce the human annotation effort. An advantage of this approach is that it can be applied to datasets with no labels by using unsupervised clustering methods to generate the groups. However, determining the annotation for a cluster of images can be more complex than per-sample based annotation, especially when unsupervised cluster decision boundaries are not aligned with the desired class boundaries, resulting in conflicted human annotations. In [50] the authors modified Gaussian mixture model based clustering to find

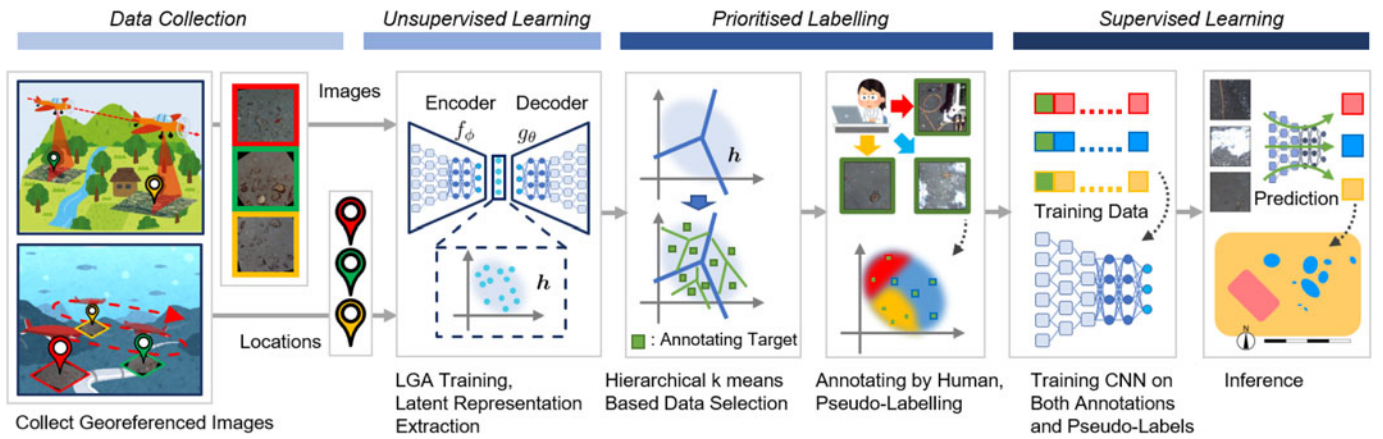


Fig. 1. A flow diagram of the proposed pipeline for LGA driven Semi-Supervised (LGA-SS) training of CNNs. Once a dataset is gathered, the latent representations of the images in the dataset are generated using the LGA [9] (section 3.1), after which hierarchical k means clustering (section 3.2) is used to identify a prioritised subset of images for human annotation. These annotations are used together with a set of algorithmically generated pseudo-labels for the remaining unannotated data to train a CNN that can be used for downstream classification tasks. The proposed LGA-SS method allows a CNN to be trained and applied to classification tasks on a per-dataset basis, making it effective in domains where there is limited transferability of learning between datasets.

clusters with high intra-cluster similarity since the samples in these clusters are considered to be more informative than others. Although these techniques have shown significant improvement in learning efficiency, the underlying assumption is that effective clustering can be achieved.

Predictive pseudo-labelling [51] reduces human effort by first training a classifier on a small subset of data that requires fewer annotations than the target dataset. An advantage of this over group labelling is that annotators consider individual images. After initial training, the classifier predicts labels for the remaining data, and these pseudo-labels are used together with the original annotations to fine-tune a classifier. Li *et al.* [47] reports that SVM and Random Forest classifiers outperform CNNs when generating pseudo-labels from an initial annotated subset. Wu *et al.* [52] uses pseudo-labelling to improve the classification performance for a hyperspectral satellite image dataset, demonstrating effective application of this approach to unstructured environmental monitoring data, where random subsets were used for initial training. The use of prioritisation methods for subset selection in pseudo-labelling has not previously been investigated.

3 EFFICIENT LEARNING IN ENVIRONMENTAL MONITORING IMAGERY

Our aim is to develop a method to efficiently learn class boundaries of interest to humans with fewer annotations than existing methods, and apply this to environmental monitoring image classification problems. Fig. 1 shows the proposed semi-supervised learning pipeline. It learns latent representations of images in a dataset using the LGA [9] (section 3.1). Next, a subset of image samples are selected based on hierarchical k means clustering (section 3.2) for prioritised annotation by humans. Pseudo-labels are assigned to all remaining images (section 3.4) based on the annotated subset. The human annotations and algorithm generated pseudo-labels are then used to fine-tune a CNN, which can be used to solve a downstream classification task. The method is designed to work off-line on a per-dataset basis, once the complete dataset has been gathered. The initial latent representation learning and

identification of prioritised images for labelling are unsupervised, where all images in the dataset are available for these steps without the need for any human input. Human input is only needed to annotate the subset of prioritised images, where the number of prioritised images can be matched and optimised according to the availability of human effort. As such, the method is compatible with post data acquisition workflows associated with environmental survey field work. The LGA driven Semi-Supervised (LGA-SS) method is versatile as it allows a CNN to be both trained and applied to classification on a per-dataset basis, making it effective in domains where the transferability of learning between datasets is limited.

3.1 Location Guided Autoencoder

Patterns of interest in environmental monitoring often occur on spatial scales larger than the image patch size considered by CNNs during their optimisation. The LGA overcomes this problem by introducing georeference regularisation in autoencoder training using a modified loss function [9]. This is designed to reflect the assumption that *two images captured within a close distance tend to look more similar than two that are far away* due to the presence of patterns beyond the footprint of a single image frame. The approach allows the LGA to recognise patterns that recur in images that are close to each other and prioritise these in its learning without introducing artefacts due to imperfect image stitching. The latent representations obtained using the LGA have been shown to perform better than those obtained using a standard convolutional autoencoder when used for clustering and content-based image retrieval [9].

3.2 Data Selection for Prioritised Labelling

The standard CNN learning process expects class-balanced distributions in training datasets. Skewed class distributions, such as those found in natural scenes on land and on the seafloor, can result in overfitting of classes with relatively large numbers of samples. If M images are randomly selected for annotation, training datasets approximate the skewed class distributions of the parent populations,

resulting in non-ideal conditions for training and carrying a risk that smaller classes may not be represented in training for small M values.

In the proposed pipeline, k means clustering is applied to the LGA's latent representation to identify densely populated regions. The number of clusters should be large enough to avoid missing small classes. As long as this condition is satisfied, the outputs are not strongly sensitive to small differences in k as the clusters attempt to evenly represent the different regions of the latent space. In this work we define, $k = \lceil k_e/10 \rceil \times 10$, where k_e is a number of clusters estimated by the elbow method [53]. The value of k is k_e rounded up to the nearest ten. Next, a subset of images for prioritised annotation are selected by taking $\lfloor M/k \rfloor$ or $\lceil M/k \rceil$ images from each cluster so that the total number of images is M . This generates a training class distribution that follows the cluster distribution, which eases the class imbalance problem as long as effective clustering is achieved. The way samples are chosen from within each cluster can also affect learning. In [46] it is assumed that the samples close to the cluster boundaries are important as they have a greater effect on classification decision boundaries. This assumption is reasonable if the boundaries of clustering and classification are comparable, but in situations where class boundaries are ambiguous, like in many environmental monitoring application, it is possible that variability in the annotations will degrade learning performance.

In this study, we consider that the samples provided for training should represent the variability within each cluster in order to deal with situations where the clustering resolution is not sufficient to resolve class boundaries. We implement two approaches to achieve this. The first approach uses k means clustering and randomly samples data from within each cluster so that each cluster in the LGA latent representation is evenly represented in the training data. We also investigate a more structured form of latent space representation, which we implement using hierarchical k means clustering. This approach is originally proposed in [54] where a multi-stage clustering process is introduced. The first stage explores the dominant patterns in the whole dataset, and the following stages attempt to select a representative set of samples from within each cluster. This approach has also been applied to extract representative data in text clustering problems [55]. In this work, we consider that it is important to guarantee that samples are selected from dense regions of the latent representation, and so after the first k means clustering, we generate $\lfloor M/k \rfloor$ or $\lceil M/k \rceil$ sub-clusters within each cluster and select samples that are closest to each sub-cluster centroid so that the total number of samples is M .

3.3 Data Augmentation

Data augmentation [56] plays an important role in reducing the risk of overfitting during CNN training. Since most visual features in downward looking images of the seafloor and of land can be considered invariant to rotation and flipping [57], we apply these augmentations randomly during the training process, together with random shift operations to account for uncertainty in position. These transformations are applied with different parameters (i.e., rotation angle and offset) that are randomly assigned every time an image is fed into the

model during training. Weighted sampling is also applied at each epoch to balance the number of samples in each class. Data augmentation is not applied to colour and scale distortions since it can be consistently corrected taking into account illumination and turbidity conditions and lens distortions [9].

3.4 Pseudo-Labeling

We predict pseudo-labels for each unseen image based on its location relative to annotated samples in the LGA latent space. Although the clustering results used to identify images for prioritised annotation can be used for this purpose, the decision boundaries of clusters and classes are not necessarily aligned. Therefore, we investigate different approaches to estimate class decision boundaries, comparing the performance of nearest neighbour (1-NN), Random Forest and SVM [58] with linear and Radial Basis Function (RBF) kernels as methods capable of expressing varying degrees of complexity of class boundaries in the latent space.

Although the original pseudo-labelling implementation for deep learning applies a single winner takes all class label to unseen data [51], recent research has demonstrated that taking the uncertainty of each pseudo-label into consideration can improve downstream classification accuracy [59], [60]. Class boundaries in environmental monitoring data are often ambiguous and so to address uncertainty near class decision boundaries, we implement probabilistic pseudo-labelling using a Gaussian Process classifier [61] to predict class conditional probability distributions for each sample in the latent space for comparison with the other methods.

Both the annotations and pseudo-labels assigned to the remaining images are used to train CNNs, where for probabilistic pseudo-labelling, the conditional probability distributions are applied to the softmax loss of CNN training in order to describe the pseudo-label uncertainty. The suitability of these classifiers for pseudo-labelling is determined through validation against human annotations.

4 EXPERIMENT

4.1 Dataset

The proposed method is applied to four different environmental monitoring image datasets. Fig. 2 shows the spatial and class distributions of the ground truth for each dataset. The Seafloor dataset (Fig. 2a, see Appendix A, available in the online supplemental material for further details) consists of seafloor visual images collected by an AUV, and the aerial image datasets (Figs. 2b, 2c and 2d, see Appendix B, available in the online supplemental material for further details) are of different types of scene (Mountain, Island and Urban). The class distributions in these spatially continuous datasets are highly skewed compared to the generic datasets that are often used in benchmarking studies. Our experiments consider each class to be of equal importance. The results are assessed based on the macro-averaged F_1 -score, where we take the mean and standard deviation (SD) of 10 repeated sets of experiments under each test configuration.

4.2 Classification With Conventional Classifiers

We investigate the performance of conventional (non-CNN) classifiers in order to generate effective pseudo-labels from a small subset of annotated examples. Five well established

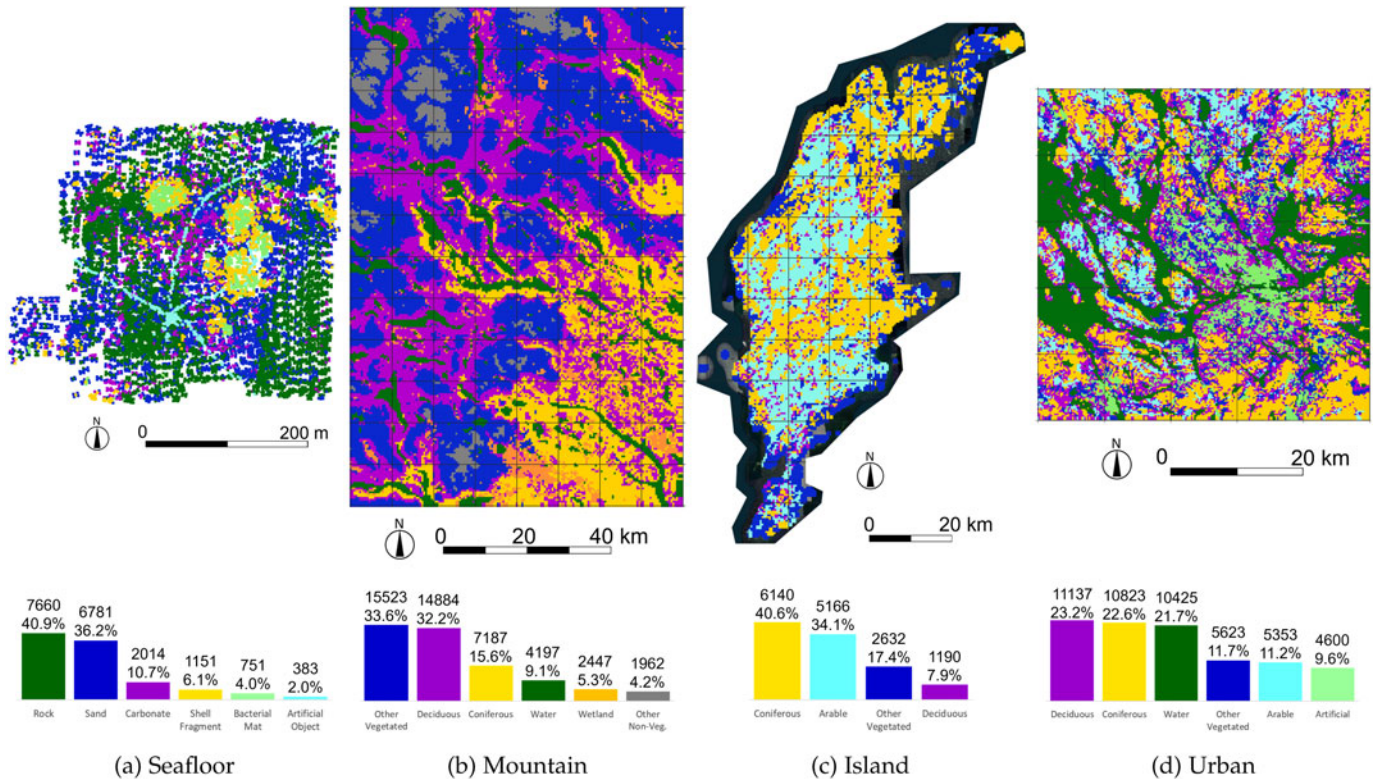


Fig. 2. Spatial patterns (top) and class distributions (bottom) of ground truth classes in four environmental monitoring datasets. Each natural or artificial object class shows a unique spatial pattern in each dataset. The class distributions are highly skewed since all the images in the corresponding areas are included in the datasets without any manual selection process. The Seafloor dataset (2a) consists of colour seafloor imagery collected by an AUV. The three aerial datasets (Mountain, Island and Urban) consist of aerial images cropped from ESRI World Imagery. Details of these datasets can be found in Appendices A and B, available in the online supplemental material, respectively.

classifiers; k -NN with $k = 1$ (1-NN), Random Forest (RF), SVM with linear (L-SVM) and RBF kernels (R-SVM) [58] and Gaussian Process (GP) [61] classifiers are applied to the latent space mapped by an LGA that has been trained on all available image patches. The results are compared with those of a standard convolutional autoencoder that uses the same architecture as the LGA except for the georeference regularisation. To evaluate the performance with a small number of annotations, an adjusted cross-validation is applied. First, half of the annotated image patches are randomly selected as a test subset, preserving the class distribution of the entire dataset in each dataset. Then M images are selected from the remaining patches based on random selection, k means based selection, and the proposed hierarchical k means based selection. Following the equation defined in section 2.2.2, $k = 20$ is used for both k means and hierarchical k means based selection for all the datasets. In k means based selection, $M/20$ images are selected randomly from each cluster. In hierarchical k means based selection, the second stage k means is applied to each cluster to find $M/20$ sub-cluster centroids, and the images closest to each centroid are selected for annotation. Training and testing are executed ten times for each configuration with $M = 20, 40, 100, 200, 400, 1000$ and 7500 (for the aerial datasets) or 9370 (for the Seafloor dataset).

Tables 1 and 2 show the mean and SD of the F_1 -scores for the ten-time cross-validation with each configuration (A1 - A20 in Table 1 and A'1 - A'20 in Table 2) on the seafloor and aerial datasets, respectively. The data selection strategy has a greater impact on performance than the choice of classifier,

with all classifiers benefiting significantly from hierarchical k means prioritisation. The relative gains in accuracy compared to random selection are especially large for small values M (20, 40 and 100), confirming the importance of the data selection strategy when training with a small number of annotations. For the Seafloor dataset (Table 1), the combination of LGA based pre-training and hierarchical k means based data selection with a R-SVM (configuration A14) performs the best among the tested cases for all values of M . The L-SVM and GP generally perform better than 1-NN and RF, where the L-SVM tends to be better for small values of M and GP better for larger M . A similar trend is observed with the aerial datasets (Table 2). For small values of M , L-SVM outperforms R-SVM; however, the difference is marginal. The largest efficiency gains are achieved in the datasets that have rare classes with the smallest number of relative observations (i.e., Seafloor and Mountain).

The standard deep learning autoencoder (configuration A16 - A20 in Table 1 and A'16 - A'20 in Table 2) is significantly less effective than the LGA for all the datasets investigated in this work. This is an expected result since our previous work has already shown that the autoencoder achieves poor clustering performance without georeference regularisation [9], and the underlying assumption behind the data selection strategies investigated here is that effective clustering can be achieved. The results demonstrate that the proposed location guided latent representation learning and representative image selection are effective for environmental applications using georeferenced image datasets across application domains.

the LGA, indicating that the advantage of LGA pre-training is lost when all the layers are trained. The fact that B4 generally outperforms these configurations demonstrates the advantage of embedding georeference information through LGA pre-training using the target dataset. When $M = 9370$, B6, corresponding to the case where all the layers of an ImageNet pre-trained ResNet18 are trained on the dataset, shows the best accuracy. This suggests that ResNet18's deeper architecture and use of residual blocks allows for better performance than AlexNet when a sufficient number of training examples is available. However, B4 is the best option overall for $M \leq 1000$, which is significant for this study since we are interested in efficient training with a small number of annotated examples.

The comparison between B1 to B3 (last layer only) and B5 to B7 (all layer) indicates that training only the last layer limits the performance of each architecture for large values of M , indicating that there is a significant difference between the low-level and mid-level features of ImageNet and the environmental monitoring dataset. In the proposed pipeline, the number of training examples can be considered large due to the use of pseudo-labels. Therefore, we choose to investigate B6 as it demonstrate the best capacity for learning among B1 to B8, and we also examine B4 since it is the most efficient learner for $M \leq 1000$.

4.3.2 Active Learning Comparison (C1-C12)

Active learning methods attempt to improve learning efficiency by training classifiers on a subset of annotated samples, and proposing which samples should be annotated next based on their prediction uncertainty [33]. CNNs are well suited to this iterative process of prediction and prioritised annotation as their outputs are already conditional probabilities against labels and so uncertainty metrics can be easily derived. Common strategies for uncertainty based prioritisation include Least Confidence (LC) sampling, margin sampling and entropy based sampling, all of which have previously been demonstrated to be effective for environmental monitoring applications [32].

Conventional active learning starts the iterative training process with a randomly selected subset of samples. However, its performance is sensitive to this initial selection and so we investigate whether an initial selection of samples nearest to the centroids of the k means clusters in the LGA latent space improves their performance. Subsequent batches of samples (20 when $M \leq 1000$ or 1000 when $M > 1000$) are selected based on the active learning query strategies and iteratively added to the subset of annotated samples for training. A training epoch of 10 was chosen so that the total number of epochs is comparable to the standard supervised learning results (B1-8) and proposed methods (D1-D8).

In our experiment, we assess two different CNN architectures (AlexNet and ResNet18), and compare the performance of three well established active learning iterative sampling techniques (LC sampling, margin sampling and entropy based sampling). The active learning process is initialised using two different initial subset selection methods; first where the initial subset is randomly sampled (corresponding to traditional active learning workflows), and

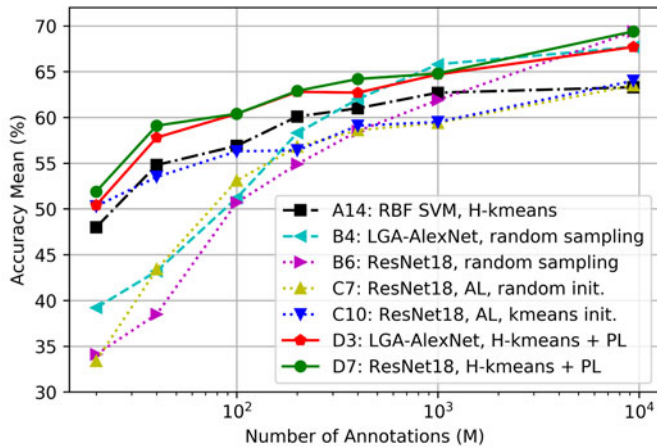
second where active learning initialised by a k means centroid based sample initialisation (taking advantage of the georeference embedded latent representations learnt during LGA pre-training).

Configuration C1 to C12 in Table 3 show the accuracy scores for CNNs trained using the different configurations for active learning. Comparing the LGA pre-trained AlexNet configurations (C1 to C3) with their transfer learning counterpart (B4) shows that the active learning reduces accuracy. However, for ResNet18, the accuracy increases when active learning is applied (B6 and C7 to C9) for small values of $M < 1000$. It is noticeable that for larger M (particularly $M = 9370$), active learning degrades performance, possibly due to overfitting of CNN weights at an early phase of the iterative learning process trapping them in local minima. This is because the CNN is trained sequentially on discrete subsets of data, where the stored weights are used to initialise the optimisation of the next subset to limit the total number of training epochs required [63]. Although overfitting is potentially mitigated by resetting the CNN weights between each training subset [64], this requires a large number of training epochs, making it impractical for use in domains that require per-dataset training.

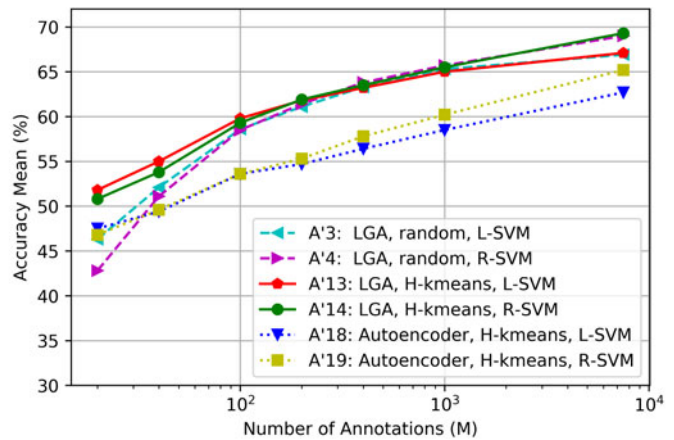
The use of the LGA k means centroids for initial sample selection significantly improves performance (C4 to C6 and C10 to C12), where the gains are largest for small numbers of training examples, i.e., $M \leq 100$. Although this advantage is lost as M increases, it does not cause any significant degradation in performance compared with the random initial subset selection. The difference between the active learning strategies is marginal for both the random and k means initial selection. Although different hyperparameters (e.g., number of epochs for each iteration) may improve active learning performance, optimisation of these is outside the scope of this work since there are no systematic methods available to determine them.

4.3.3 Data Selection Strategy Comparison (D1-D8)

Four data selection strategies; k means, hierarchical k means, and hierarchical k means with pseudo-labelling or probabilistic pseudo-labelling, are validated in this section. The previous section already confirmed that hierarchical k means based data selection is effective for small values of M when combined with conventional non-CNN classifiers. In order to allow for fair comparison, the number of training samples used by the CNN at each training epoch is fixed to the total number of available labelled training image patches (i.e., 9370 in this experiment). For configurations where all available labelled image patches are used in the training (i.e., all pseudo-label and probabilistic pseudo-label configurations and where $M = 9370$ without pseudo or probabilistic pseudo-labelling), each original labelled training image patch is used once, and these samples are individually subjected to data augmentations that randomise orientations, flipping and position offsets at each training epoch before being used by the CNN. For configurations where the number of labelled image patches used in the training is less than available labelled training image patches (i.e., $M < 9370$ with no pseudo or probabilistic pseudo labels), the selected original images are sampled multiple times (i.e.,



(a) Annotations (M) vs F_1 (macro-average) mean for the Seafloor Dataset



(b) Annotations (M) vs F_1 (macro-average) means averaged over the Mountain, Island and Urban Aerial Datasets

Fig. 3. Comparison of classification performance investigated in section 4. Mean of F_1 (macro-average) values against each M are shown. Representative configurations are chosen from Tables 1 and 3 for Fig. 3a and Table 2 for Fig.3b. The proposed georeference embedded sample selection method improves performance for all the datasets analysed in our experiments. Larger gains in learning efficiency are achieved in datasets that have a more heavily skewed class distribution.

approximately $9370/M$ times) so that a fixed number of labelled training samples are provided to the CNN, where each sample is subjected to random data augmentation before being used by the CNN at each training epoch. In [51], pseudo-labels are determined by the k means clustering results, corresponding to 1-NN in Table 1. However, Table 1 shows that R-SVM consistently estimates better class decision boundaries, and so will be used to assign predictive pseudo-labels in this work. Although the GP classifier described in section 3.4 did not perform as well as the R-SVM, the prediction uncertainty may be useful for CNN training and so experiments are also performed using these outputs as probabilistic pseudo-labels.

Configuration D1 to D4 in Table 3 shows the performance metrics for each data selection strategy with the LGA pre-trained AlexNet CNN with last layer supervised training. D5 to D8 show the same comparison for ImageNet pre-trained ResNet18 CNN with all layer supervised training, where these base configurations were chosen since they performed best in our CNN architecture comparison (B4 and B6 in section 4.3.1). For both AlexNet and ResNet18, the combination of hierarchical k means and pseudo-labelling achieves the best performance for $M \leq 200$. Comparing the cases with pseudo-labelling (D3 and D7) to the cases without (D2 and D6) shows that pseudo-labelling consistently improves classification performance. D7, which applies hierarchical k means and pseudo-labelling to ResNet18, performs the best for $M \leq 200$ among all the configurations in Table 3. The accuracies achieved by D7 with $M = 20, 40, 100$ are similar to the metrics achieved for B1 to B4 with $M = 200, 400, 1000$, which have an order of magnitude more annotations. In particular, B6 and D7 use the same CNN architecture, showing that gains in learning efficiency can be attributed to the LGA-SS training method, resulting in a significant reduction in human effort to achieve a similar level of classification accuracy. Although the efficiency gains diminish as the number of human annotations available for training increases, the LGA-SS method never degrades the CNN's performance for an equal number of

annotations. Another way to look at this is that the largest gains in learning efficiency are achieved when there is only a small amount of human effort available for annotation tasks, where D7 with 40 prioritised annotations reaches 85 % of the accuracy achieved by the best performing supervised CNN, B6, trained using 9370 human annotations, which represents just 0.4 % of the human effort. The data also shows that the combination of hierarchical k means and pseudo-labelling improves the repeatability between experiments under the same conditions, which is an important attribute for practical application of automated data interpretation.

Probabilistic pseudo-labelling outperforms pseudo-labelling only when $M = 1000$. This indicates that meaningful probabilistic expression of pseudo-labels can only be taken advantage of when a relatively large number of annotations are available. On the other hand D2, where pseudo-labelling is not applied, shows the best accuracy for $M = 1000$, and similarly D1 shows the best performance for $M = 400$ with D2 following it. This trend suggests that LGA pre-trained AlexNet is effective at describing the class boundaries when a sufficient number of annotated examples can be provided for fine-tuning. The equivalent training approach for D5 and D6 does not show this behaviour, indicating that this is a particular feature of using the LGA pre-trained network. The advantages of the proposed method with hierarchical k means for prioritised sample annotation and pseudo-labelling using R-SVM is significant for $M \leq 200$ for both CNN architectures (i.e., D3 and D7).

4.4 CNN and Conventional Classifier Comparison

Fig. 3 compares the performance metrics of several representative configurations in Tables 1, 2 and 3. The result under configuration A14 are shown as this is the best performing conventional (i.e., non-CNN) classifier. For the CNN classifiers, configurations B4, B6, C7, C10, D3 and D7 are shown to demonstrate the effectiveness of the proposed pipeline compared to other data selection strategies (random selection and active learning) in Fig. 3a. Fig. 3b shows

significantly improved even though almost 10 times the number of annotations are used for training. On the other hand, the CNNs achieves statistically significant increases from $M = 1000$ to $M = 9370$ in all cases. This supports the common understanding that deep learning CNNs are a better option than conventional classifiers when large training datasets are available, and that conventional classifiers are a reasonable option when only a small number of annotations are available for training.

Active learning (C7 and C10) benefits from LGA based k means initialisation (C10), and shows better accuracy than standard training (B4 and B6) for small M , but the performance degrades when M is large due to overfitting as discussed previously. The proposed pipeline with prioritised annotation and pseudo-labelling significantly outperforms active learning for all M and both CNN architectures (D3, D7). Pseudo-labelling is more robust to overfitting than active learning since variability within the dataset is fully represented as all the available images are used for training.

Other factors that are important for practical application include the computational cost and the requirements for human input. Compared with CNNs, conventional classifiers require less time for training once LGA latent representations are generated and annotations have been made. In active learning, the three main steps i.e., training with annotated samples, inference for prioritising samples without annotations and annotating by humans, need to be repeated in sequence. This results in a large computational cost and also leads to inefficiencies as human annotators are forced to work around classifier retraining at each iteration. On the other hand, the time investment needed for the proposed pipeline is similar to conventional CNN training, since the unsupervised training and LGA based sample prioritisation do not require any human input, and the computation time for predicting pseudo-labels is negligible.

4.5 Per-Class Performance

So far the macro-averaged F_1 score has been used as a metric to compare the overall performance of different classifiers. This is appropriate when we assume all classes in a dataset are of equal importance. However, there are applications where this is not the case, and in these scenarios it is more valuable to consider performance on a per-class basis. Figs. 4 and 5 compare the per-class confusion matrices for M values of 20, 40, 100 and 1000 for configurations B2 and D7. These represent the outputs of the best performing network, ResNet18, trained using standard transfer learning and the proposed LGA-SS pipeline, respectively. The values in each confusion matrix are normalised by the number of ground truth annotations so that the diagonal elements correspond to the recall value of each class. The confusion matrices corresponding to the trials with the closest F_1 score (macro-average) to the mean of ten repetitions (Table 3) are chosen for each value of M . The values of zeros for $M = 20$ and 40 in Fig. 4 suggest no images corresponding to 'Artificial Object' were selected in the random selection used for training and so predictions could not be made effectively for this class. On the other hand, Fig. 5 shows that all 6 classes in the dataset are predicted for all M , illustrating the advantage of using hierarchical k means based data selection to avoid minor

classes from being overlooked even when the total number of annotated images is small.

Comparing the habitat maps generated using the classification results to the ground truth annotations (Fig. 2b) shows that the random data selection (Fig. 4) requires a larger number of training samples M to capture the different spatial distribution patterns of each class. Using the proposed LGA-SS training method (Fig. 5) results in more consistent per-class performance, providing a better approximation of the ground truth class distribution patterns even for small values of M . The consistent performance for different numbers of input training data is an important attribute for practical application since the annotation resource available for different datasets is likely to vary. These points favour the proposed method over random sampling approaches that are more sensitive to the number of available annotations, and require larger amounts of training data to achieve similar performance.

5 CONCLUSION

This paper proposes a novel semi-supervised learning pipeline to classify georeferenced imagery using deep learning CNNs. The main advantage of the proposed LGA-SS method is that it can interpret images according to class boundaries of interest for environmental monitoring more efficiently than the alternative methods tested in this work, requiring less human effort and achieving better accuracy. The method is designed for per-dataset training in order to achieve high performance with a realistic investment of human effort for practical application. Experiments on four georeferenced image datasets spanning aerial and seafloor environments show that the proposed georeference embedding and sample selection methods are effective across application domains, achieving the largest gains in efficiency are achieved on datasets that have highly skewed class distributions, which are a common feature in environmental monitoring applications. Other relevant advantages include reduced variability between multiple end-to-end training and classification runs under the same configurations, and more consistent performance with different sizes of input training data compared to traditional naive (i.e., random sampling) based transfer learning methods. These properties make the LGA-SS method suitable for use in domains where there is limited transfer of learning between datasets. Our results demonstrate that:

- The proposed LGA-SS can achieve classification accuracy equivalent to naively trained CNNs with an order of magnitude fewer human annotations (i.e., tens to hundreds, as opposed to thousands). The results demonstrate improvements in accuracy by a factor of 1.2 to 1.5 when a hundred or less annotations are used, where the largest gains in learning efficiency are achieved with small numbers of annotations. The method also reduces the statistical variability between independent trials under the same learning configurations to approximately 0.6 of that when random sampling is used. The proposed method reaches 85 % of the accuracy achieved by the best performing naively trained CNN (trained using 9370 human annotations) with just 40 prioritised annotations, which represents 0.4 % of the human effort.

- The strategy to select data for human annotation affects final classification performance. On the four datasets, introducing structure to prioritise annotation effort using hierarchical k means in the latent representation shows an average of 1.12 times improvement, and leveraging LGA instead of an autoencoder with the same CNN architecture achieved 1.23 times higher accuracy in terms of R-SVM classification results when the number of annotations is less than 100. A similar gain in performance is seen when the LGA based k means selection is used to initialise active learning, with a 1.25 factor improvement compared to equivalent randomly initialised active learning setups.
- The proposed method makes more efficient use of human effort than traditional active learning based techniques tested in this work, and is less prone to overfitting, achieving a factor 1.12 and 1.22 improvement in performance for AlexNet and ResNet18 respectively when compared to randomly initialised active learning across all values of M .
- CNN architectures are able to generalise class boundaries of interest to humans even when pseudo-labels are assigned to all data in a training set. The resulting CNN is able to improve the relative classification accuracy by an average of 6.4% compared to the classification accuracy of the pseudo-labels themselves.
- The performance of conventional classifiers for pseudo-label generation is significantly improved using k means based selection compared to random selection when generating subsets of data for annotation. A factor of 1.30 improvement in classification accuracy is achieved for prioritised subsets with a hundred samples or less.
- Implementation of annotation effort prioritisation strategies relies on effective unsupervised clustering performance for seafloor images, where the use of georeferencing information by the LGA compared to an equivalent autoencoder that only uses information in images resulted in an improvement in classification accuracy by a factor of 1.4 to 8.9 (average 3.1) for the configurations tested in this work.

ACKNOWLEDGMENTS

Data was collected during the Schmidt Ocean Institute's #Adaptive Robotics campaign. The University of Southampton IRIDIS High Performance Computing Facility was used in this work.

REFERENCES

- [1] I. Colomina and P. Molina, "Unmanned aerial systems for photogrammetry and remote sensing: A review," *ISPRS J. Photogrammetry Remote Sens.*, vol. 92, pp. 79–97, 2014.
- [2] M. Bewley *et al.*, "Australian sea-floor survey data, with images and expert annotations," *Sci. Data*, vol. 2, 2015, Art. no. 150057.
- [3] B. Thornton *et al.*, "Biometric assessment of deep-sea vent megabenthic communities using multi-resolution 3D image reconstructions," *Deep Sea Res. Part I: Oceanogr. Res. Papers*, vol. 116, pp. 200–219, 2016.
- [4] S. Kentsch, M. L. Lopez Caceres, D. Serrano, F. Roure, and Y. Diez, "Computer vision and deep learning techniques for the analysis of drone-acquired forest images, a transfer learning study," *Remote Sens.*, vol. 12, no. 8, 2020.
- [5] D. Langenkämper, M. Zurowietz, T. Schoening, and T. W. Nattkemper, "BIIGLE 2.0 - browsing and annotating large marine image collections," *Front. Mar. Sci.*, vol. 4, 2017, Art. no. 83. [Online]. Available: <https://www.frontiersin.org/article/10.3389/fmars.2017.00083>
- [6] F. Althaus *et al.*, "A standardised vocabulary for identifying benthic biota and substrata from underwater imagery: The catami classification scheme," *PLoS ONE*, vol. 10, no. 10, pp. 1–18, 2015.
- [7] J. N. Gomes-Pereira *et al.*, "Current and future trends in marine image annotation software," *Prog. Oceanogr.*, vol. 149, pp. 106–120, 2016.
- [8] D. Steinberg, A. Friedman, O. Pizarro, and S. B. Williams, "A Bayesian nonparametric approach to clustering data from underwater robotic surveys," in *Proc. Int. Symp. Robot. Res.*, 2011, vol. 28, pp. 1–16.
- [9] T. Yamada, A. Prügel-Bennett, and B. Thornton, "Learning features from georeferenced seafloor imagery with location guided autoencoders," *J. Field Robot.*, vol. 38, no. 1, pp. 52–67, 2020.
- [10] K. J. Morris *et al.*, "A new method for ecological surveying of the abyss using autonomous underwater vehicle photography," *Limnol. Oceanogr.: Methods*, vol. 12, no. 11, pp. 795–809, 2014.
- [11] J. Escartín *et al.*, "Globally aligned photomosaic of the lucky strike hydrothermal vent field (mid-atlantic ridge, 3718.5' N): Release of georeferenced data, mosaic construction, and viewing software," *Geochem. Geophys., Geosyst.*, vol. 9, no. 12, pp. 1–17, 2008.
- [12] O. Beijbom, P. J. Edmunds, D. I. Kline, B. G. Mitchell, and D. Kriegman, "Automated annotation of coral reef survey images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2012, pp. 1170–1177.
- [13] U. Neettiyath *et al.*, "Deep-sea robotic survey and data processing methods for regional-scale estimation of manganese crust distribution," *IEEE J. Ocean. Eng.*, vol. 46, no. 1, pp. 102–114, Jan. 2021.
- [14] J. Inglada, "Automatic recognition of man-made objects in high resolution optical remote sensing images by SVM classification of geometric image features," *ISPRS J. Photogrammetry Remote Sens.*, vol. 62, no. 3, pp. 236–248, 2007.
- [15] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 7, pp. 971–987, Jul. 2002.
- [16] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 1794–1801.
- [17] C. Chen, B. Zhang, H. Su, W. Li, and L. Wang, "Land-use scene classification using multi-scale completed local binary patterns," *Signal, Image Video Process.*, vol. 10, no. 4, pp. 745–752, 2016.
- [18] S. Chen and Y. Tian, "Pyramid of spatial relations for scene-level land use classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 4, pp. 1947–1957, Apr. 2015.
- [19] M. Bewley, N. Nourani-Vatani, D. Rao, B. Douillard, O. Pizarro, and S. B. Williams, "Hierarchical classification in AUV imagery," in *Proc. Field. Serv. Robot.*, 2015, pp. 3–16.
- [20] D. Rao, M. De Deuge, N. Nourani-Vatani, S. B. Williams, and O. Pizarro, "Multimodal learning and inference from visual and remotely sensed data," *Int. J. Robot. Res.*, vol. 36, no. 1, pp. 24–43, 2017.
- [21] L. Ma, Y. Liu, X. Zhang, Y. Ye, G. Yin, and B. A. Johnson, "Deep learning in remote sensing applications: A meta-analysis and review," *ISPRS J. Photogrammetry Remote Sens.*, vol. 152, pp. 166–177, 2019.
- [22] A. Mahmood *et al.*, "Deep image representations for coral image classification," *IEEE J. Ocean. Eng.*, vol. 44, no. 1, pp. 121–131, Jan. 2019.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [24] D. Hou, Z. Miao, H. Xing, and H. Wu, "V-RSIR: An open access web-based image annotation tool for remote sensing image retrieval," *IEEE Access*, vol. 7, pp. 83852–83862, 2019.
- [25] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [26] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *Int. J. Robot. Res.*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [27] A. Van Etten, D. Lindenbaum, and T. M. Bacastow, "SpaceNet: A remote sensing dataset and challenge series," *arXiv:1807.01232*.

- [28] D. Langenkämper, R. van Kevelaer, A. Purser, and T. W. Nattkemper, "Gear-induced concept drift in marine images and its effect on deep learning classification," *Front. Mar. Sci.*, vol. 7, 2020, Art. no. 506.
- [29] M. Zurowietz and T. W. Nattkemper, "Unsupervised knowledge transfer for object detection in marine environmental monitoring and exploration," *IEEE Access*, vol. 8, pp. 143558–143568, 2020.
- [30] D. Steinberg, "An unsupervised approach to modelling visual data," PhD Thesis, Univ. Sydney, Camperdown, NSW, Australia, 2013.
- [31] T. Vigneshl and K. Thyagarajan, "Local binary pattern texture feature for satellite imagery classification," in *Proc. Int. Conf. Sci. Eng. Manage. Res.*, 2014, pp. 1–6.
- [32] A. Friedman, D. Steinberg, O. Pizarro, and S. B. Williams, "Active learning using a variational dirichlet process model for pre-clustering and classification of underwater stereo imagery," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2011, pp. 1533–1539.
- [33] B. Settles, "Active learning literature survey," Univ. Wisconsin-Madison, Dept. Comput. Sci., Madison, WI, USA, Tech. Rep. TR1648, 2009.
- [34] J. W. Kaeli and H. Singh, "Online data summaries for semantic mapping and anomaly detection with autonomous underwater vehicles," in *Proc. OCEANS 2015-Genova*, 2015, pp. 1–7.
- [35] J. Shields, O. Pizarro, and S. B. Williams, "Towards adaptive benthic habitat mapping," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2020, pp. 9263–9270.
- [36] N. Jean, S. Wang, A. Samar, G. Azzari, D. Lobell, and S. Ermon, "Tile2Vec: Unsupervised representation learning for spatially distributed data," *Proc. AAAI Conf. Artif. Intell.*, vol. 33, pp. 3967–3974, 2019.
- [37] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu, "A survey on deep transfer learning," in *Proc. Int. Conf. Artif. Neural Netw.*, 2018, pp. 270–279.
- [38] H.-C. Shin *et al.*, "Deep convolutional neural networks for computer-aided detection: CNN architectures, dataset characteristics and transfer learning," *IEEE Trans. Med. Imag.*, vol. 35, no. 5, pp. 1285–1298, May 2016.
- [39] X. Yao, J. Han, G. Cheng, X. Qian, and L. Guo, "Semantic annotation of high-resolution satellite images via weakly supervised learning," *IEEE Trans. Geosci. Remote Sen.*, vol. 54, no. 6, pp. 3660–3671, Jun. 2016.
- [40] M. Zurowietz, D. Langenkämper, B. Hosking, H. A. Ruhl, and T. W. Nattkemper, "MAIA—A machine learning assisted image annotation method for environmental monitoring and exploration," *PLoS One*, vol. 13, no. 11, 2018, Art. no. e0207498.
- [41] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [42] T. Lin *et al.*, "Microsoft COCO: common objects in context," 2014, *arXiv:1405.0312* [Online]. Available: <http://arxiv.org/abs/1405.0312>
- [43] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal visual object classes challenge: A retrospective," *Int. J. Comput. Vis.*, vol. 111, no. 1, pp. 98–136, Jan. 2015.
- [44] D. Marmanis, M. Datcu, T. Esch, and U. Stilla, "Deep learning earth observation classification using imagenet pretrained networks," *IEEE Geosci. Remote Sen. Lett.*, vol. 13, no. 1, pp. 105–109, Jan. 2016.
- [45] A. Lapedriza, H. Pirsiavash, Z. Bylinskii, and A. Torralba, "Are all training examples equally valuable?," 2013, *arXiv:1311.6510*.
- [46] S. Paul, J. H. Bappy, and A. K. Roy-Chowdhury, "Efficient selection of informative and diverse training samples with applications in scene classification," in *Proc. IEEE Int. Conf. Image Process.*, 2016, pp. 494–498.
- [47] Z. Li, B. Ko, and H.-J. Choi, "Naive semi-supervised deep learning using pseudo-label," *Peer-to-Peer Netw. Appl.*, vol. 12, no. 5, pp. 1358–1368, 2019.
- [48] D. Dai, M. Prasad, C. Leistner, and L. Van Gool, "Ensemble partitioning for unsupervised image categorization," in *Proc. 12th Eur. Conf. Comput. Vis.*, 2012, pp. 483–496.
- [49] M. Wigness, B. A. Draper, and J. R. Beveridge, "Efficient label collection for unlabeled image datasets," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 4594–4602.
- [50] Y. Tian, W. Liu, R. Xiao, F. Wen, and X. Tang, "A face annotation framework with partial clustering and interactive labeling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2007, pp. 1–8.
- [51] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," *Proc. Workshop Challenges Representation Learn.*, vol. 3, no. 2, p. 896, 2013.
- [52] H. Wu and S. Prasad, "Semi-supervised deep learning using pseudo labels for hyperspectral image classification," *IEEE Trans. Image Process.*, vol. 27, no. 3, pp. 1259–1270, Mar. 2018.
- [53] V. Satopaa, J. Albrecht, D. Irwin, and B. Raghavan, "Finding a 'kneedle' in a haystack: Detecting knee points in system behavior," in *Proc. 31st Int. Conf. Distrib. Comput. Syst. Workshops*, 2011, pp. 166–171.
- [54] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 2, pp. 2161–2168, 2006.
- [55] H. S. Gowda, M. Suhil, D. Guru, and L. N. Raju, "Semi-supervised text categorization using recursive k-means clustering," in *Proc. Int. Conf. Recent Trends Image Process. Pattern Recognit.*, 2016, pp. 217–227.
- [56] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *J. Big Data*, vol. 6, no. 1, 2019, Art. no. 60.
- [57] G. Cheng, P. Zhou, and J. Han, "Learning rotation-invariant convolutional neural networks for object detection in VHR optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 12, pp. 7405–7415, Dec. 2016.
- [58] J. Friedman, T. Hastie, and R. Tibshirani, *The Elements of Statistical Learning*. Berlin, Germany: Springer, 2001.
- [59] A. Iscen, G. Toliás, Y. Avrithis, and O. Chum, "Label propagation for deep semi-supervised learning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 5070–5079.
- [60] E. Arazo, D. Ortego, P. Albert, N. E. O'Connor, and K. McGuinness, "Pseudo-labeling and confirmation bias in deep semi-supervised learning," in *Proc. Int. Joint Conf. Neural Netw.*, 2020, pp. 1–8.
- [61] C. K. Williams and C. E. Rasmussen, *Gaussian Processes for Machine Learning*. Cambridge, MA, USA: MIT Press, 2006.
- [62] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.
- [63] Z. Zhou, J. Shin, L. Zhang, S. Gurudu, M. Gotway, and J. Liang, "Fine-tuning convolutional neural networks for biomedical image analysis: Actively and incrementally," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 7340–7351.
- [64] Y. Gal, R. Islam, and Z. Ghahramani, "Deep Bayesian active learning with image data," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 1183–1192.



Takaki Yamada received the BEng and MSc degrees from the University of Tokyo, Japan, in 2009 and 2011, respectively, and the PhD degree from the University of Southampton, U.K., in 2021. He is currently a member of Centre for In Situ and Remote Intelligent Sensing. His research interests include sensing and perception for autonomous underwater vehicles.



Miquel Massot-Campos received the MEng degree from BarcelonaTech, Spain, in 2011, and the MSc and PhD degrees from the University of the Balearic Islands, Spain, in 2013 and 2019, respectively. He is currently with the University of Southampton, U.K. and the member of the Centre for In Situ. His research interests include intelligent sensing and scalability of autonomous underwater vehicle's missions.



Adam Prügel-Bennett received the BSc degree in physics from Southampton University in 1984 and the PhD degree in theoretical physics from Edinburgh University in 1989. He worked in research jobs in Oxford, Paris, Manchester, Copenhagen and Dresden before finally returning to Southampton. He is currently a professor of electronics and computer science. His research interests include machine learning and particularly deep learning.



Oscar Pizarro (Member, IEEE) received the BSc degree in electronic engineering from the Universidad de Concepcion in 1997, and the dual MSc degree in ocean engineering and electrical engineering and computer science and PhD degree in oceanographic engineering from the MIT-WHOI Joint Program, in 2003 and 2004, respectively. He joined the University of Sydney's Australian Centre for Field Robotics in 2005, where he is currently a principle research fellow. His research interests include scalable approaches to seafloor imaging and habitat characterisation.



Stefan B. Williams (Senior Member, IEEE) received the BSc degree in systems engineering design from the University of Waterloo in 1997 and the PhD degree in field robotics from the University of Sydney in 2002. He is currently a professor of marine robotics with the University of Sydney's Australian Centre for Field Robotics and the head of Australia's Integrated Marine Observing System AUV Facility. His research interests include simultaneous localization and mapping in underwater environments, autonomous navigation, and data interpretation.



Blair Thornton (Member, IEEE) received the BEng degree in naval architecture and the PhD degree in underwater robotics from Southampton University in 2002 and 2006, respectively. In 2003, he joined the Underwater Robotics and Application Lab, Institute of Industrial Science, UTokyo, before rejoining Southampton in 2016 where he is currently a professor of marine autonomy. He has more than 450 days with sea deploying robotic systems and is dedicated to generating data and insight in marine science through improved sensing and autonomy.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.