

RESEARCH ARTICLE

ML-morph: A fast, accurate and general approach for automated detection and landmarking of biological structures in images

Arthur Porto  | Kjetil L. Voje 

Centre for Ecological and Evolutionary
Synthesis, University of Oslo, Oslo, Norway

Correspondence

Arthur Porto

Email: agporto@wustl.edu

Present address

Arthur Porto, Department of Biosciences,
University of Oslo, P.O. Box 1066 Blindern,
NO-0316, Oslo, Norway

Funding information

Research Council of Norway, Grant/Award
Number: 249961; European Research
Council, Grant/Award Number: 724324

Handling Editor: Samantha Price**Abstract**

1. Morphometrics has become an indispensable component of the statistical analysis of size and shape variation in biological structures. Morphometric data have traditionally been gathered through low-throughput manual landmark annotation, which represents a significant bottleneck for morphometric-based phenomics. Here we propose a machine-learning-based high-throughput pipeline to collect high-dimensional morphometric data in two-dimensional images of semi-rigid biological structures.
2. The proposed framework has four main strengths. First, it allows for dense phenotyping with minimal impact on specimens. Second, it presents landmarking accuracy comparable to manual annotators, when applied to standardized datasets. Third, it performs data collection at speeds several orders of magnitude higher than manual annotators. And finally, it is of general applicability (i.e. not tied to a specific study system).
3. State-of-the-art validation procedures show that the method achieves low error levels when applied to three morphometric datasets of increasing complexity, with error varying from 0.57% to 2.2% of the structure's length in the automated placement of landmarks. As a benchmark for the speed of the entire automated landmarking pipeline, our framework places 23 landmarks on 13,686 objects (zooids) detected in 1,684 pictures of fossil bryozoans in 3.12 min using a personal computer.
4. The proposed machine-learning-based phenotyping pipeline can greatly increase the scale, reproducibility and speed of data collection within biological research. To aid the use of the framework, we have developed a file conversion algorithm that can be used to leverage current morphometric datasets for automation, allowing the entire procedure, from model training all the way to prediction, to be performed in a matter of hours.

KEYWORDS

automation, biological structures, dense phenotyping, landmarks, machine learning, morphometrics, phenotyping pipeline, zooids

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2020 The Authors. *Methods in Ecology and Evolution* published by John Wiley & Sons Ltd on behalf of British Ecological Society

1 | INTRODUCTION

In the past 20 years, genomics has revolutionized our understanding of biology, leading to the discovery of a wealth of novel phenomena (Koboldt, Steinberg, Larson, Wilson, & Mardis, 2013). These discoveries were only possible through the development of a technological infrastructure that allowed us to acquire genomic information at a large scale (Schuster, 2007). While many researchers have argued that phenomics—large-scale phenotyping—will bring about a similar revolution in biology, most approaches for collecting high-throughput phenotypic data developed so far are system-specific and difficult to generalize (see Boyer et al., 2011; Falkingham, 2012; Hsiang et al., 2018; Kristensen, Parsons, Hallgrímsson, & Boyd, 2008; Manacorda & Asurmendi, 2018), with the notable exception of subcellular phenotypes (e.g. Clish, 2015). As a consequence, for most study systems, phenotyping methods remain low-throughput and can only be applied on a small scale (Houle, Govindaraju, & Omholt, 2010). These limitations have important consequences for biological research. First, the majority of the history of life remains inaccessible to genomic information, due to DNA degradation (Allentoft et al., 2012), and therefore needs to be integrated into current evolutionary theory through its phenotype. Second, even with the advances that came with genomics, we still cannot understand a wealth of important biological phenomena, such as disease and evolutionary fitness, suggesting that we have not been able to fully characterize such phenomena (Houle et al., 2010). Third, phenomic and genomic research are largely synergistic, and the end product of being able to characterize both aspects of biological variation simultaneously is likely to increase the power of each approach (Houle et al., 2010).

The low-throughput nature of current phenotyping methods is particularly problematic in morphometrics. Geometric morphometrics in particular has become an indispensable component of statistical analyses of size and shape variation in biological structures, with thousands of papers that make use of it being published every year (Zelditch, Swiderski, & Sheets, 2012). In the last two decades, with the growth of geometric morphometrics, geneticists, evolutionary biologists, ecologists and paleobiologists have accumulated dense landmark datasets, often collected from thousands of specimen images. Despite considerable progress in multivariate statistical analyses of morphometric data (Adams, Collyer, Kaliontzopoulou, & Sherratt, 2016) and attempts of incorporating it into phylogenies (Parins-Fukuchi, 2018), we are still far from a comprehensive understanding of multivariate patterns of morphometric variation. Understanding multivariate patterns of variation requires large sample sizes (Grabowski & Porto, 2017). One of the main impediments to the acquisition of large landmark datasets is the manual collection of landmark data, which is both time- and labour-intensive. Given the recent explosion of semi-landmark use in morphometrics (Watanabe, 2018), data collection time has only increased. Depending on the number of landmarks and the necessary steps to prepare a specimen for landmark data collection, this manual annotation can take months, if not years.

A promising way to collect high-throughput phenotypic data is to automate landmark data collection using computer vision techniques (e.g. Manacorda & Asurmendi, 2018). Automated landmarking has become the gold standard in human facial landmarking for both biomedicine (Porto et al., 2019) and, more notoriously, social networking websites and software developed for mobile phones (Kazemi & Sullivan, 2014). Its application in geometric morphometrics has, however, remained restricted, likely due to technical barriers, such as non-overlapping software traditions (but see Manacorda & Asurmendi, 2018; Vandaele et al., 2018). However, the explosion in machine-learning algorithms for computer vision represents an important technological leap, which lays the foundation for the development of general methods of high-throughput high-dimensional morphometrics (He, Gkioxari, Dollar, & Girshick, 2017; Kazemi & Sullivan, 2014; Voulodimos, Doulamis, Doulamis, & Protopapadakis, 2018).

Here we develop a supervised learning-based phenotyping pipeline (ML-morph) to collect high-dimensional morphometric data in two-dimensional images of semi-rigid biological structures. This pipeline is based on adapting methods currently being used in computer vision research to morphometrics and allows for dense and accurate landmarking at low cost, high speed and with minimal impact on specimens. While morphometrics has traditionally relied on specialized software (Rohlf, 2006) and R packages (Adams et al., 2016; Dryden, 2018; Olsen & Westneat, 2015; Schlager, 2017), most computer vision libraries are implemented in Python or C++ (e.g. King, 2009). We therefore also develop file conversion algorithms to leverage current morphometric datasets for training the machine-learning phenotyping pipeline. Development of a machine-learning infrastructure will enable biologists to get the necessary data to investigate and tackle important theoretical and empirical challenges within the field, and will greatly increase the scale, reproducibility and reach of biological research.

2 | MATERIALS AND METHODS

2.1 | Description of the framework

We approach the problem of automated landmark detection using a supervised learning approach. In this approach, automated landmarking is performed by a combination of object detection (Dalal & Triggs, 2005) followed by shape prediction (Kazemi & Sullivan, 2014). In other words, our models are trained to (a) predict the location of a biological structure we intend to landmark in images (object detection) and to (b) predict the shape of each detected structure (i.e. annotate landmarks). Object detection and shape prediction models are trained using a dataset of manually annotated (i.e. digitized) images (see Section 2.3). In the following sections, we present and describe in detail the training process that can be used to generate detectors and predictors, the set of images and landmarks used to explore the performance of the method, and finally the metrics employed to validate the approach and to quantify how reliable it is in comparison with system-specific tools that have been published in the literature.

2.2 | Samples

We tested the below proposed framework on three morphometric datasets (fly wings, sea basses and bryozoan colonies) with different levels of complexity. Figure 1 shows the landmarks we analysed in each of the three datasets.

2.2.1 | Low complexity image set

The drosophilid wing is a structure commonly used in morphometric studies in ecology and evolution (e.g. Houle, Mezey, Galpern, & Carter, 2003; Palaniswamy, Thacker, & Klingenberg, 2010). We consider it to be a relatively simple dataset to develop automated landmarking algorithms for, in large part because (a) the structure is translucent, creating clear contrast between structure and background, (b) the structure itself is highly conserved across species (Houle, Bolstad, Linde, & Hansen, 2017, Figure S1) and (c) landmarks are positioned in clear intersections of wing veins. Indeed, image analysis has already been automated for drosophilid wings using a B-spline model of wing shape (Houle et al., 2003), which makes this an ideal structure to investigate the performance of our machine-learning phenotyping pipeline. More specifically, we obtained a total

of 280 images of drosophilid wings from Morphobank (Table S1) and placed a total of 12 landmarks in each image, following Houle et al. (2003). Images have 632×480 pixels resolution and have been randomly selected from a larger pool of images. The reason to limit the number of images is to mimic the size of datasets that are usually used in morphometric studies. Figure 1a shows all 12 landmarks collected in each specimen.

2.2.2 | Intermediate complexity image set

We analysed 180 lateral images of sea basses belonging to 42 species of the genus *Pseudanthias* (Randall, 1997; Table S2). This genus was chosen because it is composed by a large array of species with diverse morphologies, not only in terms of body shape but also in terms of colouration. The images have been captured by Dr. John E. Randall and are deposited at the Bernice Pauahi Bishop Museum (Honolulu, Hawaii). They represent an intermediate level of complexity for the pipeline because the morphological (and shape) diversity is much larger among these organisms than what is observed in the drosophilid dataset (Figure S1). Background colour and image resolution are also more varied across image files compared to the drosophilids wings. Following the approach used for drosophilids, a

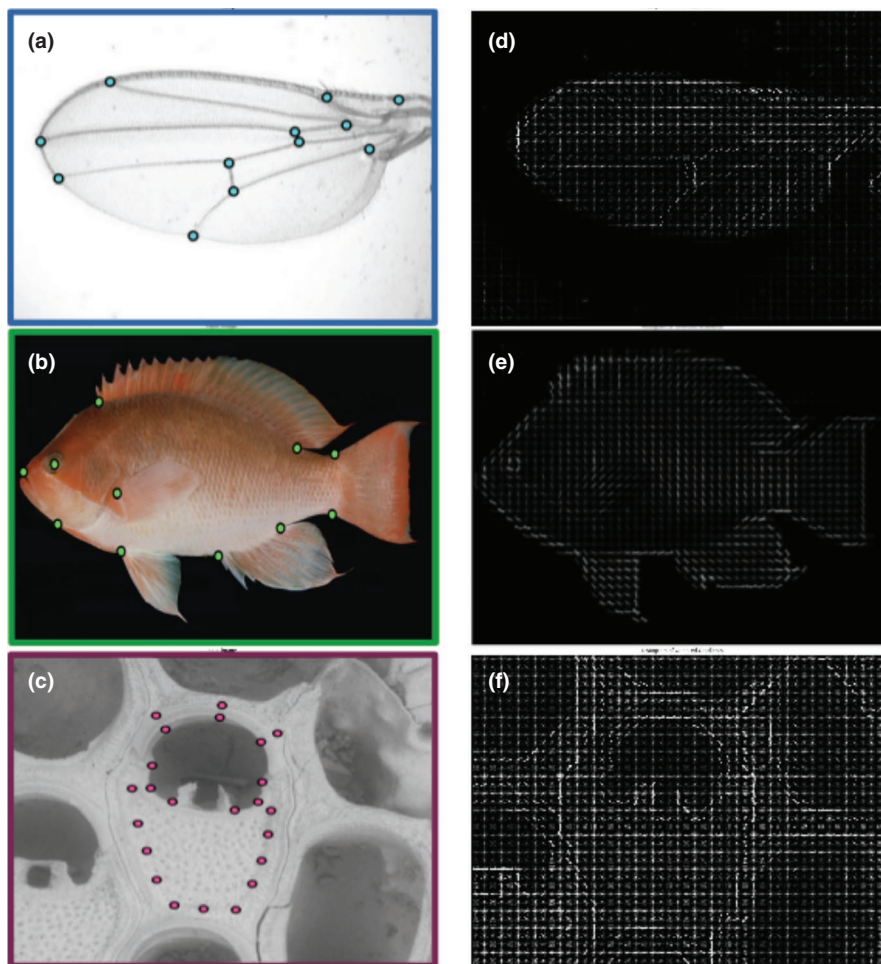


FIGURE 1 Landmark datasets and corresponding histogram of gradient (HOG) features. Three landmark datasets of increased difficulty were used to test the performance of the automated landmarking framework. Our framework uses HOG features (second column) to detect objects within an image (first column). (a) Drosophilid wing, (b) Sea bass (genus *Pseudanthias*), (c) Bryozoan colony (*Steginoporella magnifica*) and (d–f) Histogram of gradient features of images (a–c)

total of 12 landmarks were manually annotated in each fish by the same expert (Figure 1b).

2.2.3 | High complexity image set

Our high complexity dataset was collected in-house (UiO) and consists of 400 scanning electron microscope (TM-4000; Hitachi) images of bryozoan colonies belonging to the species *Steginoporella magnifica* (Table S3). These data have been acquired as part of a fossil time-series project and is composed exclusively of Plio-Pleistocene fossil specimens collected from the Wanganui basin, New Zealand (Carter & Naish, 1998; Liow et al., 2017). We consider these pictures more complex than the fish and wing pictures, as the material contains colonies with different levels of fossil degradation. In addition, colonies contain multiple zooids with extreme plasticity in morphological traits (Figure S1), requiring the identification of individual zooids from whole-colony image data, while the fish and fly wing pictures contained only one specimen per image. We imaged all *S. magnifica* colonies at $1,280 \times 960$ resolution and annotated individual zooids for a total of 23 landmarks (Figure 1c). Broken, rotated and/or partially missing zooids were marked as 'ignore' to prevent them from being used to train the object detectors and shape predictors. Each colony was imaged twice, on average, depending on the degree of degradation of each specimen. Different images of the same colony do not overlap.

2.3 | Image annotation: Defining ground truths

An essential component of supervised learning is image annotation (i.e. digitizing). In our case, prior to model fitting, we annotated all images belonging to the three datasets using both bounding box annotations and individual landmark XY coordinates using *imglab* (King, 2009). Bounding box annotations are used to locate objects within images to train object detectors, while XY coordinates are used as landmark positional ground truths by the shape prediction part of the pipeline. A total of 10% of all images from each of the three datasets were annotated twice (3 months apart), allowing us to estimate intra-observer measurement error. The intra-observer error is a key parameter, as it limits the maximum accuracy of any machine-learning algorithm. Here, we use intra-observer error as the theoretical minimum error our pipeline can achieve.

Once annotated, we randomly split all three datasets into training and validation sets, representing 80% and 20% of the images, respectively. We used this initial training and validation sets to explore the parameter space of the object detectors and shape predictors. Once the initial exploration of parameter space was completed and the parameters that resulted in the best model performance had been detected, we used a 10-fold cross validation design to evaluate the performance of these algorithms at those parameter values.

We employ metrics that are standard in problems of such kind to evaluate the performance of our models (see Sections 2.6 and 2.7 for details).

2.4 | Training object detectors

The first step in any computer vision algorithm for automated land-marking is the detection of the presence and position of the structure of the interest within the image. For example, if we want to identify landmarks in still images of drosophilid wings, we need the ability to identify the number and position of wings within an image. While a large number of object detection algorithms have been published in the past few years (Girshick, 2015; King, 2009; Ren, He, Girshick, & Sun, 2015; Viola & Jones, 2001 to name a notable few), the basic procedure for object detection remains largely the same. First, a set of positive (object, e.g. wing) and negative (not object, e.g. background) image windows are produced based on annotated training images and then a binary classifier is trained based on these windows. Finally, this classifier is tested on images in which the main model has not been trained on.

While the current state of the art in object detection relies on convolutional neural networks (CNN, e.g. Ren et al., 2015), these can be overpowered for standard biological applications. They also require large training samples (one to two orders of magnitude higher than the one proposed here), more fine-tuning of training parameters, specialized hardware (e.g. graphics processing unit, GPU), and usually perform at lower image processing rates (images per second) than simpler models (Suleiman, Chen, Emer, & Sze, 2017).

Our implementation of object detection is based on histogram of oriented gradients (HOG) features (Figure 1d–f) within a sliding window framework (Dalal & Triggs, 2005; Figure 2). In brief, based on each training image, we randomly extract image regions that contain and do not contain the object of interest using bounding box annotations. From each randomly extracted region, we extract HOG features and create a training set. We then train a structural support vector machine (SVM) classifier to classify images according to these two labels (object vs. no object). To perform detection on the test set, we scan this classifier over an image pyramid using a sliding window, and, whenever a certain window passes a threshold test, it is output as being an object, after non-maximum suppression is performed. HOG feature generation follows the general approach developed by Dalal and Triggs (2005). We start by dividing each image into 8×8 cells and compute the gradient vector at each pixel. Following this step, we end up with a $64 (8 \times 8)$ gradient vectors for each cell, which are then represented as histograms. These histograms compress the information in each cell by splitting that information into angular bins, where each bin corresponds to a gradient direction (20 degrees each). We then normalize the gradients using block normalization.

During training, two parameters of SVMs require particular attention: the soft margin parameter C and the insensitivity

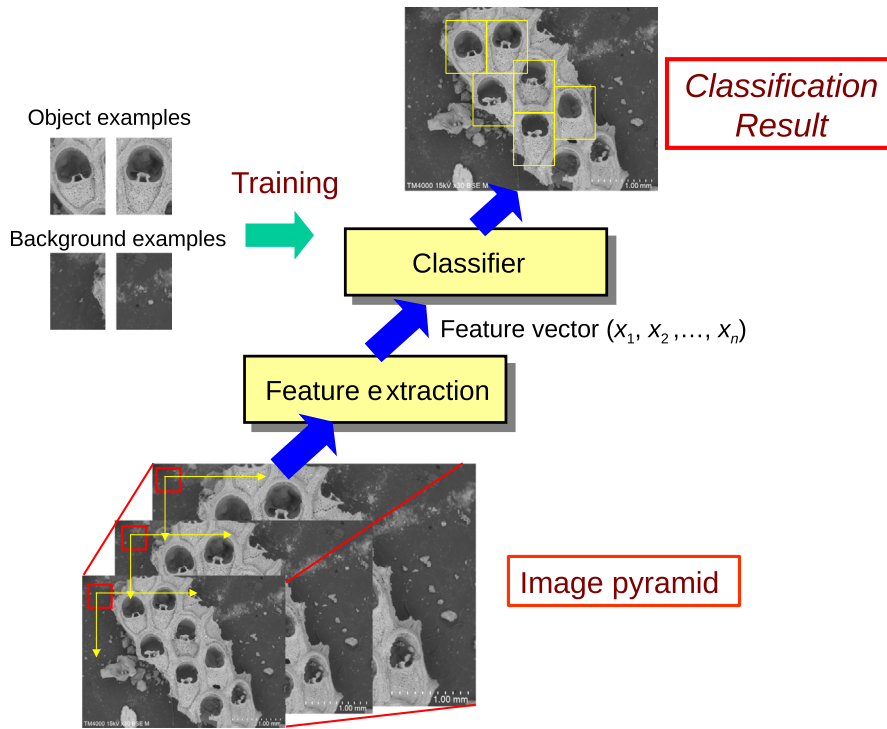


FIGURE 2 Diagram of object detection framework. In the proposed framework, a training set containing object and background examples is used to train a support vector machine classifier. A sliding window is then scanned over an image pyramid and its features are then classified as either object or background. Detected objects are output by the model after non-maximum suppression of overlapping windows is performed

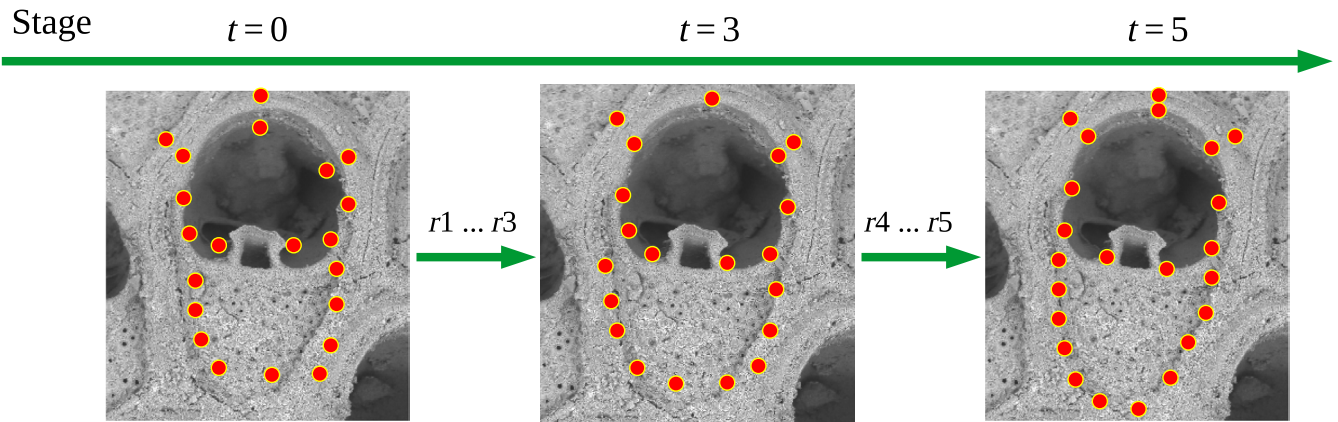


FIGURE 3 Diagram of shape prediction framework. In the proposed framework, a training set containing positional ground truths is used to train a cascade shape regression model. Shape prediction of unannotated objects is then performed using a sparse subset of pixel intensities collected on the area of the image where each object was identified. Cascade shape regression uses an iterative procedure (green arrows) to predict the shape of the object (orange landmarks). Starting with an initial shape 'guess' (landmarks at $t=0$), the predicted shape is refined through shape increments (using r regressors) in a cascade (from $t=0$ to $t=5$)

zone (ϵ ; Kecman, 2001, pp. 182–183). The C parameter regulates the size of the margin, which is the distance between the hyper-plane that separates the two classes and the closest data point. The ϵ parameter, on the other hand, regulates the penalty associated with errors in classification. To evaluate the impact various C and ϵ parameters (the training parameters) have on the performance of the final model, we performed an exhaustive grid search, in which a model is trained for each possible combination of hyperparameters, given a certain parameter range. In our case, we varied both the SVM C (from 1 to 7) and ϵ parameters (from 10^{-2} to 10^{-4}). The combination of hyperparameters that results in the best model performance is then used for the

10-fold cross validation, preventing errors that are incurred by over-fitting.

2.5 | Training shape predictors

We performed shape prediction by feeding objects detected using the above procedure to a cascade regression algorithm (Kazemi & Sullivan, 2014) that performs object-specific automated landmark detection.

In brief, cascade shape regression predicts shape (S) in an iterative procedure using a sparse subset of pixel intensities collected on the

area of the image where the object was identified. Starting with an initial shape 'guess' (S^0), the predicted shape (S) is refined through shape increments in a cascade of depth K (Figure 3) so that

$$\hat{S}_i^{(t+1)} = \hat{S}_i^{(t)} + r_t(I_{\pi_i}, \hat{S}_i^{(t)}),$$

where r_t represents a trained regressor, I_{π_i} represents an image and $\hat{S}_i^{(t)}$ represents the previous stage's shape estimate (Kazemi & Sullivan, 2014). Each regressor in the cascade is learnt using a gradient boosting algorithm with a square error loss function (Hastie, Tibshirani, & Friedman, 2009). The gradient boosting algorithm uses a user-defined number of regression trees of depth T for which decisions at each node are based on thresholding the difference in intensity values at a random pair of pixels, given an exponential prior over the distance between pixels used in a split. More details of the gradient boosting algorithm can be found in Kazemi and Sullivan (2014). During training, we performed data augmentation of the training set by adding in 300 random deformations of each training object, effectively boosting the number of training examples. Similar to object detection, we carry out an exhaustive grid search of the training parameters T (from 1 to 8) and K (from 10 to 30), and use the best performing model to define the training parameters during the 10-fold cross validation.

2.6 | Testing object detectors: Precision and recall

We evaluated the performance of object detectors using a 10-fold cross validation design. We used three metrics that are standard in the field: precision, recall and mean average precision (Powers, 2011). Precision is defined here as the ratio of true positive object predictions to all positive object predictions. Recall, also known as the true positive rate, refers to the ratio of true positive object predictions to all true positive objects present in the data. Finally, mean average precision is defined as the area under the recall-precision curve. All three performance measurements vary from 0 to 1 and can be directly compared across datasets. Note that the performance of the classifier on the test sets will only correctly reflect the performance of the model on unlabelled data if these pictures were taken with the same general setup as in the training data (i.e. similar background and resolution), and not necessarily if image capture methods change dramatically (e.g. extremely low resolution images, such as 80×80).

To allow for a more in-depth evaluation of the performance of object detectors trained above, we also report the confusion matrix of this binary classification system as well as the receiver operating characteristic (ROC) curve of the k -fold with worst performance. The ROC curve was generated following King (2009).

2.7 | Testing shape predictors: Euclidean distance

To quantify the accuracy in placement of landmarks in each k -fold, we measured the normalized Euclidean distance in pixels between each landmark's location in the ground truth (test) set and as

predicted by the model (Kazemi & Sullivan, 2014). The normalization process allows direct comparison of the error between the three datasets (and across studies), since image parameters are different in each study. Here, landmark distances were normalized by the total length of the structure. For the most complex dataset, we also break down the total measurement error by image file and landmark, allowing us to evaluate potential sources of error in landmark predictions. Finally, in the Supplemental Text 1, we analyse the effects of prediction errors in possible downstream morphometric analyses.

2.8 | Implementation and file conversion

All algorithms were implemented in Python using the following libraries: numpy 1.13.3, pandas 0.22.0, dlib 19.7.0 and opencv-python 3.4.0.12. Our current implementation can be found at <https://github.com/agporto/ml-morph>. The scripts can be used out-of-the-box and can be run on any operating system. On the github page, we also provide a detailed vignette and an example dataset to illustrate the use of this software. Among the capabilities of this software, we should note that there is a preprocessing script that converts traditional landmark files from the most common morphometric packages (tps format; Rohlf, 2006) to standard input files used in training and testing of the object detectors and shape predictors used in our machine-learning phenotyping pipeline. This script allows previously landmarked datasets to be immediately used in automation.

3 | RESULTS

3.1 | Object detection

The object detectors learnt based on HOG features achieved a high degree of recall and precision in all three datasets (Table 1). To a large extent, the training parameters had no major impact on the performance of the final model, with several different training parameter combinations resulting in similar performance (Figure 4). As expected, recall and precision were highest in the low and intermediate complexity datasets (100%), and slightly lower in the high complexity dataset (~96%). Using the high complexity dataset as a benchmark, we can infer that false positives occur at low frequency (~4%; Table 1) and that only a small amount of true positives are missed by the pipeline (~5%; Table 1). It is worth noting, however, that other properties of objects can be used to further improve these results. For example, false positives can be further eliminated using size filters or through multivariate outlier analyses. Closer inspection of the confusion matrix and the ROC curve of the worst performing k -fold model of the most complex dataset (Figure S2) reveal that increasing the precision (i.e. removing false positives) based on detector scores alone would require driving the recall rate (i.e. true positive rate) close to 75%. In other words, one would have to be willing to ignore 25% of all true objects to make the detector output only real objects.

TABLE 1 Training and testing parameters of all k -fold object detection models. In this table, we detail the training parameters of all k -fold object detection models for each of the landmark datasets used in this manuscript, together with their performance on the test sets. Precision is defined here as the ratio of true positives to all positives. Recall refers to the ratio of true positives to all true positives. Mean average precision is defined as the area under the recall–precision curve. Details of the training parameters can be found in the main text and in Powers (2011)

Datasets	N	Training		Testing (10-fold design)								
				Precision			Recall			Average precision		
		C	ϵ	M	Minimum	Maximum	M	Minimum	Maximum	M	Minimum	Maximum
Low	280	1	10^{-4}	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Intermediate	180	3	10^{-2}	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
High	400	5	10^{-2}	0.96	0.95	0.99	0.95	0.92	0.97	0.94	0.92	0.97

Abbreviations: C, soft margin parameter; epsilon, penalty parameter; MAP, mean average precision; N, sample size.

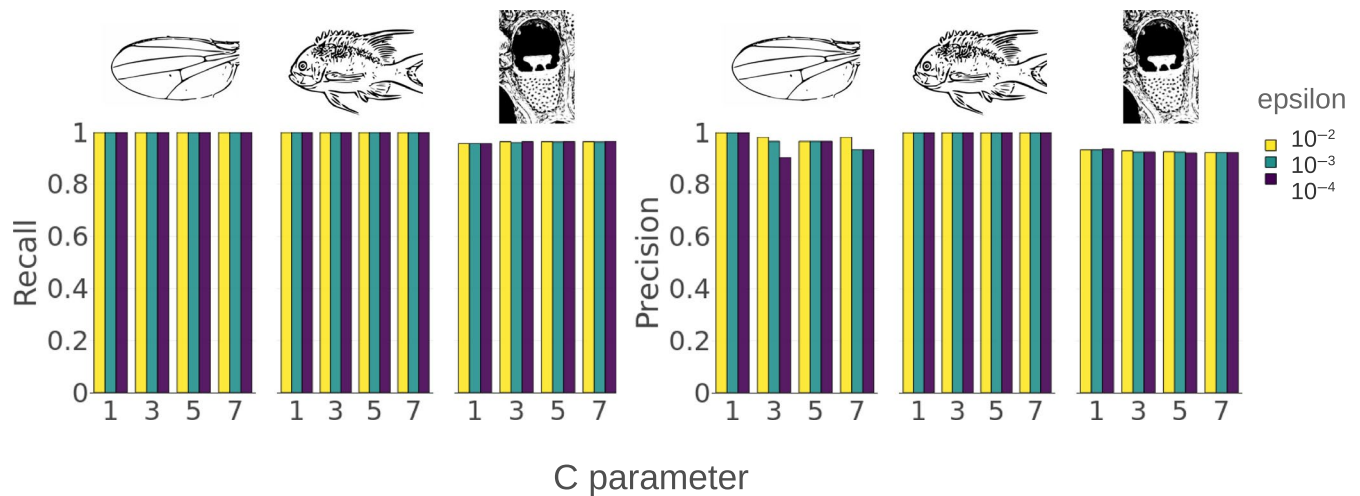


FIGURE 4 Grid search for optimal object detection training parameters. Barplots illustrating the recall and precision obtained by object detection models trained with varying training parameters. Each dataset is represented by a drawing of an example specimen. In our case, we varied the soft margin parameter C (X-axis, from 1 to 7) and the insensitivity zone (coloured bars, from 10^{-2} to 10^{-4}). The C parameter regulates the distance between the hyperplane that separates the two classes and the closest data point. The insensitivity zone (ϵ) regulates the penalty associated with errors in classification

3.2 | Shape prediction

We report shape prediction results for the three datasets in Table 2. The table contains comparisons of the average, minimum and maximum observed error in the 10-fold cross-validation design against the theoretical minimum (intra-observer error) and against the best performing semi-automatic methods available in the literature. All errors are reported as normalized mean Euclidean distances. The magnitude of error is low across all datasets. In the low complexity dataset, the average k -fold model presents a normalized error of 0.57% of the wing length (3 pixels in raw values). This magnitude of measurement error is remarkably close to the theoretical maximum accuracy that is possible given intra-observer measurement error (0.33% or 1.8 pixels). The average k -fold model performs equivalently to the best semi-automatic approach available in the literature (0.67% or 3.54 pixels; Loh, Ogawa, Kawana, Tamura, & Lee, 2017, table 5) and considerably better than others (e.g. >2% or >10.6 pixels for Houle et al., 2003). Search of the training parameter space

reveals that similar results can be obtained for different tree depths and cascade depths (Figure 5). In other words, prediction performance is generally high, regardless of model fine-tuning. Relatively, small models (low depths) still provide robust results. Over-fitting can only be observed for tree depths larger than six.

While the degree of error observed for the intermediate complexity dataset is larger compared to the low complexity dataset (0.87% or 5.14 pixels, Table 2), it is still similar to the intra-observer error (0.54% or 3.2 pixels). Search of the training parameter space reveals that greater tree depths lead to over-fitting and pronounced performance loss (Figure 5). The depth of the cascade has only a mild impact on the results.

Finally, the degree of error observed for the high complexity dataset is higher than for the other two datasets (2.2% or 6.5 pixels, Table 2), in a large part because intra-observer error is also higher (1.3% or 3.8 pixels). Notably, the ratio of the prediction error relative to intra-observer error is remarkably similar across all datasets, with the shape predictors producing datasets with a degree of error

TABLE 2 Training and testing parameters of all k -fold shape prediction models. In this table, we detail the training parameters of all k -fold shape prediction models for each of the landmark datasets used in this manuscript, together with their performance on the test set. Error is measured here as the normalized Euclidean distance. The literature tab refers to the performance of the best semi-automatic or fully automatic model present in the literature for each study system. Given that drosophilids are the only structure that was automated in the past, only results for drosophila algorithm are detailed. Details of the training parameters can be found in the main text and in Dalal and Triggs (2005) and King (2009)

									Testing (10-fold design)				Literature
									Error (normalized Euclidean distance)				
Datasets	Training								M	Minimum	Maximum	Theoretical minimum	Error
	N	nu	T	K	Trees (N)	Feature pool size	Test splits	Oversampling amount (N)					
Low	280	0.1	3	30	500	1,200	20	300	0.57%	0.47%	0.79%	0.330%	0.6% (Loh et al., 2017) >2% (Houle et al., 2003) ^a
Intermediate	180	0.1	2	30	500	1,200	20	300	0.87%	0.68%	1.04%	0.540%	None
High	400	0.1	2	25	500	1,200	20	300	2.22%	1.86%	2.45%	1.310%	None

Abbreviations: N , sample size; nu , regularization parameter; T , tree depth; K , cascade depth; Theoretical minimum, intra-observer error.

^aAs reported in Loh et al. (2017).

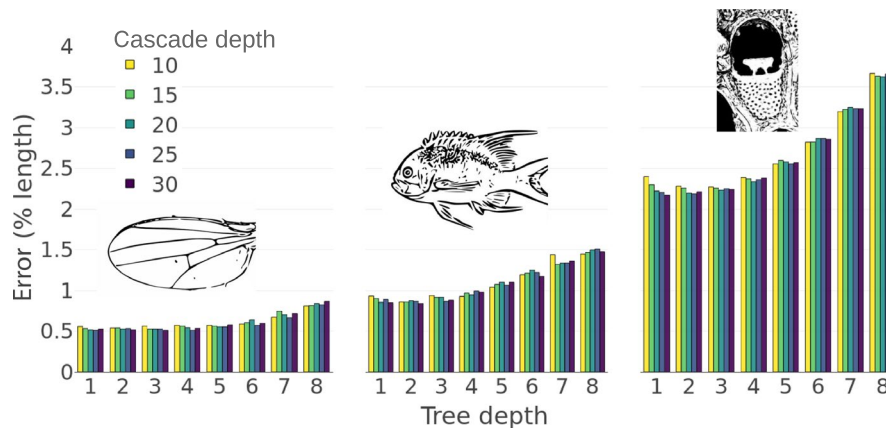


FIGURE 5 Grid search for optimal shape prediction training parameters. The barplot illustrates the normalized mean error obtained by shape prediction models trained with varying training parameters. Each dataset is represented by a drawing of an example specimen. In our case, we varied the regression tree depth (X-axis, from 1 to 8) and the cascade depth (coloured bars, from 10 to 30). Tree depth refers to the depth of the decision trees used by the gradient boosting algorithm. The cascade depth refers to the number of shape updates that are used to predict each object's shape, given local pixel patterns

65% higher than the theoretical minimum. While 65% might seem like a substantial difference, note that in the drosophilid dataset this is equivalent to a difference of approximately one pixel.

An examination of the sources of error in each shape prediction algorithm reveals that the degree of shape variation—measured in terms of Procrustes distances—is either not or only marginally associated with prediction error in all three datasets ($r^2 = .12$, 4×10^{-4} and 8×10^{-5} for the low, medium and high complexity datasets, respectively; Figure S3). In other words, larger prediction errors are not generally associated with particularly divergent shapes, suggesting that other sources of error are more important in determining the algorithm's accuracy.

Representation in the training set is also not found to have an effect: in the medium and highly complex datasets, species or colonies lacking (by chance) representation on the training sets are predicted

with virtually identical error rates compared to species and colonies that are represented by multiple specimens in the training set, indicating that these models are learning generalizable features, rather than over-fitting to specific morphologies (Figure S4).

We have revealed two important sources of error for the complex dataset in our approach. First, mean average error varies from 1.1% to 3.5% across image files (Figure 6). Differences among image files are associated with differences in the quality of preservation of fossil specimens and with the quality of image capture. Figure 6 illustrates one of the best and worst performing images, which vary considerably in the level of contrast and taphonomy. Similarly, accuracy in landmark predictions varies significantly across landmarks, with certain landmarks presenting mean average error around 0.91%, while others present values of 4.2% (Figure 7). The presence of clear edges and high contrast is clearly associated with accuracy in the predictions (Figure 7).

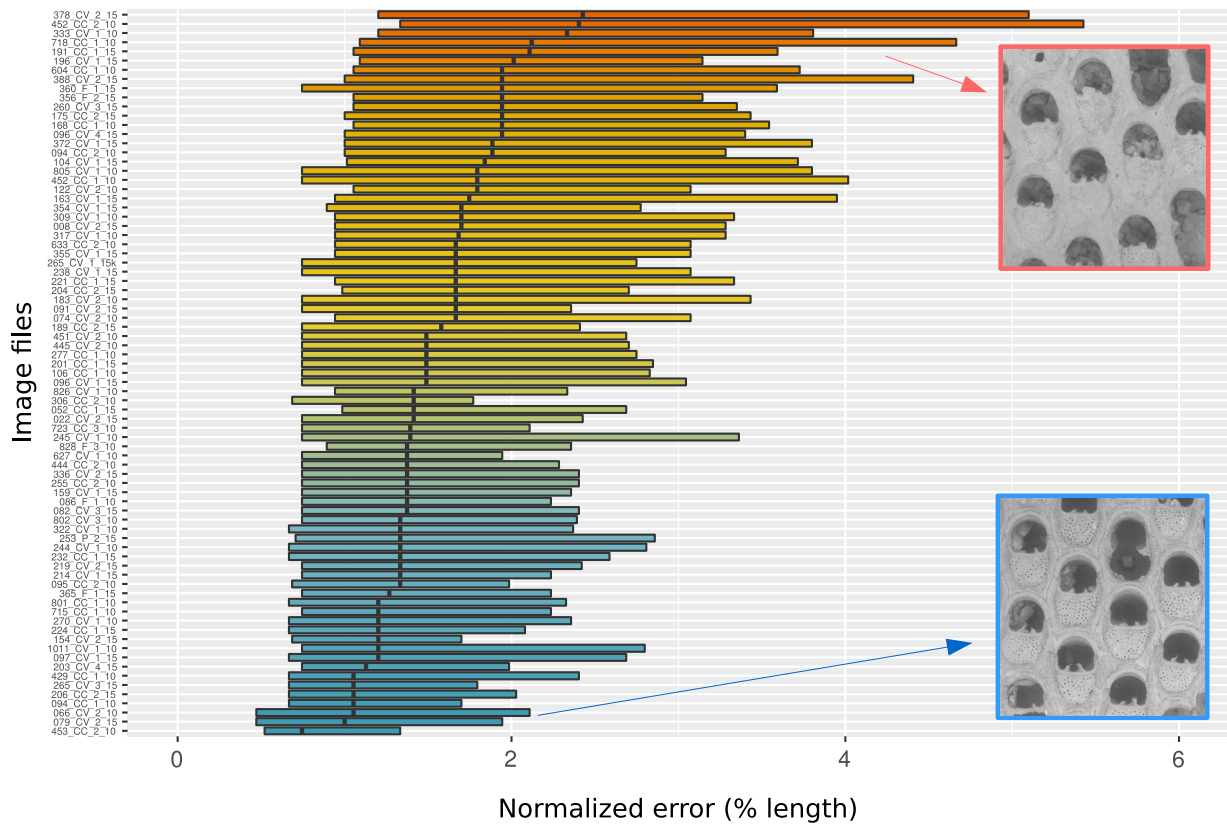


FIGURE 6 Breaking down the prediction error according to image quality. Each boxplot illustrates the distribution of landmark prediction errors of our best fit model on each image of the test set (whiskers omitted for clarity). Image files are sorted according to their median prediction error, from low (blue) to high (orange). Representative specimens from the extremes of the distribution are shown near the upper and lower margins. Note the difference in fossil degradation and image contrast between the two specimens. Note also that images with poorer image quality (orange bars) have more variance in prediction

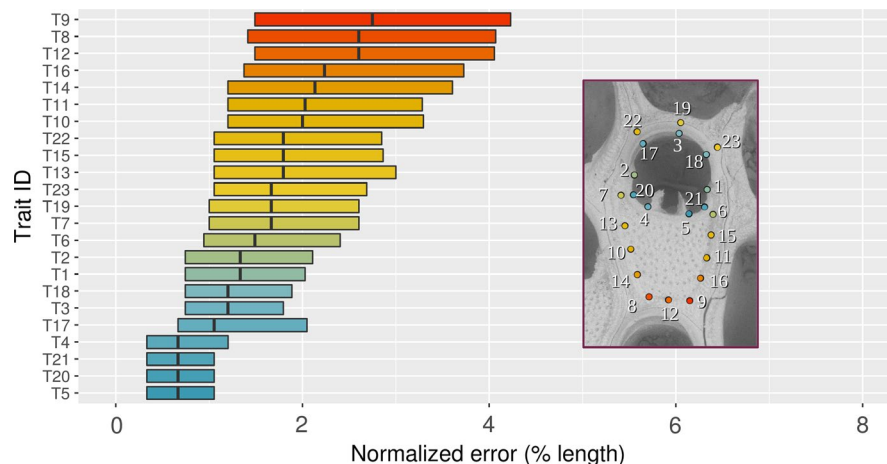


FIGURE 7 Breaking down the prediction error among the landmarks. Each boxplot illustrates the distribution of landmark prediction errors of our best fit model according to each landmark (whiskers omitted for clarity). Landmarks are sorted according to their median prediction error, from low (blue) to high (orange). We also show one representative specimen with landmarks coloured according to the degree of prediction error. Note the difference in the degree of error between landmarks in low (e.g. T8, T9) and high (e.g. T4, T5) contrast areas. Note also that the variance in the distribution of errors per landmark is correlated with the median

3.3 | Implementation speed

We used the high complexity dataset to benchmark the speed of the pipeline, when run on an Intel Core i7 2.7 Ghz with 16 GB of RAM.

Using the training parameters of the best fit models (Tables 1 and 2), object detector training occurred in 9.3 min per run, and testing took a total of 15 s per run. Shape predictor training occurred in 9.23 hr per run, while testing took only 7 s per run. To a large extent, the amount

of oversampling applied to the training set is the key parameter regulating the training times for shape predictors, with higher oversampling amount requiring longer training.

To benchmark the speed at which the pipeline can generate predictions for new images, we applied the final object detector and shape predictor to a larger sample of *Steginoporella* images ($N = 1,684$). In this larger sample, the entire automated landmarking pipeline took 3.12 min to process all images, during which it identified 23 landmarks per zooid in a total of 13,686 zooids at a resolution of 800×600 pixels. Based on generous assumptions (3 min per zooid), an experienced morphometrician would spend 684 hr (about eighty-five 8-hr uninterrupted workdays) to collect a comparable dataset.

4 | DISCUSSION

Morphometric characterization of biological structures has become an essential component in studies of morphology. However, morphometric data collection has remained mostly manual, in large part due to the lack of a general framework for automation (see Hsiang et al., 2018 for a notable exception). In this study, we propose a simple, fast and accurate pipeline for automation of landmark collection in any semi-rigid biological structure. Our approach is based on supervised learning and uses a combination of object detection (Dalal & Triggs, 2005) and shape prediction (Kazemi & Sullivan, 2014) to accurately place landmarks of interest on one or several objects in an image.

4.1 | Model performance

Object detectors performed well in all three datasets, making no mistakes in the low and intermediate complexity datasets (Table 1, Figure 4). Even in the most complex dataset consisting of pictures of fossil specimens in various taphonomic states, object detection rates were high (~96%), indicating detection should not be a concern for the standardized images datasets typically analysed in biological studies. Note, however, that the high performance of object detectors in the low and intermediate complexity datasets is due to the high standardization of image capture method and should not be used as the null expectation for less standardized data.

As expected, the shape prediction part of the pipeline is where differences in accuracy among datasets were the largest (Table 2, Figure 5). In drosophilid wings, accuracy is loosely associated with the degree of shape disparity of each specimen ($r^2 = 0.12$; Figure S3) and very similar to the scale typically found among human observers. In *Steginoporella*, heterogeneous levels of accuracy depend instead on the landmark and the individual specimen (Figures 6 and 7). We highlight two main sources for the heterogeneity in accuracy in the *Steginoporella* dataset. First, image quality varied considerably across specimens, largely due to fossil degradation. Since the *Steginoporella* samples are composed exclusively of fossil specimens (see Liow et al., 2017), some of which were collected at harsher depositional

environments than others, the samples vary in the degree of damage and cementation, leading to a reduction in contrast between the morphological features of interest and the background. Additionally, some landmarks are placed in contrast poor locations (Figure 6), adding noise to automation based on HOG features, as this is informed by local image contrast (Dalal & Triggs, 2005).

4.2 | Advantages of the framework

Given the relative small size of our training sets and how well the pipeline performs on datasets of different complexity, we argue here that this pipeline has general applicability in biology. The main advantages of our framework are its accuracy, speed and infrastructure requirement. Even in the most complex dataset, the landmarking pipeline being proposed here exceeds the accuracy reported in many system-specific landmark approaches observed in the literature (Houle et al., 2003; Loh et al., 2017). Note also that the intra-observer error that we define as the theoretical minimum error when evaluating our method would probably be larger if multiple people had annotated such datasets, as commonly done in morphometric studies. Furthermore, potential errors committed by the pipeline can be easily corrected manually after prediction using the image software of choice (e.g. *imglab*: King, 2009; and/or *tps-Dig*: Rohlf, 2006). When comparing time-weighted performance with manual annotation, the benefits of automation are even clearer. Manual annotation of the entire *Steginoporella* dataset would take around 85 days, using generous assumptions of the work-load (3 min per zooid). In about the same amount of time (3.12 min), our pipeline is able to process 1,684 images, representing 13,686 zooids and 23 landmarks per zooid. This time can even be reduced to 1.5 min if a lower image resolution is used, though that is likely to increase error levels. Moreover, if a previously annotated dataset has already been collected, the whole procedure, from producing the training sets to training the object detection and shape predictors to obtaining an automated landmarking framework, can be performed on a personal computer in the timespan of two workdays. The proposed methodology for object detection, while not in our purview, can also be used for other purposes in biology. For example, based on the detected objects, one could adapt this framework to count objects and to obtain, for example, the spacing of objects on an image, such as the distance between the zooids of *Steginoporella*.

In our view, this automated landmarking framework opens up the possibility of development of truly high-throughput high-dimensional phenotyping procedures, which will propel biology into the age of phenomics (Houle et al., 2010).

4.3 | Limitations of the framework

The most important limitation in the proposed pipeline is that HOG-based object detectors are partially sensitive to changes in orientation (Dalal & Triggs, 2005). As a consequence, objects cannot be

detected when upside down, for example. In our view, these limitations can be easily overcome using one of the following: (a) training of multiple object detectors, one for each position or (b) through the use of CNN-based detectors, such as Ren et al. (2015). While CNN-based detectors are insensitive to changes in orientation, it is worth pointing out that such algorithms require a larger training dataset (more images), higher hardware specifications, are less portable (file size), and require more fine-tuning (Suleiman et al., 2017).

Another limitation of the proposed pipeline is that it can only extract data from 2D images. A non-trivial portion of geometric morphometrics is done in 3D, and while the techniques presented here could be expanded to 3D objects, we currently do not have an efficient implementation of it. At most, our proposed framework can be applied to 2D slices of 3D structures, but while this might prove useful in some systems (Hsiang et al., 2018), our framework cannot be considered an explicit 3D approach. Finally, although we provide python code that can be used out-of-the-box, we strongly recommend that other authors explore all training parameters for object detection and shape prediction, as those can have significant impact on the accuracy of the final model (Kazemi & Sullivan, 2014). This is especially true if image parameters (e.g. resolution) are significantly different from the three datasets presented here.

4.4 | Larger context: Alternatives to supervised learning

In recent years, we have witnessed an increase in the diversity of high-throughput high-dimensional morphometric approaches. Below, we briefly consider the advantages and disadvantages of other approaches that can be used for high-throughput morphometrics when compared to the supervised approach being proposed here. We should note, however, that these approaches are not mutually exclusive and that it is quite likely that a synergistic approach to high-throughput morphometrics would likely be the most beneficial of all (e.g. Keshavan, Yeatman, & Rokem, 2019).

One important alternative to supervised learning approaches is the use of crowd-sourcing, in which data annotation tasks are distributed to a large pool of remote workers through online servers (e.g. Chang & Alfaro, 2016). While still uncommon, crowd-sourcing approaches are starting to be effectively used in biological research (Bender, 2016; Chang & Alfaro, 2016; Willis et al., 2017). The advantages of crowd-sourcing over manual annotation are very similar to the advantages of supervised learning. In particular, crowd-sourcing allows for large amounts of morphometric data to be collected in the span of hours. Crowd-sourcing also has two additional benefits. By relying on multiple annotators, it allows researchers to control for inter-observer error. Crowd-sourcing approaches also allow for data collection at larger taxonomic scopes (Chang & Alfaro, 2016). It also has some disadvantages. Crowd-sourcing tends to be much more expensive and requires much more setup time than supervised learning approaches. The expenses are related to the payment of remote workers and for the use of the server. The setup time, for its turn, is

related to the development of training materials that effectively teach remote workers on how to digitize the landmarks, and also on the development of the annotation front-end that is used by each worker when annotating images. For certain anatomical landmarks of difficult recognition, crowd-sourcing also leads to much higher error rates when compared to expert morphologists (Chang & Alfaro, 2016).

Another alternative to supervised learning is the use of citizen science (Keshavan et al., 2019). Citizen science refers to the use of trained volunteers to collect data on the images of interest. Usually, these volunteers are considerably better trained than remote workers and, therefore, produce higher quality landmark annotations (Chang & Alfaro, 2016). However, they are also on a lower throughput range due to the manual nature of the data collection and are unlikely to scale up to the extent of supervised learning algorithms, since that would require hundreds of volunteers. Moreover, the availability of volunteers is often a limiting factor outside of large cities and/or outside of natural history museum environments. For that reason, citizen science has been most effective when combined with computer vision tools (Keshavan et al., 2019).

Finally, a third alternative to the supervised approach proposed here would be the use of unsupervised learning algorithms (e.g. Jakab, Gupta, Bilen, & Vedaldi, 2018). In this class of machine-learning algorithms, landmark positions are learnt from the data itself. This class of algorithms not only provides the possibility of an unbiased choice of landmarks but also effectively removes the need for user annotation. While past unsupervised approaches performed considerably more poorly than their supervised counterparts, recent works have reached state-of-the-art results in human face datasets (Jakab et al., 2018), indicating that there is an untapped potential for unsupervised approaches to become relevant for morphometrics. In our view, there are three main obstacles to be overcome for unsupervised approaches to have a significant role to play in morphometrics. First, unsupervised approaches require much larger samples sizes to train, a requirement that is unlikely to be met by most biological applications. Second, unsupervised approaches also make 'naïve' topological assumptions (sensu Fukushima, Funai, & Iida, 2019) about the homology of landmarks, which could lead to substantial complications when applied to biological datasets in which homology cannot be easily inferred from topology alone. Third, there is no guarantee that an unsupervised approach places landmarks that represent traits of interest for a particular research question.

5 | CONCLUSIONS

We have developed a machine-learning pipeline (ML-morph) for automated detection and landmarking of biological structures in images, which can be used to collect morphometric data at a large scale. ML-morph opens up the number of possibilities for automation within the morphometric community, greatly increasing the scale of the questions being asked and opening up new research avenues that previously faced sample size barriers.

ACKNOWLEDGEMENTS

A.P. and K.L.V. were funded by the Research Council of Norway—grant no. 249961. We wish to thank Emma Sherratt and an anonymous reviewer for thoughtful comments in previous versions of this manuscript. We also thank Mali Ramsfjell for helping with the *Steginoporella* sample preparation and imaging, and Lee Hsiang Liow for access to the scanning electron microscope (TM-4000). The scanning electron microscope (TM-4000) used in this study has been funded by The European Research Council—grant no. 724324.

AUTHORS' CONTRIBUTIONS

A.P. and K.L.V. conceived the project; A.P. collected the data and designed the methodology; A.P. analysed the data and wrote the code; all authors contributed to writing the manuscript, with A.P. having the lead on the writing.

DATA AVAILABILITY STATEMENT

The necessary image and landmark files are available on Morphobank https://morphobank.org/index.php/Projects/ProjectOverview/project_id/3598 (Porto & Voje, 2020). All necessary scripts can be found on Zenodo <https://doi.org/10.5281/zenodo.3634588> (Porto, 2020) or GitHub (<https://github.com/agporto/ml-morph>).

ORCID

Arthur Porto  <https://orcid.org/0000-0002-9210-8750>

Kjetil L. Voje  <https://orcid.org/0000-0003-2556-3080>

REFERENCES

- Adams, D. C., Collyer, M., Kaliontzopoulou, A., & Sherratt, E. (2016). *Geomorph: Software for geometric morphometric analyses*. Retrieved from <https://rune.une.edu.au/web/handle/1959.11/21330>
- Allentoft, M. E., Collins, M., Harker, D., Haile, J., Oskam, C. L., Hale, M. L., ... Bunce, M. (2012). The half-life of DNA in bone: Measuring decay kinetics in 158 dated fossils. *Proceedings of the Royal Society B: Biological Sciences*, 279(1748), 4724–4733. <https://doi.org/10.1098/rspb.2012.1745>
- Bender, E. (2016). Crowdsourced solutions. *Nature*, 533(3), S62–S64.
- Boyer, D., Lipman, Y., St Clair, E., Puente, J., Funkhouser, T., Patel, B., ... Daubechies, I. (2011). Algorithms to automatically quantify the geometric similarity of anatomical surfaces. *Proceedings of the National Academy of Sciences of the United States of America*, 108(45), 18221–18226. <https://doi.org/10.1073/pnas.1112822108>
- Carter, R. M., & Naish, T. R. (1998). A review of Wanganui Basin, New Zealand: Global reference section for shallow marine, Plio–Pleistocene (2.5–0 Ma) cyclostratigraphy. *Sedimentary Geology*, 122(1–4), 37–52. [https://doi.org/10.1016/S0037-0738\(98\)00097-9](https://doi.org/10.1016/S0037-0738(98)00097-9)
- Chang, J., & Alfaro, M. E. (2016). Crowdsourced geometric morphometrics enable rapid large-scale collection and analysis of phenotypic data. *Methods in Ecology and Evolution*, 7(4), 472–482. <https://doi.org/10.1111/2041-210X.12508>
- Clish, C. B. (2015). Metabolomics: An emerging but powerful tool for precision medicine. *Molecular Case Studies*, 1(1). <https://doi.org/10.1101/mcs.a000588>
- Dalal, N., & Triggs, B. (2005). *Histograms of oriented gradients for human detection*, June. Retrieved from <https://hal.inria.fr/inria-00548512>
- Dryden, I. L. (2018). *shapes package*. Vienna, Austria: R Foundation for Statistical Computing. Contributed package. Version 1.2.4. Retrieved from <http://www.R-project.org>
- Falkingham, P. L. (2012). Acquisition of high resolution three-dimensional models using free, open-source, photogrammetric software. *Palaeontologia Electronica*. <https://doi.org/10.26879/264>
- Fukushima, K., Funai, S. S., & Iida, H. (2019). Featuring the topology with the unsupervised machine learning. *arXiv preprint arXiv:1908.00281*.
- Girshick, R. (2015). Fast R-CNN. *ArXiv:1504.08083 [Cs]*, April. Retrieved from <http://arxiv.org/abs/1504.08083>
- Grabowski, M., & Porto, A. (2017). How many more? Sample size determination in studies of morphological integration and evolvability. *Methods in Ecology and Evolution*, 8(5), 592–603. <https://doi.org/10.1111/2041-210X.12674>
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction*, Springer series in statistics. New York, NY: Springer New York.
- He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). Mask R-CNN. In *2017 IEEE International Conference on Computer Vision (ICCV)* (pp. 2980–2988). Venice, Italy: IEEE. <https://doi.org/10.1109/ICCV.2017.322>
- Houle, D., Bolstad, G. H., van der Linde, K., & Hansen, T. F. (2017). Mutation predicts 40 million years of fly wing evolution. *Nature*, 548(7668), 447–450. <https://doi.org/10.1038/nature23473>
- Houle, D., Govindaraju, D. R., & Omholt, S. (2010). Phenomics: The next challenge. *Nature Reviews Genetics*, 11(12), 855–866. <https://doi.org/10.1038/nrg2897>
- Houle, D., Mezey, J., Galpern, P., & Carter, A. (2003). Automated measurement of drosophila wings. *BMC Evolutionary Biology*, 3(1), 25. <https://doi.org/10.1186/1471-2148-3-25>
- Hsiang, A. Y., Nelson, K., Elder, L. E., Sibert, E. C., Kahanamoku, S. S., Burke, J. E., ... Hull, P. M. (2018). AutoMorph: Accelerating morphometrics with automated 2D and 3D image processing and shape extraction. *Methods in Ecology and Evolution*, 9(3), 605–612. <https://doi.org/10.1111/2041-210X.12915>
- Jakab, T., Gupta, A., Bilén, H., & Vedaldi, A. (2018). Unsupervised learning of object landmarks through conditional image generation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in neural information processing systems* (pp. 4016–4027). Red Hook, NY: Curran Associates Inc.
- Kazemi, V., & Sullivan, J. (2014). One millisecond face alignment with an ensemble of regression trees. In *2014 IEEE conference on computer vision and pattern recognition* (pp. 1867–1874). Columbus, OH: IEEE. <https://doi.org/10.1109/CVPR.2014.241>
- Kecman, V. (2001). *Learning and soft computing: Support vector machines, neural networks, and fuzzy logic models*. Cambridge, MA: MIT Press.
- Keshavan, A., Yeatman, J. D., & Rokem, A. (2019). Combining citizen science and deep learning to amplify expertise in neuroimaging. *Frontiers in Neuroinformatics*, 13, 29. <https://doi.org/10.3389/fninf.2019.00029>
- King, D. E. (2009). Dlib-MI: A machine learning toolkit. *Journal of Machine Learning Research*, 10(Jul), 1755–1758.
- Koboldt, D. C., Steinberg, K. M., Larson, D. E., Wilson, R. K., & Mardis, E. R. (2013). The next-generation sequencing revolution and its impact on genomics. *Cell*, 155(1), 27–38. <https://doi.org/10.1016/j.cell.2013.09.006>
- Kristensen, E., Parsons, T. E., Hallgrímsson, B., & Boyd, S. K. (2008). A novel 3-D image-based morphological method for phenotypic analysis. *IEEE Transactions on Biomedical Engineering*, 55(12), 2826–2831. <https://doi.org/10.1109/TBME.2008.923106>
- Liow, L. H., Di Martino, E., Krzeminska, M., Ramsfjell, M., Rust, S., Taylor, P. D., & Voje, K. L. (2017). Relative size predicts competitive outcome through 2 million years. *Ecology Letters*, 20(8), 981–988. <https://doi.org/10.1111/ele.12795>
- Loh, S. Y., Ogawa, Y., Kawana, S., Tamura, K., & Lee, H. K. (2017). Semi-automated quantitative drosophila wings measurements. *BMC Bioinformatics*, 18(1), 319. <https://doi.org/10.1186/s12859-017-1720-y>
- Manacorda, C. A., & Asurmendi, S. (2018). Arabidopsis phenotyping through geometric morphometrics. *GigaScience*, 7(7). <https://doi.org/10.1093/gigascience/giy073>
- Olsen, A. M., & Westneat, M. W. (2015). StereoMorph: An R package for the collection of 3D landmarks and curves using a stereo camera

- set-up. *Methods in Ecology and Evolution*, 6, 351–356. <https://doi.org/10.1111/2041-210X.12326>
- Palaniswamy, S., Thacker, N. A., & Klingenberg, C. P. (2010). Automatic identification of landmarks in digital images. *IET Computer Vision*, 4(4), 247–260. <https://doi.org/10.1049/iet-cvi.2009.0014>
- Parins-Fukuchi, C. (2018). Bayesian placement of fossils on phylogenies using quantitative morphometric data. *Evolution*, 72(9), 1801–1814. <https://doi.org/10.1111/evo.13516>
- Porto, A. (2020). Data from: agporto/ml-morph: ml-morph v.1.0.0 (Version v.1.0.0). *Zenodo*, <https://doi.org/10.5281/zenodo.3634588>
- Porto, L. F., Lima, L. N., Flores, M. R., Valsecchi, A., Ibanez, O., Palhares, C. E., & de Barros Vidal F. (2019). Automatic cephalometric landmarks detection on frontal faces: An approach based on supervised learning techniques. *ArXiv:1904.10816 [Cs]*, April. <http://arxiv.org/abs/1904.10816>
- Porto, A., & Voje, K. L. (2020). Data from: ML-morph: A fast, accurate and general approach for automated detection and landmarking of biological structures in images. *Morphobank*, https://morphobank.org/index.php/Projects/ProjectOverview/project_id/3598
- Powers, D. M. W. (2011). Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and correlation. *Journal of Machine Learning Technologies*, 2(1), 37–63.
- Randall, J. E. (1997). Randall's tank photos. Collection of 10,000 large-format photos (slides) of dead fishes. Unpublished.
- Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards real-time object detection with region proposal networks. *ArXiv:1506.01497 [Cs]*, June. <http://arxiv.org/abs/1506.01497>
- Rohlf, F. J. (2006). *TpsDig, version 2.10*. Retrieved from <http://Life.Bio.Sunysb.Edu/Morph/Index.Html>. <https://ci.nii.ac.jp/naid/10022020610/>
- Schlager, S. (2017). Morpho and Rvcg – Shape analysis in R. In G. Zheng, S. Li, & G. Székely (Eds.), *Statistical shape and deformation analysis* (pp. 217–256). Cambridge, MA: Academic Press. ISBN 9780128104934.
- Schuster, S. C. (2007). Next-generation sequencing transforms today's biology. *Nature Methods*, 5(1), 16–18. <https://doi.org/10.1038/nmeth1156>
- Suleiman, A., Chen, Y. H., Emer, J., & Sze, V. (2017). Towards closing the energy gap between HOG and CNN features for embedded vision. *ArXiv:1703.05853 [Cs]*, March. <http://arxiv.org/abs/1703.05853>
- Vandaele, R., Aceto, J., Muller, M., Peronnet, F., Debat, V., Wang, C. W., ... Marée, R. (2018). Landmark detection in 2D bioimages for geometric morphometrics: A multi-resolution tree-based approach. *Scientific Reports*, 8(1), 538. <https://doi.org/10.1038/s41598-017-18993-5>
- Viola, P., & Jones, M. (2001). Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE Computer Society conference on computer vision and pattern recognition. CVPR 2001* (Vol. 1, pp. I-511–I-518). Kauai, HI: IEEE Computer Society. <https://doi.org/10.1109/CVPR.2001.990517>
- Voulodimos, A., Doulamis, N., Doulamis, A., & Protopapadakis, E. (2018). Deep learning for computer vision: A brief review. *Computational Intelligence and Neuroscience*, 2018, 13. <https://doi.org/10.1155/2018/7068349>
- Watanabe, A. (2018). How many landmarks are enough to characterize shape and size variation? *PLoS ONE*, 13(6), 1–17. <https://doi.org/10.1371/journal.pone.0198341>
- Willis, C. G., Law, E., Williams, A. C., Franzone, B. F., Bernardos, R., Bruno, L., ... Davis, C. C. (2017). CrowdCurio: An online crowdsourcing platform to facilitate climate change studies using herbarium specimens. *New Phytologist*, 215(1), 479–488.
- Zelditch, M. L., Swiderski, D. L., & Sheets, H. D. (2012). *Geometric morphometrics for biologists: A primer*. Cambridge, MA: Academic Press.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

How to cite this article: Porto A, Voje KL. ML-morph: A fast, accurate and general approach for automated detection and landmarking of biological structures in images. *Methods Ecol Evol*. 2020;11:500–512. <https://doi.org/10.1111/2041-210X.13373>