

# 1 Annotation-free Learning of Plankton for 2 Classification and Anomaly Detection

3 Vito P. Pastore<sup>1,2</sup>, Thomas G. Zimmerman<sup>1,2</sup>, Sujoy Biswas<sup>1,2</sup>, and Simone Bianco<sup>1,2\*</sup>

4 <sup>1</sup>Industrial and Applied Genomics, S2S - Science to Solution, IBM Research – Almaden, San  
5 Jose, CA USA.

6 <sup>2</sup>NSF Center for Cellular Construction, University of California San Francisco, San Francisco,  
7 CA USA.

8

9 \*To whom correspondence should be addressed

10

## 11 **Abstract**

12 The acquisition of increasingly large plankton digital image datasets requires automatic methods  
13 of recognition and classification. As data size and collection speed increases, manual annotation  
14 and database representation are often bottlenecks for utilization of machine learning algorithms  
15 for taxonomic classification of plankton species in field studies. In this paper we present a novel  
16 set of algorithms to perform accurate detection and classification of plankton species with minimal  
17 supervision. Our algorithms approach the performance of existing supervised machine learning  
18 algorithms when tested on a plankton dataset generated from a custom-built lensless digital device.  
19 Similar results are obtained on a larger image dataset obtained from the Woods Hole  
20 Oceanographic Institution. Our algorithms are designed to provide a new way to monitor the  
21 environment with a class of rapid online intelligent detectors.

22

23

## 24 **Author Summary**

25 Plankton are at the bottom of the aquatic food chain and marine phytoplankton are estimated to be  
26 responsible for over 50% of all global primary production [1] and play a fundamental role in  
27 climate regulation. Thus, changes in plankton ecology may have a profound impact on global  
28 climate, as well as deep social and economic consequences. It seems therefore paramount to collect  
29 and analyze real time plankton data to understand the relationship between the health of plankton  
30 and the health of the environment they live in. In this paper, we present a novel set of algorithms  
31 to perform accurate detection and classification of plankton species with minimal supervision. The  
32 proposed pipeline is designed to provide a new way to monitor the environment with a class of  
33 rapid online intelligent detectors.

## 34 **Introduction**

35 Plankton are a class of aquatic microorganisms, composed of both drifters and swimmers, which  
36 can vary significantly in size, morphology and behavior. The exact number of plankton species is  
37 not known, but an estimation of oceanic plankton puts the number between 3444 and 4375 [2].  
38 Traditionally, plankton are surveyed using either satellite remote sensing, where leftover biomass  
39 is inferred indirectly through measurement of total chlorophyll concentration, or with large net  
40 tows via oceanic vessels [3], with subsequent microscopic analysis of the preserved samples.  
41 Satellite imaging methods are extremely accurate in terms of global geographic association and  
42 very useful for broad species characterization but may present practical challenges in terms of  
43 accuracy of the performed counts, species preservation and fine-grained characterization. The  
44 analysis of preserved samples, instead, allows for fine grained classification and accurate counting  
45 with narrow spatial sampling. More recently, real time observation of plankton species has been

46 made possible by novel instruments for high-throughput *in situ* autonomous and semi-autonomous  
47 microscopy [4]. Such high-resolution imaging instruments make it possible to observe and study  
48 spatio-temporal changes in plankton morphology and behavior, which can be correlated with  
49 environmental perturbations. Sudden or unexpected changes in number, shape, aggregation  
50 patterns, population composition or collective behavior may be used to infer anomalous conditions  
51 related to potentially catastrophic events, either natural, like harmful algal blooms, or man-made,  
52 like industrial run offs or oil spills. Intelligent systems trained on curated data could help establish  
53 the characteristics of a healthy ecosystem and detect perturbations that may represent potential  
54 threats. More importantly, given the diversity of plankton morphology and behavior across species  
55 and the growing but still limited availability of high-quality labeled data sources, there is a need  
56 for algorithms which require minimal supervision to classify and monitor plankton species with a  
57 performance approaching that of supervised algorithms. Moreover, it is also desirable for such  
58 algorithms to aid the discovery of new plankton classes, which cannot generally happen with  
59 supervised classification techniques.

60 In this paper we propose a set of novel algorithms to reliably characterize and classify plankton  
61 data. Our method is based on an unsupervised approach to overcome the limits of supervised  
62 machine learning techniques, and designed to dynamically classify plankton from instruments that  
63 continuously acquire plankton images. First, we evaluate the performances of our algorithms on  
64 a mixture of ten freshwater plankton species imaged with a lensless microscope designed for *in*  
65 *situ* data collection [5]. Next, we evaluate the performance of our algorithms on an image dataset  
66 extracted from the Woods Hole Oceanographic Institution (WHOI) plankton database [6].  
67 Machine learning methods are becoming a popular way to characterize and classify plankton [7]–  
68 [14]. A recent paper [15] explores the use of Convolutional Neural Networks to classify species of

69 zooplankton, by introducing an architecture named ZooplanktoNet. The authors claim that their  
70 customized architecture can reach higher accuracy compared to standard deep learning  
71 configurations, like VGG, AlexNet, CaffeNet, and GoogleNet. In [16] and [17], the authors use  
72 an SVM based algorithm to classify species with high accuracy from the WHOI dataset. In a recent  
73 Kaggle competition contest (<http://www.kaggle.com/c/datasciencebowl>), the authors developed a  
74 deep learning architecture named DeepSea [18] to perform accurate classification of plankton  
75 collected with an underwater camera. In [19] the authors combine features obtained with multiple  
76 kernel learning to achieve higher accuracy than classic machine learning algorithms. However, all  
77 these advancements use supervised learning algorithms that rely on large labeled training sets  
78 which are very difficult and time consuming to create. Although recent computational advances  
79 may reduce the annotation burden for large biological datasets [20], a high-performance  
80 unsupervised learning algorithm can provide an alternative for real time unbiased in situ analysis.

81

82

83

84

85

86

87

88

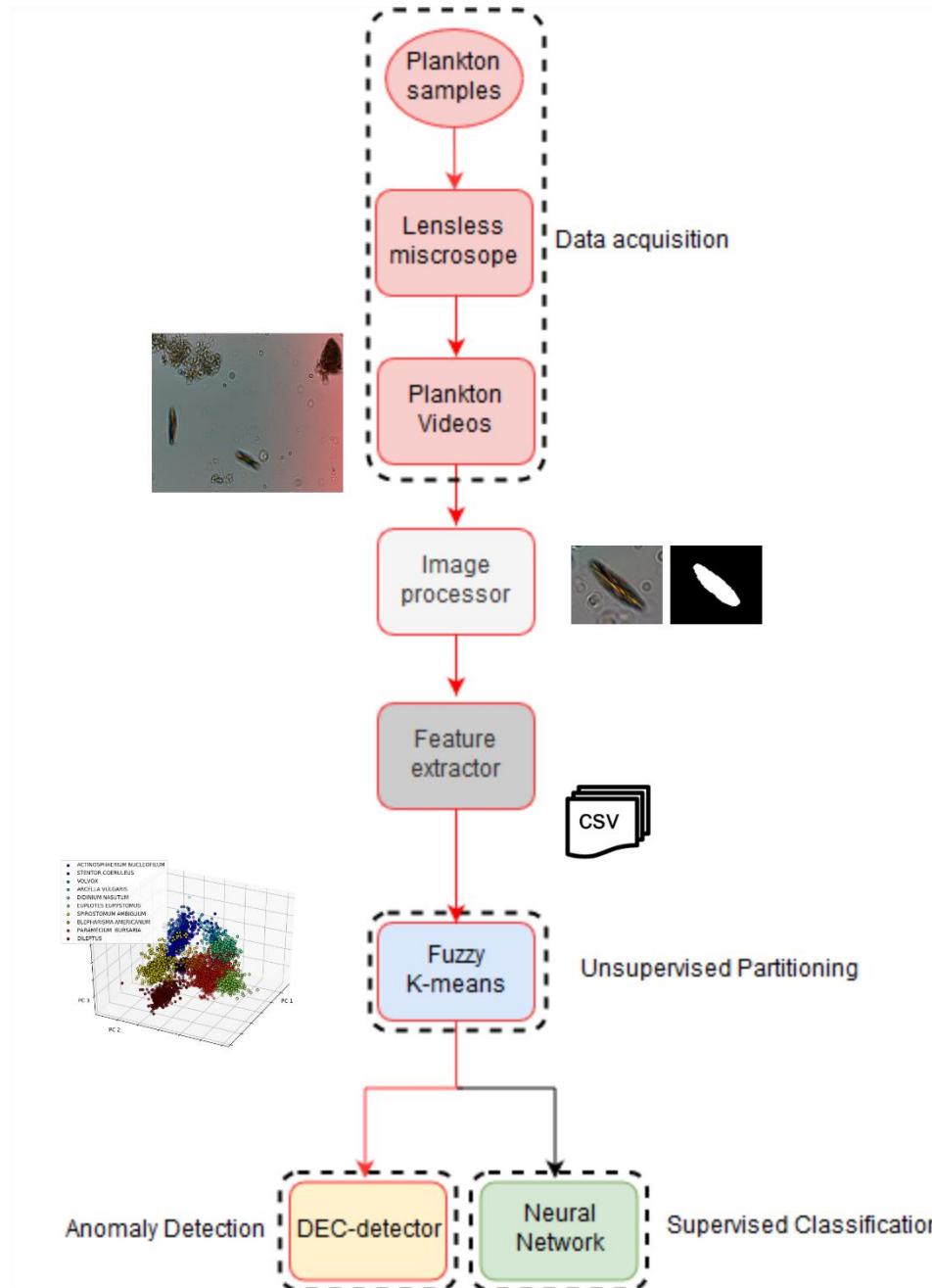
89

90

## 91 **Results**

### 92 **Plankton Classifier**

93 We developed an unsupervised customized pipeline for plankton classification and anomaly  
94 detection, that we named *plankton classifier*. The pipeline, shown in Fig 1, is tested on a collection  
95 of videos containing ten fresh water species of plankton captured with a lensless microscope [5].  
96 Each video is ten seconds long and contains one or more species. As the method is unsupervised,  
97 no labels are provided to the classifier during training. The plankton classifier consists of four  
98 modules: an **image processor**, a **feature extractor**, an **unsupervised partitioning module** and a  
99 **classification module**. The **image processor** examines each frame of video and generates cropped  
100 images of each plankter. The **feature extractor** examines each plankter image and generates a  
101 collection of features. The **unsupervised partitioning module** clusters samples by features into  
102 classes. The **classification** module comprises of a neural network-based **anomaly detector** to both  
103 perform classification based on the inferred labels and provide information to extend the database  
104 in an unsupervised manner. A sample is considered an anomaly with respect to a class if the  
105 extracted features are significantly different from the class average, as described below. The  
106 classification module also includes a standard neural network classifier, for performance  
107 comparison. See section materials and methods for a description of the modules in more details,  
108 along with the methods considered and tested that led to our final design.



109

110 *Fig 1. Schematic overview of the pipeline used to detect and classify plankton species with minimal supervision. Our preferred*  
111 *embodiment is represented by the red lines.*

112

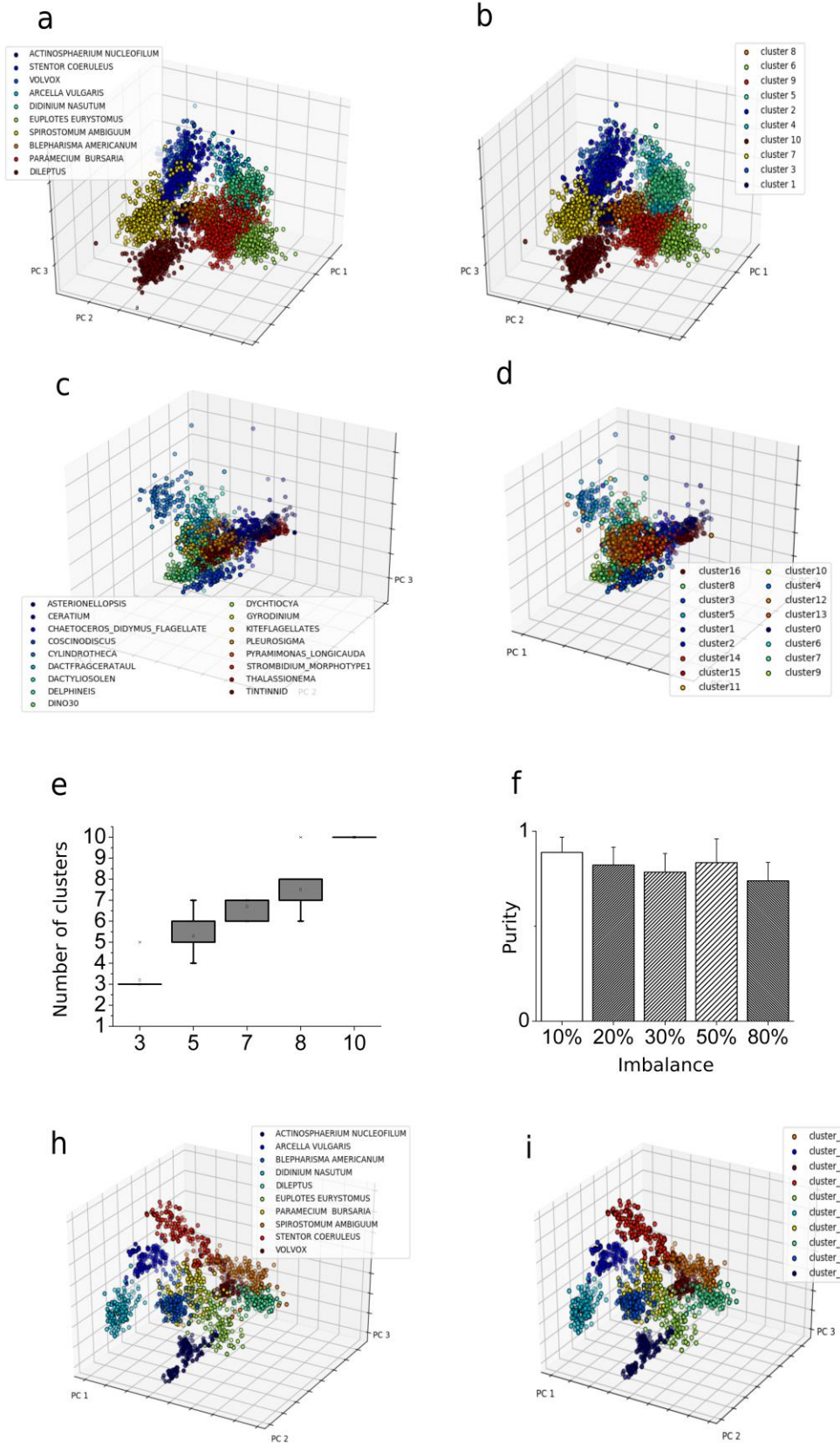
113

## 114 **Unsupervised partitioning performance**

115 First, the plankton classifier examines each frame of an acquired video and generates cropped  
116 images of each plankter. A set of 131 features is then extracted, as described in Materials and  
117 Methods. The unsupervised partitioning module uses such features to place each plankton sample  
118 into one of  $Z$  classes. To automatically obtain the number of classes from the dataset, we have  
119 designed a custom algorithm based on partition entropy (see Materials and Methods). We  
120 evaluated the robustness of the implemented method on random subsets of the lensless dataset  
121 with different sizes, ranging from three to ten species. The box plot indicating the distribution for  
122 the estimated number of clusters  $Z$  among ten iterations can be observed in Fig 2e. The inferred  
123 number of classes,  $Z$ , is correctly identified in every case. A comparison of the performance of this  
124 algorithm against other existing methods is reported in the Supporting Information. Once we have  
125 obtained the number of clusters, we compared three clustering algorithms (see Supporting  
126 Information): k-Means, Fuzzy k-Means and Gaussian Mixture Model (GMM). Clustering  
127 accuracy is evaluated using purity (see materials and methods). The Fuzzy k-Means algorithm  
128 reaches a purity value of 0.934 (see Figs 2a, 2b), outperforming the standard k-Means (purity value  
129 = 0.887) and GMM [21] (purity value = 0.886). A posterior analysis of the results of the GMM  
130 reveals that this algorithm is not able to distinguish between *Blepharisma americanum* and  
131 *Paramecium bursaria*, due to their nearly identical appearance in the acquired videos. The Fuzzy  
132 k-Means algorithm is able to match the fuzziness exhibited by the plankton classes in parameter  
133 space which explains the lower accuracy of the crisp algorithms (k-Means and GMM). Therefore,  
134 we use the Fuzzy k-Means for our unsupervised classifier. A potentially important effect on the  
135 performance of any clustering algorithm is the class imbalance. The lensless microscope dataset is  
136 composed of 500 training samples for each of the ten considered species. To evaluate the impact

137 of class imbalance, we performed the following experiment: We have built a dataset where the  
138 number of images of a species is a fraction (between 10% and 80%) of the number of images of  
139 the other species. We then evaluate the purity of this dataset and repeat the procedure for all the  
140 other species. Fig 2f reports the average performance over the ten datasets obtained as described  
141 above, as measured by the purity. The algorithm is always able to infer the correct number of  
142 species, without any overlap, with a minimum average purity value of  $0.74 \pm 0.09$  (corresponding  
143 to 80% of class imbalance) and a maximum average purity value equal to  $0.90 \pm 0.08$   
144 (corresponding to 10% of class imbalance), with a maximum purity value of 0.972. This result  
145 shows that our pipeline can accurately cluster the data even in the case of strong class imbalance.





147 **Fig 2. Unsupervised clustering results. a, b** We performed a PCA analysis on the lensless digital microscope dataset to provide  
148 a graphical representation of the data distribution into the features space. We plot the first three principal components that account  
149 for ~67% of the total variance. We assigned different colors to the different plankton species. **a** Species are assigned using ground  
150 truth labels. **b** Species are assigned to the most overlapping cluster resulting from the unsupervised partitioning procedure. **c, d**  
151 Same analysis and procedure applied on the WHOI dataset. **c** Species are assigned using ground truth labels. **d** Species are assigned  
152 to the most overlapping cluster, resulting from the unsupervised partitioning procedure. **e** Distribution of number of clusters  
153 computed using our PE algorithm for a random subset of species in the lensless microscope dataset. Results are reported for different  
154 initial number of species. **f** Effect of class imbalance. For each of the ten species included into the lensless microscope dataset, we  
155 simulated class imbalance by increasing the number of images available to the clustering algorithm for the considered species. **h, i**  
156 PCA analysis on the lensless digital microscope dataset provides a graphical representation of the data distribution into the deep  
157 features space. The unsupervised partitioning using deep features is highly accurate. The first three principal components are plotted  
158 and different colors to the different plankton species are assigned. **h** Species are assigned using ground truth labels. **i** Species are  
159 assigned to the most overlapping cluster resulting from the unsupervised partitioning.

160

## 161 **Algorithm performance on features extracted using deep feature extraction**

162 Feature selection is an important part of any unsupervised learning pipeline. Indeed, hand  
163 engineering features introduces a degree of arbitrariness, which can be removed using a method  
164 of automated feature selection. Deep feature extraction, which consists in training a neural network  
165 architecture on either in- or out-of-domain data and use the last layer before prediction to extract  
166 features [9][22], is one such method. We trained the model described in section *Convolutional*  
167 *Neural Network (CNN) for deep features extraction* using the ten classes included in our lensless  
168 microscope dataset. The model reached 99% of training accuracy, 99% of validation accuracy and  
169 98% of testing accuracy on the dataset obtained using our lensless microscope. Finally, the 128  
170 neurons from the fully connected layers preceding the output are extracted and used as features for  
171 our pipeline. The PCA computed for the lensless microscope testing set among these features can  
172 be visualized in Fig 2h. Fig 2i shows the results of the unsupervised partitioning procedure. The

173 underlying structure of the data set is very accurately captured, with a purity value of 0.98. Despite  
174 the fact that the accuracy obtained using deep feature extraction is slightly higher than the one  
175 obtained using the hand engineered features (purity of 0.980 vs 0.934), we decide to use the  
176 interpretable features described in Table 1. In fact, we think it is important that interpretability is  
177 maintained for the purpose of establishing a causal link between environmental perturbations and  
178 morphological modifications. However, for the purpose of organism classification, the customized  
179 deep feature extraction algorithm we implemented is a very viable alternative to the one proposed.

180

181

## 182 **Classification**

183 **Supervised Classifier.** At this stage of the pipeline, all samples have been assigned labels which  
184 have no correspondence to the actual plankton classes. We use the same trained clustering  
185 algorithm to classify the test samples, assigning each sample to the closest centroid. Using the  
186 trained Fuzzy k-means algorithm we reach a testing accuracy of 89%. Alternatively, one can use  
187 the labels obtained by our unsupervised partitioning algorithm to train a supervised classifier. We  
188 evaluated two algorithms: An Artificial Neural Network (ANN) and a Random Forest (RF)  
189 classifier. Our ANN architecture consists of a collection of classifiers, each trained to detect one  
190 plankton class. The RF approach consists in a set of decision trees to separate the training step  
191 samples into the correct classes.

192 For comparison, a simple ANN classifier is trained using the labels provided by the unsupervised  
193 partitioning algorithm. The ANN is a massive parallel combination of single processing units  
194 which can learn the structure of the data and store the knowledge in its connections [23]. See

195 Materials and Methods for further information and for a detailed description of the implemented  
196 architecture. The network is very shallow, providing an efficient feature selection process. The  
197 ANN classifier reaches a validation accuracy of 99% and a testing accuracy of 94.5%. Figs 3c and  
198 3d report the ROC curves and the confusion matrix obtained by testing the trained ANN classifier  
199 on our ten species plankton dataset. The ROC curves are close to a perfect classifier and the  
200 confusion matrix is almost diagonal with minor overlap between two pairs of species: *Blepharisma*  
201 *americanuum-Paramecium bursaria* and *Spirostomum ambiguum-Stentor coeruleus*. This  
202 misclassification is primarily due to the similarity in the shape, size and texture of the two pairs of  
203 species, influencing both the unsupervised training clustering and the subsequent testing of the  
204 supervised classifier.

205 An alternative classifier method employs a Random Forest (RF) approach, a popular ensemble  
206 learning method used for classification and regression tasks.

207 We train an RF algorithm using the labels provided by the unsupervised classifier and reach an  
208 accuracy of 94%. For comparison, we train the same RF algorithm using the actual labels (ground  
209 truth) of the training set and reach an accuracy around 98%, proving that our unsupervised  
210 classification approach performs comparably well with respect to the correspondent supervised  
211 approaches for the trained classifier. Since the ANN performs marginally better than the RF  
212 classifier, we propose the former for a pipeline. In the next section, we will present an alternative  
213 classification method

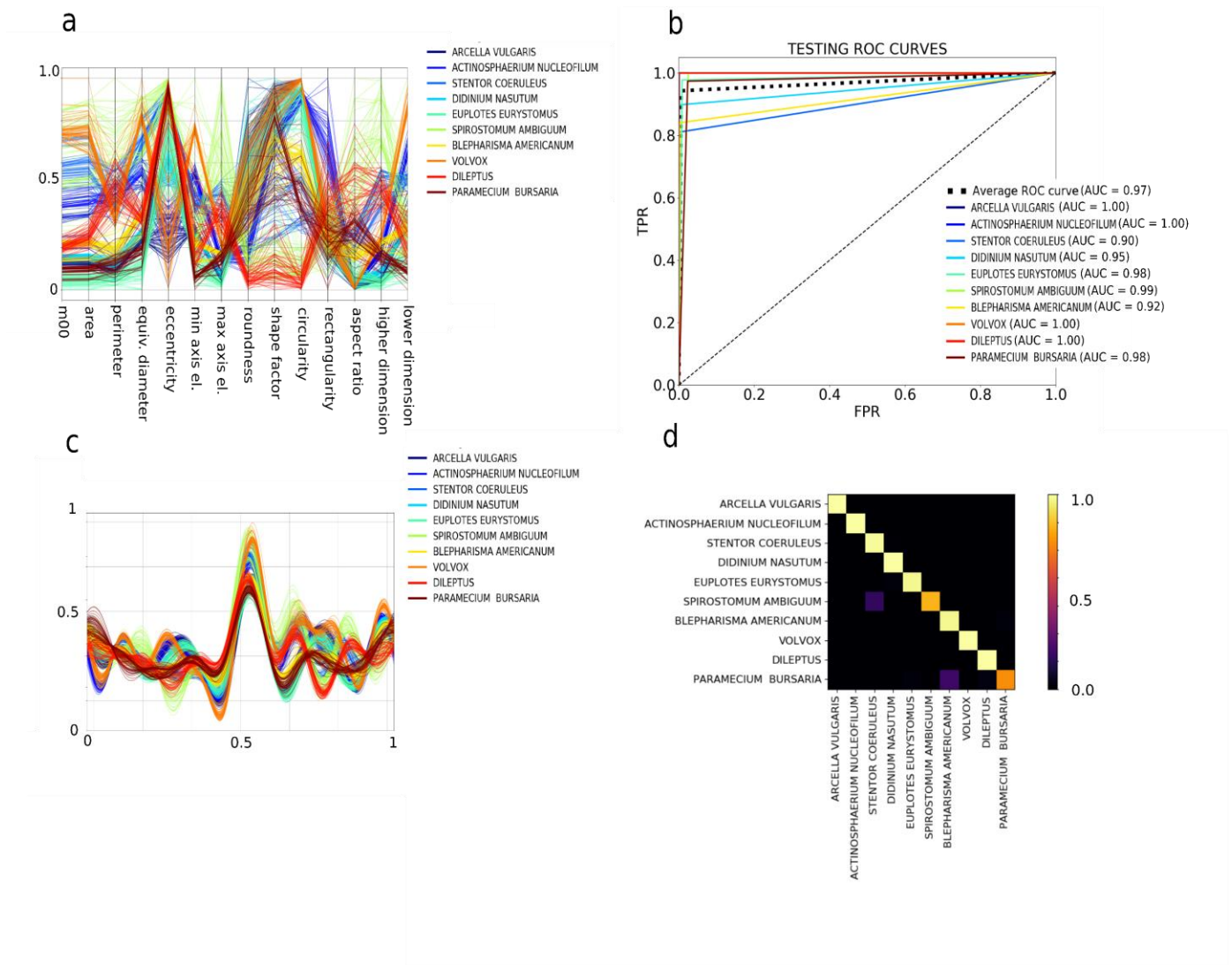
214

## 215 **Anomaly Detector**

216 When deployed in the field, microscopes will encounter species that have never been seen before,  
217 so it is essential that such samples are detected and correctly identified as anomalies. For a given

218 class, a sample is considered an anomaly if the sample features are significantly different from the  
219 feature average for the class. Algorithms for anomaly detection based on the separation of the  
220 features space have been successfully used to identify the intrusion in computer networks for  
221 security purposes [24]. Two anomaly detectors are implemented and compared; a state of the art  
222 one-class SVM<sup>15</sup> and a customized neural network we call a Delta-Enhanced Class (DEC) detector  
223 that combines classification with anomaly detection. The one-class SVM algorithm uses a kernel  
224 to project the data onto a multidimensional space and can be interpreted as a two class SVM  
225 assigning the origin to one class and the rest of the data to another class. It then solves an  
226 optimization problem determining a hyperplane with maximum geometric margin, i.e., a surface  
227 where the separation between the two sets of points is maximal, that will be used as decision rule  
228 during the testing step.

229 A customized one-class SVM is implemented by normalizing the testing samples using the training  
230 data belonging to a single class. In this way, there will be a significant difference in the absolute  
231 value obtained for the anomaly (out-of-class) samples compared to the in-class samples, improving  
232 the accuracy of the SVM. The one-class SVM so designed reaches an average testing accuracy of  
233  $(93.5 \pm 6.0) \%$ , with high accuracy in both anomaly detection and classification.



234

235 **Fig 3. Feature space representation and classification performances. a, b** Multidimensional visualization of the geometric  
 236 subset of the ten species in the lensless microscope dataset, obtained using the following methods (see Supporting Information): **a**  
 237 Andrew's curve. **b** Parallel coordinates. **c** ROC curves obtained for the neural network classifier trained on the labels provided by  
 238 the clustering algorithm for the lensless microscope dataset. **d** Corresponding confusion matrix.

239

240

241 We now describe an alternative ANN-based approach that simultaneously performs classification  
 242 and anomaly detection. As demonstrated above, a single layer ANN is able to satisfactorily classify  
 243 plankton data from our in-house dataset. However, to effectively approach the anomaly detection

244 step, we designed a deep neural network called Delta-Enhanced Class (DEC) detector (see  
245 materials and methods for further details). One DEC detector must be trained for each of the  
246 training species. Therefore, we train ten DEC detectors, one for each of the species of plankton  
247 identified in the unsupervised learning step. This procedure affords excellent accuracy on both  
248 classification and anomaly detection, on both real and simulated plankton data (see Fig 4), with an  
249 average testing accuracy on real data of  $98.8 \pm 2.4$  %, an average anomaly detection testing  
250 accuracy of  $99.2 \pm 0.7$  % and an average overall testing accuracy of  $99.1 \pm 0.9$  % (see Fig 4b for  
251 details). The confusion matrices in Fig 4a demonstrate the discrimination power of our algorithm.  
252 The DEC detector outperforms the alternative one-class SVM classifier in both supervised  
253 (average accuracy equal to 95%) and unsupervised (average accuracy equal to 93.5%)  
254 configurations. It is worth reporting that the unsupervised one-class SVM reached a minimum  
255 overall accuracy of 79%, compared to 97.2% for the DEC detector (minimum values correspond  
256 to *Paramecium bursaria* detector). To test the overall performance of our method, we produce a  
257 dataset of surrogate plankton organisms. For each different species, we test the corresponding DEC  
258 detector architecture using a surrogate species created with a feature-by-feature weighted average  
259 of all the species in our dataset. Starting with a uniform weight distribution, we increase the weight  
260 for the species corresponding to the trained DEC detector architecture up to 0.9 (steps of 0.1),  
261 obtaining 9 different surrogate species (see Fig 4d for an average parallel coordinates plot, showing  
262 the resulting distributions for the species *Spirostomum ambiguum*). The aim of this robustness test  
263 is to simulate the acquisition of an unknown species, whose features are increasingly closer to the  
264 features of the class correspondent to the detector, up to a maximum of 90% similarity. As Fig 4e  
265 shows, our classifier can recognize the synthetic species as an anomaly with an average accuracy  
266 higher than 98% if the similarity between the synthetic and the real species is up to 30%, and it

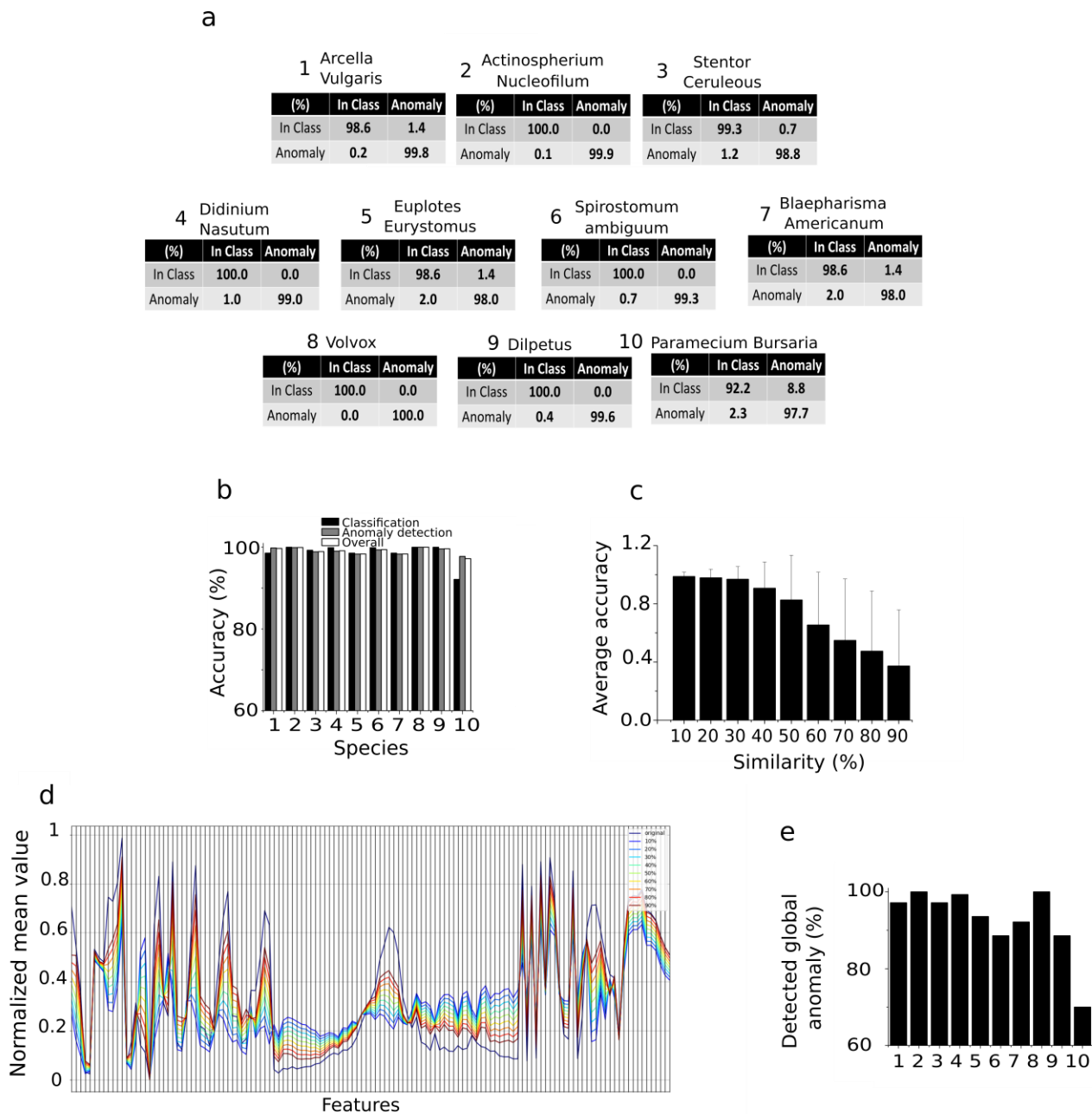
267 can maintain an average accuracy of over 82.6% if the species similarity is up to 50%. Accuracy  
268 of anomaly detection severely decreases if the species similarity is over 50%, reaching the  
269 minimum value of 37.5%.

## 270 **Plankton classifier performance on the WHOI dataset**

271 The WHOI provides a public dataset comprising millions of still monochromatic images of  
272 microscopic marine plankton, captured with an optical Imaging FlowCytobot  
273 (<https://mclanelabs.com/imaging-flowcytobot/>). To use this dataset as a benchmark to test our  
274 unsupervised classifier, we extract a set of 128 features from a collection of 40 species of plankton  
275 (100 images per species, randomly selected), using both the segmented binary image and the  
276 portion of the gray-scale image containing the plankton cell body. A full description of the species  
277 selection process is reported in the Supporting Information. The features set is identical to the one  
278 used for the lensless microscope dataset, except for the absence of three-color features, as the  
279 lensless microscope is a color-based sensor, while the Imaging FlowCytobot is monochromatic.  
280 Figs 2c, 2d show the results of our pipeline applied on the normalized features set. The algorithm  
281 reaches an overall purity value of 0.715 for the 40 WHOI species that we selected. The ability of  
282 our pipeline to distinguish between inter-species plankton morphology can be further observed  
283 comparing Fig 2c, which represents the PCA space corresponding to a subset of 18 of the 40  
284 species for the ground truth dataset, and Fig 2d, which represents the corresponding PCA space  
285 resulting from the unsupervised partitioning algorithm. A complete PCA representation for the 40  
286 species can be found in Supporting Information. We trained a random forest algorithm using the  
287 labels provided by the unsupervised partitioning with a train-test ratio of 80:20, obtaining a  
288 classification accuracy around 63%. For comparison, we have trained a supervised random forest



289 algorithm using the ground truth labels on the extracted features, obtaining a classification  
 290 accuracy around 79%.



291  
 292 **Fig 4. Delta-Enhanced Class detector performances and results.** **a** Confusion matrix corresponding to each of the ten neural  
 293 networks trained on the lensless microscope dataset. **b** Overall testing accuracy performances for each of the ten testing classes.  
 294 The number used on x axis to label each species correspond to the species number in panel **a**. **c-d** DEC detector anomaly detection

295 performances tested on in silico generated data. **d** Testing accuracy performances for varying percentage values of in silico species  
296 similarity with the trained species. **e** Example of average features space parallel coordinates plot for the in-silico species obtained  
297 using the species *Spirostomum Ambiguum*. By increasing the similarity, the features of the surrogate species approach the features  
298 of the real species, resulting in an increased average anomaly misclassification rate, decreasing the overall accuracy levels. **e**  
299 Detection of unknown species. The panel shows the percentage of samples detected by all the DEC detectors as anomaly, when  
300 removing one training species from the set, for each of the ten training species. These numbers reflect the level of accuracy of the  
301 proposed algorithm in detecting unseen species. The number used on x axis to label each species correspond to the species number  
302 in panel **a**.

### 303 **The plankton classifier can reveal unseen species**

304 We have demonstrated that our DEC neural networks are able to classify a sample as either a  
305 training class (i.e., the plankton species used to train the detector) or as an anomaly. If a sample is  
306 discarded by all the implemented detectors, it could either represent an intra-species anomaly (i.e.,  
307 species included into the training set) or a sample belonging to an unseen species (i.e., species not  
308 included in the training set). The former represents the basis for using the proposed pipeline for  
309 real-time environmental monitoring, and its implications are discussed in the next section. We now  
310 test the potential of our pipeline to detect new species. We remove one class from our unsupervised  
311 partitioning ensemble set, consider it as never before seen and compute the number of testing  
312 samples detected as anomaly by all the remaining DEC detectors. This number indicates the  
313 algorithm accuracy in detecting new species. We repeat the procedure for each class. The average  
314 detection accuracy is  $98.3 \pm 10.1$  % (see Fig 4e), demonstrating the ability of the pipeline to detect  
315 the presence of a new species. If two or more unseen species are detected, they will be stored as  
316 anomalies. As this group of anomalies grows, a human expert may determine offline the actual  
317 labels for these new species, thus allowing a DEC detector to be trained for each new species.  
318 Alternatively, the samples corresponding to unseen species may be clustered and classified by the

319 unsupervised partitioning step of our pipeline, reducing the number of new species that must be  
320 examined by a human.

## 321 **Discussion**

322 The plankton classifier described in this paper provides the foundation for a robust, accurate and  
323 scalable mean to autonomously survey plankton in the field. We have identified interpretable and  
324 non-interpretable image features that work with our algorithms to perform an efficient clustering  
325 and classification on plankton data using minimal supervision and with a performance accuracy  
326 comparable to supervised learning algorithms [16]. Instead of labeling thousands of samples, an  
327 expert need only identifying one member of cluster to label all the samples of the cluster.

328 We introduced a neural network that performs classification by learning the shape of the feature  
329 space and uses this information to identify anomalies. The network uses a novel unbiased  
330 methodology of feature-to-feature comparison of a test sample to a random set of training samples.  
331 While most of the existing classification methods require various degrees of user input, our method  
332 is automated, without sacrificing performance accuracy or efficiency.

333 All features the plankton classifier relies upon are extracted from static images. However, our  
334 custom lensless microscope captures 2D and 3D dynamic of plankton. While this dynamic  
335 information is not considered in the analysis presented here, motion data can increase the  
336 dimensionality of the feature space, by adding spatio-temporal “behavioral” components, and may  
337 improve the performance of classifiers and anomaly detectors. This is particularly valuable in cases  
338 where species have considerable overlap in morphology feature space, as seen with *Blepharisma*  
339 *americanum* and *Paramecium bursaria*, and *Spirostomum ambiguum* and *Stentor coeruleus*,

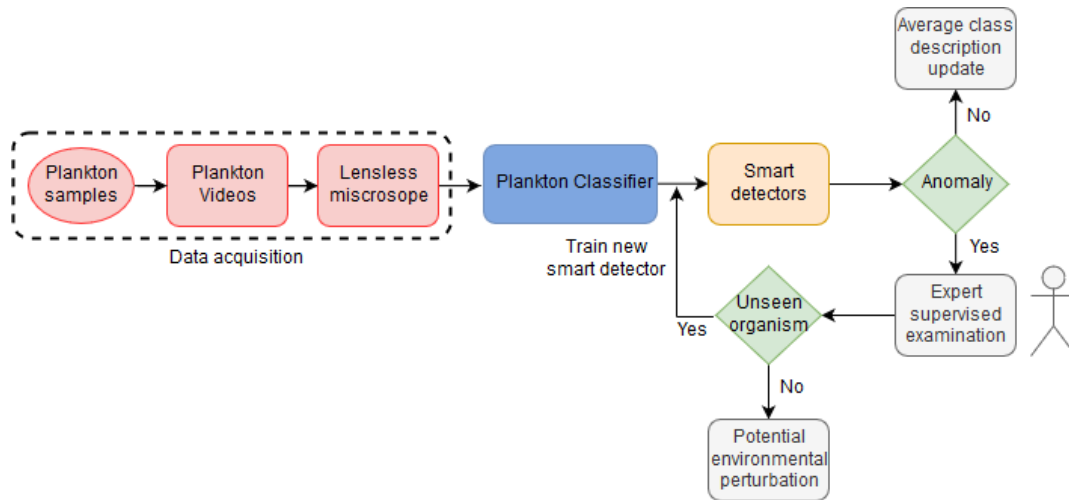
340 shown in the confusion matrices in Fig 3d. Currently, existing large plankton datasets, like the  
341 WHOI used in our validation experiments, are based on static images, but as the cost of video-  
342 based *in situ* microscopes drops and their deployment increases, we believe datasets that include  
343 spatio-temporal data will become available and the use of such features will gain importance.

344 Deploying smart microscopes capable of real-time continuous monitoring will give biologist an  
345 unprecedented view of plankton *in situ*. The adoption of an unsupervised unbiased pipeline is a  
346 significant step ahead in the development of a real-time “smart” detector for environmental  
347 monitoring. Several high-resolution acquisition systems for real-time plankton imaging already  
348 exist [25] and could adopt the pipeline proposed into this paper. Fig 5 shows a high-level  
349 representation of a continuous environmental monitoring system in the form of a flow chart,  
350 showing an example of how the detector could be coupled to the computational pipeline we  
351 designed. Once the descriptors have been extracted from the acquired videos, it is possible to use  
352 them to build a set of DEC detectors. It is important to stress that the size of the data likely to be  
353 acquired, or already present in the databases, makes neural networks the obvious choice to carry  
354 out the analysis due to their unsurpassed scalability. Our newly designed and customized DEC  
355 detector neural architecture for plankton classification and anomaly detection is a functional and  
356 efficient example of such algorithm. Moreover, neural algorithms can infer non-linear  
357 relationships between features (input) and correlate them with the class description (output)  
358 without making any assumptions on the underlying learning model. Hence, the classification  
359 depends only on the extracted features. Every time the network identifies a species belonging to a  
360 specific class, the average set of morphological features is then updated, thereby further qualifying  
361 the class morphology phase space. If an anomaly is detected, it may be sent to an expert for a  
362 supervised examination. The expert will determine whether that sample could be a species not

363 represented in the training set, or if it belongs to an existing training class, but its morphological  
364 features deviate significantly from the average features space of the corresponding class. In the  
365 former case, a new smart detector will be trained offline, so that the training set is dynamically  
366 expanded, and the system will provide a continuous monitoring of the aquatic environment using  
367 the human expert-in-the-loop paradigm. In the latter case, the identified anomalies may represent  
368 local environmental perturbations, either natural or man-made. Further work is needed to assess  
369 the validity of such hypothesis. An additional re-training step may be necessary to update the  
370 algorithms. Our pipeline is based on local analysis using a low powered device, capable of image  
371 capture and processing, classification and anomaly detection. Coupling such platform with a local  
372 (laptop, server) or cloud-based system where the training step may occur could provide the  
373 flexibility and resources needed to close the loop and generate the training data the low power  
374 platform can use for classification. Examples of systems that use this paradigm are already present  
375 in the literature [26], and we hope the availability of computational paradigms like the one we  
376 propose may increase the research in the field. A high-resolution plankton acquisition system  
377 placed in the water and powered with our unsupervised pipeline may enable the development of  
378 real time continuous smart environmental monitoring systems that are fundamentally needed to  
379 stakeholders and decision-making bodies to monitor plankton microorganisms and, consequently,  
380 the entire aquatic ecosystem [27].

381 Finally, it is interesting to consider if such unsupervised approach can be utilized for different data  
382 types, thus widening the potential applicability and interest of the technique. While an extensive  
383 analysis of the performance of our pipeline on diverse set of data is beyond the scope of this work,  
384 it is worth commenting that the algorithms we use are general and pose no evident drawback to  
385 their application to other cell types. Particularly, the features our classifier uses to cluster the

386 images do not include anything specific to plankton species (e.g. detection and estimation of  
387 number of flagella or other organelles.) Moreover, the proposed Deep Feature extraction method  
388 is even less dependent on the kind of data under study and may increase the applicability to other  
389 cell types. Thus, we expect the method to be potentially useful to other biological imaging fields.



390

391 **Fig 5. Proposed real-time smart environmental monitoring pipeline.**

392

393

394

395

396

397

398

399

400

## 401 **Material and methods**

402 The proposed unsupervised pipeline (i.e., the plankton classifier) shown in Fig 1, consists of four  
403 modules: an **image processor**, a **feature extractor**, an **unsupervised partitioning module** and a  
404 **classification module**. In the following paragraphs we provide a description of the modules in  
405 more details, along with the methods considered and tested that led to our final design.

### 406 **Image Processing**

407 Each video consists of ten seconds of color video (1920x1080) captured at 30 frames per second.  
408 Background subtraction is applied to each frame to detect the swimming plankton in the image. A  
409 contour detector is applied to the processed image to create a bounding box around each plankter.  
410 Because of instrument design, organisms can swim in and out of the field of view (FOV) during  
411 acquisition. Our algorithm automatically selects organisms which are fully contained inside the  
412 FOV by checking whether the bounding box touches the borders of the FOV. In this way, the  
413 images we obtain will be only of fully visible organisms. The resulting cropped image is then  
414 saved. From this collection of images, a training set of 640 images (500 training and 140 testing)  
415 is selected for each class. An image processor module for static images has also been implemented  
416 for benchmarking the plankton classifier on existing plankton datasets (e.g., the WHOI dataset;  
417 See Supporting Information for further details.).

### 418 **Feature Extraction**

419 For each plankter image, 131 features are extracted from four categories: geometric (14), invariant  
420 moments (32), texture (67) and Fourier descriptors (10). Geometric features include area,  
421 eccentricity, rectangularity and other morphological descriptors, that have been used to distinguish

422 plankton by shape and size [16]. The invariant Hu [28](7) and Zernike moments [29] (25) are  
 423 widely used in shape representation, recognition and reconstruction. Texture based features encode  
 424 the structural diversity of plankton. Fourier Descriptors (FD) are widely used in shape analysis as  
 425 they encode both local fine-grained features (high frequency FD) and global shapes (low frequency  
 426 FD). A full list of the features we have selected is reported in Table 1. These features span a 131-  
 427 dimensional space, capturing the biological diversity of the acquired plankton images. Figs 3a and  
 428 3b demonstrate as an example, the discriminating power of the geometrical features for the ten  
 429 evaluated species.

Class	Number	Description
<b>Geometric feature</b>	14	Area (pixels), area (0-th order moment), perimeter, eccentricity, rectangularity, roundness, shape factor, width and height (minimum fitting rectangle), circularity, major and minor axis (fitting ellipse), equivalent diameter, convexity.
<b>Hu moments</b>	7	Hu moments computed from the normalized central image moments.
<b>Zernike moments</b>	25	Zernike moments up to order 5.
<b>Image intensity features</b>	8	Blue/green channels ratio, red/green channels ratio, red/blue channels ratio, gray levels histogram statistical features (skewness, kurtosis, mean value, standard deviation and entropy)
<b>Haralick features</b>	13	The first 13 features as proposed by Haralick in <sup>13</sup> , computed from the Gray Scale Co-occurrence Matrix (GSCM).
<b>Local binary patterns</b>	54	Local binary patterns summarize structures of the image comparing each pixel to its neighborhood
<b>Fourier descriptors</b>	10	Fourier descriptors are contour-based features invariant with respect to rotation, scaling and translation.

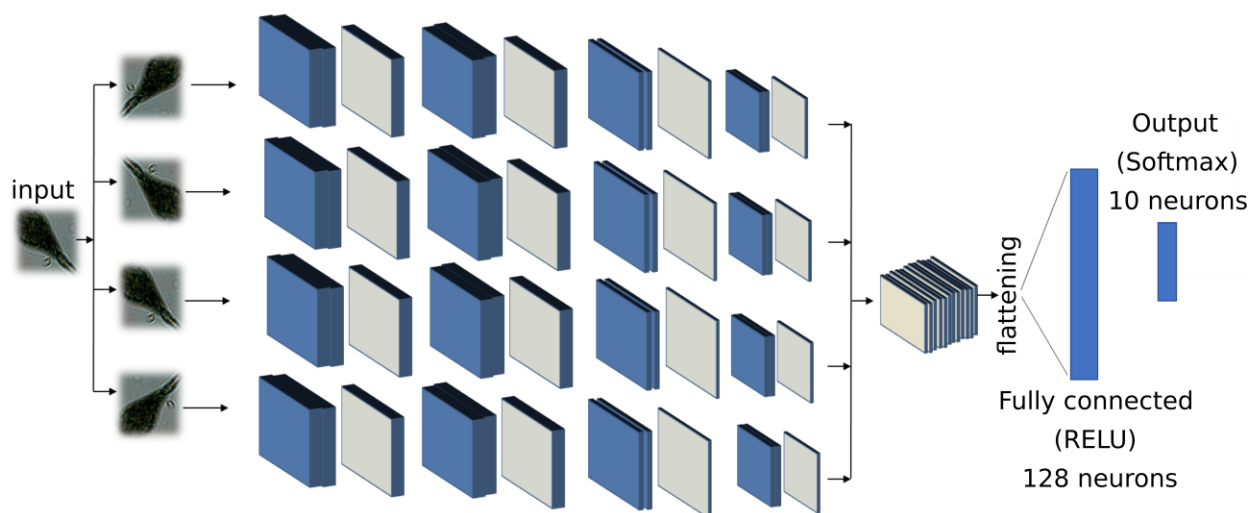
430



431 **Table 1: List of morphological features extracted from the processed images.** See Supporting Information for a detailed  
432 explanation.

### 433 **Convolutional Neural Network (CNN) for deep features extraction**

434 We implemented a deep CNN using eight convolutional layers and two fully connected layers, as  
435 described in Fig 6. We customized our architecture to be invariant with respect to rotation, similar  
436 to what has been done in [18]. Each input sample is rotated four times at multiples of 90 degrees,  
437 and all the tensors resulting from the features extraction module are concatenated and used to train  
438 the fully connected layers. The neural network has been trained for 60 epochs, using stochastic  
439 gradient descent with learning rate equal to  $10^{-5}$ , using data augmentation by means of translation,  
440 zooming, and rotation. It is worth noticing that the implemented rotational invariance module  
441 actually performs a data augmentation operation, and it is indeed useful when partial training data  
442 are available.



444 **Fig 6. Deep features extraction.** Deep CNN implemented for the purpose of deep features extraction. The blue layers represent  
445 convolutional layers, the grey ones represent a max pooling 2D operation. The fully connected layer with 128 neurons output has  
446 been used as feature set to the subsequent modules in our pipeline.

## 447 **Unsupervised Partitioning**

### 448 **Partition Entropy (PE)**

449 The Partition Entropy (PE) coefficient is defined as:

$$450 \quad PE = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K u_{ij} * \log(u_{ij}) \quad (1)$$

451  
452  
453 The coefficient is computed for every  $j$  in  $[0, K]$  and takes values in range  $[0, \log(K)]$ . The  
454 estimated number of clusters is assigned to the index  $j^*$  corresponding to the maximum PE value,  
455  $PE(j^*)$ . The lower the  $PE(j^*)$ , the higher the uncertainty of the clustering. We repeat this procedure  
456 ten times and obtain a distribution of  $j^*$ . Finally, the estimation of the number of clusters  $Z$  is the  
457 mode of this distribution.

### 458 **Clustering accuracy**

459 Clustering accuracy is evaluated using purity:

$$460 \quad purity = \frac{1}{N} \sum_k \max_j |w_k \cap c_j| \quad (1)$$

461 where the class  $k$  is associated to the cluster  $j$  with the highest number of occurrences. A purity  
462 value of one corresponds to clusters that perfectly overlap the ground truth. Purity decreases when  
463 samples belonging to the same class are split between different clusters, or when two or more  
464 clusters overlap with the same species. We have implemented a purity algorithm capable of  
465 checking for these occurrences and automatically adapt to the correct number of non-overlapping  
466 clusters (see Supporting Information).

## 467 **Classification algorithms**

### 468 **Random Forest**

469 Random Forests (RF) is a popular ensemble learning method [30] used for classification and  
470 regression tasks, introduced in 2001 by Breiman. Random forests model providing estimators of  
471 either the Bayes classifier or the regression function. Basically, RF work building several binary  
472 decision trees using a bootstrap subset of samples coming from the learning sample and choosing  
473 randomly at each node a subset of features or explanatory variables [31]. Random forests are often  
474 used for classification of large set of observations. Each observation is given as input at each of  
475 the decision tree, which will output a predicted class. The model outputs the class that is the mode  
476 of the class output by individual trees [32].

477 Let us consider a set of observations  $x_1, x_2, \dots, x_N$ , with  $x \in R^m$ . The decision tree is designed as  
478 follows: we extract N times from the set of training observations (with replacement), for a each of  
479 the total number of decision tree. We specify the number of features  $m^*$  to consider for the tree  
480 growing, with  $m^* \ll m$ . For each of the nodes in the tree, the algorithm randomly selects  $m^*$   
481 features and calculates the best split for that node. The trees are only grown and not pruned (as in  
482 a normal tree classifier [33]). The split's aim is to reduce the classification error at each branch. In  
483 detail, the algorithm considers an entropy-based measure trying to reduce the amount of entropy  
484 at each branch, selecting, with such a procedure, the best split. A possible choice is the Gini index:

485

$$486 \quad G_m = \sum_{i=1}^K p_{im}(1 - p_{im}) \quad (27)$$

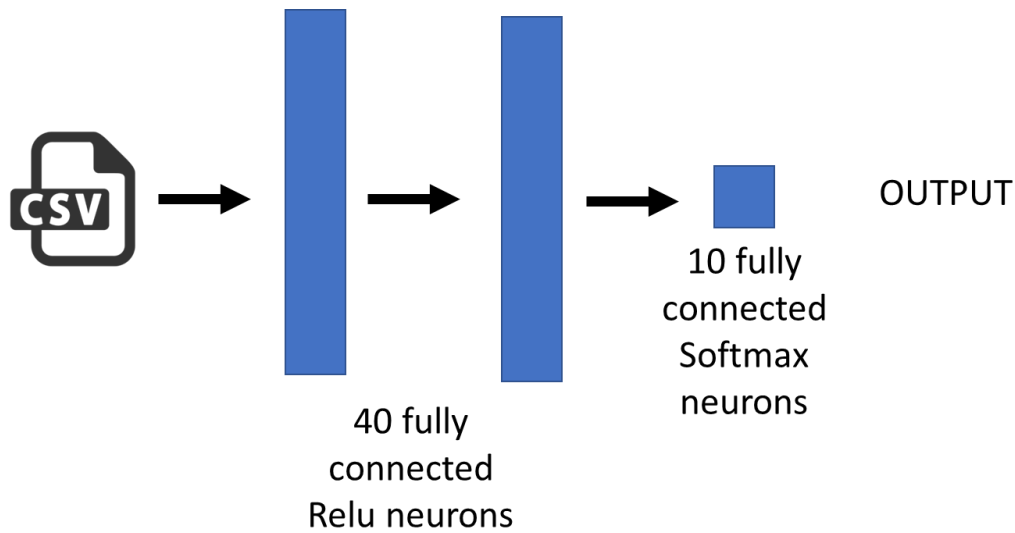
487

488 Where  $G_m$  is the Gini Index for branch at level  $m$  in the decision tree, and  $p_{im}$  is the proportion of  
489 observations assigned to class  $i$ . Minimizing  $G_m$ , means to decrease the heterogeneity at each  
490 branch, i.e., a best split will correspond to a lower number of class in the children nodes. The  
491 algorithms continue in growing trees until convergence on the entropy-based on the generalization  
492 error [32].

### 493 **Neural Networks**

494 An artificial neural network (or multi-layer perceptron) is a massive parallel combination of single  
495 processing unit which can acquire knowledge from environment through a learning process and  
496 store the knowledge in its connections [23]. Classification is one of the most active research and  
497 application areas of neural networks. In this work we used an artificial neural network to build a  
498 classifier able to predict the species for each observation extracted using the shadow microscope.  
499 Fig. 2 shows the developed architecture. The network is very shallow, with two hidden layers of  
500 40 neurons and an output layer with as much neurons as the number of species to classify. As  
501 reported in the main text of this manuscript, we used a training dataset with 10 species, thus the  
502 output layer is made up of  $k$  neurons, where  $k$  is the number of clusters obtained using the  
503 unsupervised clustering. As Fig 7 shows, the developed NN uses RELU activation function and  
504 dropout to reduce the overfitting. The network was trained using 200 epochs, Root mean square  
505 as an optimizer, a learning rate  $\lambda = 0,005$  and categorical cross-entropy as loss function. The  
506 training requires 50 seconds on a MAC book PRO, core i7 – 2.9 GHz, solid state disk and 16 GB  
507 of RAM. The neural network has been implemented using KERAS, a powerful high-level neural  
508 network API running on top of TensorFlow.

509



510

511 **Fig 7. ANN architectures implemented for classification based on the extracted features.**

## 512 **Anomaly Detection**

### 513 **One Class SVM**

514 We adopted the one class SVM described by Scholpoff in [34]. Let us consider a set of N  
515 observations:  $\{x_i, y_i | x_i \in R^m, y_i = +1\}$ . Where  $x_i$  is a m-dimensional real vector and  $y_i = +1$   
516 simply imply that the set contains normal observations belonging to a certain class. The one-class  
517 SVM is a classification algorithm returning a function which takes +1 in a “small” region capturing  
518 most of the data points, and -1 elsewhere. Let  $\phi$  be a feature map that map our observations set  $x_i$ ,  
519 into an inner product space such as the inner product for the image of  $\phi$  can be evaluated using  
520 some simple kernel:

521

$$522 \quad k(x, y) = \phi(x)\phi(y) \quad (28)$$

523

524 The strategy of the one class SVM is to map the data into the kernel space and separate the data  
525 from the origin with maximum margin, defining a hyperplane as:

526

$$527 \quad H(x) = w \phi(x) - \rho \quad (29)$$

528

529 Meaning that we want to maximize the ratio  $\frac{\rho}{\|w\|}$ , corresponding to the hyperplane's distance  
530 from the origin. In order to solve this maximization problem, we have to solve a quadratic  
531 problem:

532

$$533 \quad \min_{w, \xi, \rho} \frac{1}{2} \|w\|^2 + \frac{1}{\nu l} \sum_i \xi_i - \rho \quad (30)$$

534

535 subject to  $w \phi(x) \geq \rho - \xi_i, \xi_i \geq 0$ .

536

537 Where  $\phi(x)$  is the feature mapping function that maps observations  $x$  into a feature space,  $\xi_i$  is a  
538 slack variable for outlier that allows observations to fall on the other side of the hyperplane  
539 ,  $\nu \in [0,1)$  is a regularization parameter determining the bounding for the fractions of outliers  
540 and support vectors.

541 If  $w$  and  $\rho$  solve this problem, then the decision function:

542

$$543 \quad f(x) = \text{sgn}(H(x)) \quad (31)$$

544

545 will be positive for most of the training observation, while  $w$  will be still small. The parameter  
546 influences the trade-off between the reported properties. To solve the quadratic form, we can use  
547 Lagrangian multipliers, obtaining:

548

$$549 \quad L(w, \xi, \rho, \alpha, \beta) = \frac{1}{2} \|w\|^2 + \frac{1}{vl} \sum_i \xi_i - \rho - \sum_{i=1}^l \alpha_i (w * \phi(x_i) - \rho + \xi_i - \sum_{i=1}^l \beta_i \xi_i) \quad (32)$$

550 And set the derivatives with respect to  $w$ ,  $\xi$  and  $\rho$  and expanding using the kernel expression  
551 yields:

552

553

$$554 \quad f(x) = \text{sgn} \left( \sum_i \alpha_i k(x_i, x) - \rho \right) \quad (a)$$

555

$$556 \quad \alpha_i = \frac{1}{vl} - \beta_i \leftrightarrow 0 \leq \alpha_i \leq \frac{1}{vl} \quad (b) \quad (33)$$

557

$$558 \quad \sum_{i=1}^l \alpha_i = 1 \quad (c)$$

559

560

561 We used a Radial Basis Function kernel (RBF):

562

$$563 \quad k(x_i, x) = e^{-\frac{\|x_i - x\|^2}{2\sigma^2}} \quad (34)$$

564

565 And then the original quadratic problem is solved substituting Eq. 16 into Eq. 15, yielding:

566

567

$$568 \quad \min_{\alpha} \sum_{i=1}^l \alpha_i \alpha_j k(x_i, x_j) \quad (35)$$

569

570 under the constraint of Eq. (16b) and (16c).

571

572

573 We finally use the support vectors  $x_i$  to recover the parameter  $\rho$  needed to compute the

574 hyperplane:

575

576

577

$$578 \quad \rho = w \phi(x_i) = \sum_j \alpha_j k(x_j, x_i) \quad (36)$$

579

## 580 **DEC detectors**

581 We designed a deep neural network that we named Delta-Enhanced Class (DEC) detector for the

582 purpose of anomaly detection. The DEC detector's architecture is represented in Fig 8, and shows

583 a 2-neurons output, indicating that the sample is a member of the class or is an anomaly (i.e. not a

584 member of the class). For each observation, we train such neural network with the actual features

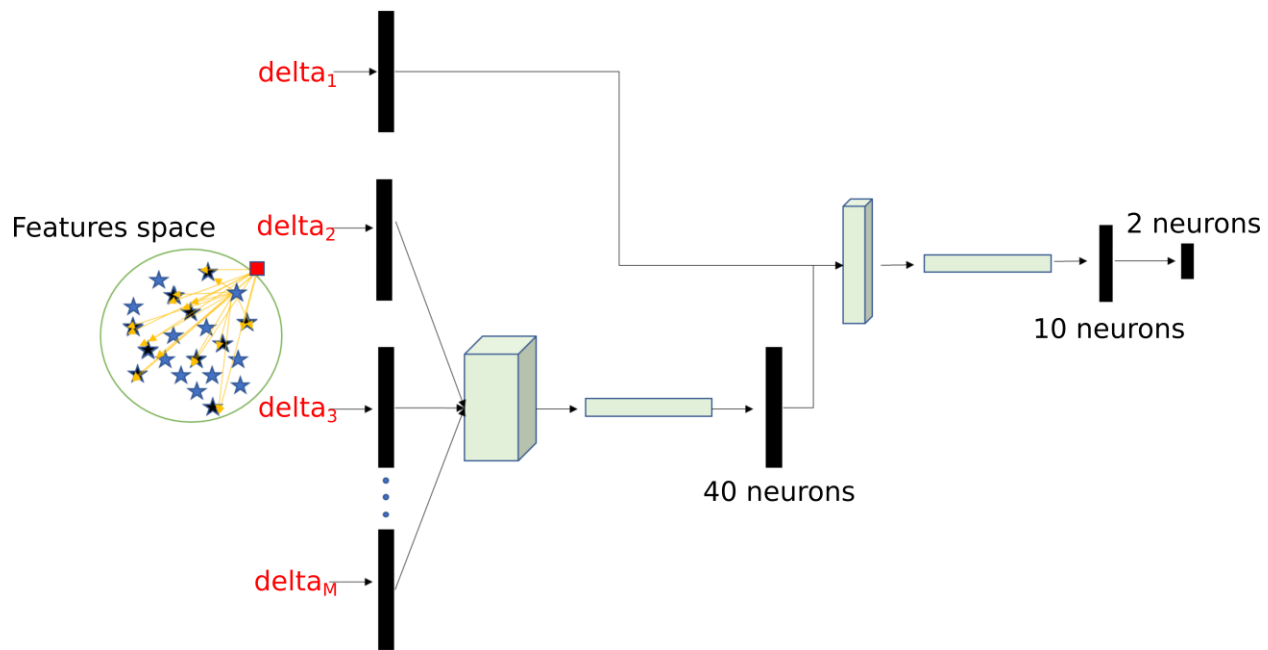
585 vector and extract randomly select a set of points from the training class in our dataset. For each

586 of these selected points, we define a custom network layer (delta layer) that computes the

587 difference in absolute value (as a vector, feature by feature) between the actual observation and



588 the extracted random set. The vector of differences and the actual observations are used as inputs  
589 to the neural network (Fig 8), which assigns the proper weights to either one during training. The  
590 set of points to select is a hyperparameter which needs to be tuned. Through testing we determine  
591 that 25 points is the optimal tradeoff accuracy and computational cost.



593 **Fig 8. Schematic representation of DEC detector architecture.**

594

595

596

### 597 **Code availability**

598 The full source code accompanying this paper has been made available under EPL license at the  
599 following link: <https://github.com/sbianco78/UnsupervisedPlanktonLearning>.

600

## 601 **Supporting information**

602 **S1 Data. The lensless microscope dataset and the dataset extracted from the WHOI used in**  
603 **this paper is available at the following link:**

604 <https://ibm.ent.box.com/s/8g2mp5knl2by7cv0ie0fx60mlb3rs6v3>

605 **S1 Text. Supplementary Information include: S1.** Implemented detector to extract plankton  
606 images from the acquired videos **S2.** Evaluation of purity with respect to the number of samples  
607 using the lensless microscope dataset **S3.** Example images from the considered datasets **S4.**  
608 Example images from the considered datasets **S5.** *Estimated number of clusters adopting the*  
609 *partition coefficient* **S6.** *Local Binary Pattern computation.* **S7.** *Multi-dimensional representation*  
610 *for the Haralick subset of features* **S8.** *Multi-dimensional representation for the Hu-moments*  
611 *subset of features* **S9.** *Multi-dimensional representation for the features extracted from the gray*  
612 *values histogram* **S10.** *Multi-dimensional representation for the LBP subset of features* **S11.**  
613 *Multi-dimensional representation for the Fourier Descriptors subset of features* **S12.** *Multi-*  
614 *dimensional representation for the Zernike moments subset of features* **S13.** *Histogram reporting*  
615 *the normalized ranking score for the set of designed descriptors* **S14.** Schematic work flow  
616 describing how an observation is associated to the three possible outputs of the developed system:  
617 retraining class, anomaly or belonging to a trained class

618 **S1 Fig. Implemented detector to extract plankton images from the acquired videos.** The  
619 bounding box corresponding to the final detected contour is used to crop the plankton image.

620 **S2 Fig. Evaluation of purity with respect to the number of samples using the lensless**  
621 **microscope dataset.** The results are very accurate with number of images per sample higher or

622 equal to 100. Using 50 images results in an overlap between two clusters (corresponding to the  
623 species *Paramecium bursaria* and *Blepharisma americanum*), and in a decrease of the  
624 performances (light gray bar). The corrected purity algorithm introduced in this supplement (see  
625 Customized purity algorithm section), allows for a more accurate result (patterned bar).

626 **S3 Fig. Example images from the considered datasets. a-z13** WHOI dataset (names as they are  
627 labeled in the dataset) **z14-z23** lensless microscope dataset. **a** *Ceratium* **b** *Chrysochromulina* **c**  
628 *Coscinodiscus* **d** *Dactyliosolen* **e** *Gyrodinium* **f** *Strombidium\_morphotype1* **g** *Dino30* **h** *Euglena*  
629 **i** *Eucampia* **j** *Flagellate\_sp3* **k** *Pyramimonas\_longicauda* **l** *Thalassionema* **m** *Delphineis* **n**  
630 *Pleurosigma* **o** *Chaetoceros\_didymus\_flagellate* **p** *Dictyocha* **q** *DactFragCerataul* **r**  
631 *Emiliana\_huxleyi* **s** *Corethron* **t** *Kiteflagellates* **u** *Tintinnid* **v** *Dinobryon* **w** *Ephemera* **x**  
632 *Thalassiosira\_dirty* **y** *Skeletonema* **z** *Pseudochattonella\_farcimen* **z0** *Proterothropsis\_sp* **z1**  
633 *Heterocapsa\_triquetra* **z2** *Rhizosolenia* **z3** *Prorocentrum* **z4** *Pleurosigma* **z5** *Phaeocystis* **z6** *Laboea*  
634 *Strobila* **z7** *Katodinium\_or\_Torodinium* **z8** *Mesodinium\_sp* **z9** *Paralia* **z10** *Guinardia\_striata* **z11**  
635 *Asterionellopsis* **z12** *Amphidinium\_sp* **z13** *Pennate\_morphotype1* **z14** *Blaepharisma Americanum*  
636 **z15** *Euplotes* *Eurystomus* **z16** *Spirostomum ambiguum* **z17** *Volvox* **z18** *Arcella Vulgaris* **z19**  
637 *Actinosphaerium Nucleofilum* **z20** *Dileptus* **z21** *Stentor Coeruleous* **z22** *Paramecium Bursaria* **z23**  
638 *Didinium nasutum*.

639 **S4 Fig. Examples of species that are incorrectly assigned to the same cluster by our algorithm**  
640 **because of their morphological similarity in our feature space.** Similarity is intended from left  
641 to right **a** *Proterothropsis\_sp* **b** *Heterocapsa\_triquetra* **c** *Amphidinium\_sp* **d**  
642 *Pseudochattonella\_farcimen* **e** *Gyrodinium* **f** *Prorocentrum*

643 **S5 Fig. Estimated number of clusters adopting the partition coefficient. a** and the XIE-BENI  
644 index **b** as a function of sample size (species). The results are less precise if compared with the

645 partition entropy (see fig 2e in the main text). However, both the algorithms can reconstruct  
646 correctly the number of clusters for subset of 3 species and 5 species. The number of clusters on  
647 the y axis is the distribution of ten runs on random subsets of all species. For example, for the  
648 leftmost box, 3 species have been randomly chosen from the lensless microscope database. This  
649 procedure is repeated ten times and the mode is then used as the estimated number of clusters.

650 **S6 Fig. Local Binary Pattern computation.**

651 **S7 Fig. Multi-dimensional representation for the Haralick subset of features. a** Andrew's  
652 curve. **b** Parallel coordinates

653 **S8 Fig. Parallel coordinate for the Hu-moments subset of features. a** Andrew's curve. **b**  
654 Parallel coordinates

655 **S9 Fig. Multi-dimensional representation for the features extracted from the gray values**  
656 **histogram. a** Andrew's curve. **b** Parallel coordinates

657 **S10 Fig. Multi-dimensional representation for the LBP subset of features. a** Andrew's curve.  
658 **b** Parallel coordinates

659 **S11 Fig. Multi-dimensional representation for the Fourier Descriptors subset of features. a**  
660 Andrew's curve. **b** Parallel coordinates

661 **S12 Fig. Multi-dimensional representation for the Zernike moments subset of features. a**  
662 Andrew's curve. **b** Parallel coordinates

663 **S13 Fig. Histogram reporting the normalized ranking score for the set of designed**  
664 **descriptors.**

665 **S14 Fig. Schematic work flow describing how an observation is associated to the three**  
666 **possible outputs of the developed system: retraining class, anomaly or belonging to a**  
667 **trained class**

668 **S1 Table. Computational time on raspberry pi for the analysis of one sample.** The standard  
669 deviation is computed among the objects contained into the 60 frames of the analyzed video.

670

## 671 **Acknowledgment**

672 We thank Amanda K. Paulson and Aleksandar Godjoski for critical reading of the manuscript.

673 We also thank all faculty and students in the National Science Foundation Center for Cellular

674 Construction for discussion and critical feedback on the general idea and pipeline.

## 675 **Author contribution**

676 **Conceptualization:** Vito Paolo Pastore, Simone Bianco.

677 **Data curation:** Vito Paolo Pastore, Simone Bianco.

678 **Funding acquisition:** Simone Bianco.

679 **Investigation:** Vito Paolo Pastore, Simone Bianco and Thomas Zimmerman.

680 **Methodology:** Vito Paolo Pastore, Sujoy K. Biswas, Thomas Zimmerman and Simone Bianco.

681 **Project administration:** Simone Bianco.

682 **Resources:** Simone Bianco.

683 **Software:** Vito Paolo Pastore

684 **Supervision:** Simone Bianco and Thomas Zimmerman.

685 **Validation:** Vito Paolo Pastore, Sujoy K. Biswas, Thomas Zimmerman and Simone Bianco.

686 **Visualization:** Vito Paolo Pastore

687 **Writing ± original draft:** Vito Paolo Pastore, Thomas Zimmerman and Simone Bianco.

688 **Writing ± review & editing:** Vito Paolo Pastore, Thomas Zimmerman and Simone Bianco

689

690

## 691 **REFERENCES**

- 692 [1] M. J. Behrenfeld *et al.*, “Biospheric primary production during an ENSO transition,”  
693 *Science*, vol. 291, no. 5513, pp. 2594–2597, Mar. 2001.
- 694 [2] A. Sournia, M.-J. Chrdtiennot-Dinet, and M. Ricard, “Marine phytoplankton: how many  
695 species in the world ocean?,” *J. Plankton Res.*, vol. 13, no. 5, pp. 1093–1099, Jan. 1991.
- 696 [3] A. J. Richardson *et al.*, “Using continuous plankton recorder data,” *Prog. Oceanogr.*, vol.  
697 68, no. 1, pp. 27–74, Jan. 2006.
- 698 [4] T. O. Fossum *et al.*, “Toward adaptive robotic sampling of phytoplankton in the coastal  
699 ocean,” *Sci. Robot.*, vol. 4, no. 27, p. eaav3041, Feb. 2019.
- 700 [5] T. G. Zimmerman and B. A. Smith, “Lensless Stereo Microscopic Imaging,” in *ACM*  
701 *SIGGRAPH 2007 Emerging Technologies*, New York, NY, USA, 2007.
- 702 [6] Heidi M. Sosik, Emily E. Peacock, Emily F. Brownlee, “Annotated Plankton Images - Data  
703 Set for Developing and Evaluating Classification Methods.”
- 704 [7] M. S. Schmid, C. Aubry, J. Grigor, and L. Fortier, “The LOKI underwater imaging system  
705 and an automatic identification model for the detection of zooplankton taxa in the Arctic  
706 Ocean,” *Comput. Vis. Oceanogr.*, vol. 15–16, pp. 129–160, Apr. 2016.
- 707 [8] Culverhouse, P. F., Ellis, R. E., Simpson, R. G., Williams, R., Pierce, R. W., Turner, J. T.,  
708 “Categorisation of five species of *Cymatocylis* (Tintinidae) by artificial neural network,”  
709 *1994*, pp. 107:273–280.
- 710 [9] E. C. Orenstein and O. Beijbom, “Transfer Learning and Deep Feature Extraction for  
711 Planktonic Image Data Sets,” in *2017 IEEE Winter Conference on Applications of*  
712 *Computer Vision (WACV)*, 2017, pp. 1082–1088.
- 713 [10] Lumini, Alessandra & Nanni, Loris, “Deep learning and transfer learning features for  
714 plankton classification,” *2019*.
- 715 [11] Qiao Hu and Cabell Davis, “Automatic plankton image recognition with co-occurrence  
716 matrices and Support Vector Machine,” *Marine Ecology Progress Series*, vol. 295, pp. 21--  
717 31, 2005.
- 718 [12] M. C. B. | D. of Oceanography *et al.*, “RAPID: Research on Automated Plankton  
719 Identification,” *Oceanography*, vol. 20, Jun. 2007.

- 720 [13] Vito P. Pastore, Thomas Zimmerman, Sujoy K. Biswas, and Simone Bianco,  
721 “Establishing the baseline for using plankton as biosensor,” presented at the Proc.SPIE,  
722 2019, vol. 10881.
- 723 [14] Sujoy Kumar Biswas *et al.*, “High throughput analysis of plankton morphology and  
724 dynamic,” presented at the Proc.SPIE, 2019, vol. 10881.
- 725 [15] J. Dai, R. Wang, H. Zheng, G. Ji, and X. Qiao, “ZooplanktoNet: Deep convolutional  
726 network for zooplankton classification,” 2016, pp. 1–6.
- 727 [16] H. M. Sosik and R. J. Olson, “Automated taxonomic classification of phytoplankton  
728 sampled with imaging-in-flow cytometry,” *Limnol. Oceanogr. Methods*, vol. 5, no. 6, pp.  
729 204–216, 2007.
- 730 [17] M. B. Blaschko *et al.*, “Automatic In Situ Identification of Plankton,” in *2005 Seventh*  
731 *IEEE Workshops on Applications of Computer Vision (WACV/MOTION’05) - Volume 1*,  
732 2005, vol. 1, pp. 79–86.
- 733 [18] S. Dieleman, J. De Fauw, and K. Kavukcuoglu, “Exploiting Cyclic Symmetry in  
734 Convolutional Neural Networks,” *ArXiv E-Prints*, p. arXiv:1602.02660, Feb. 2016.
- 735 [19] H. Zheng, R. Wang, Z. Yu, N. Wang, Z. Gu, and B. Zheng, “Automatic plankton image  
736 classification combining multiple view features via multiple kernel learning,” *BMC*  
737 *Bioinformatics*, vol. 18, no. 16, p. 570, Dec. 2017.
- 738 [20] A. Hughes, J. D. Mornin, S. K. Biswas, D. P. Bauer, S. Bianco, and Z. J. Gartner,  
739 “Quantius: Generic, high-fidelity human annotation of scientific images at 105-clicks-per-  
740 hour,” *bioRxiv*, p. 164087, Jul. 2017.
- 741 [21] D. A. Reynolds, “Gaussian Mixture Models,” in *Encyclopedia of Biometrics*, 2009.
- 742 [22] A. Romero, C. Gatta, and G. Camps-Valls, “Unsupervised Deep Feature Extraction for  
743 Remote Sensing Image Classification,” *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 3,  
744 pp. 1349–1362, Mar. 2016.
- 745 [23] S. Haykin, *Neural Networks: A Comprehensive Foundation*, 1st ed. Upper Saddle River,  
746 NJ, USA: Prentice Hall PTR, 1994.
- 747 [24] M. H. Bhuyan, D. K. Bhattacharyya, and J. K. Kalita, “Network Anomaly Detection:  
748 Methods, Systems and Tools,” *IEEE Commun. Surv. Tutor.*, vol. 16, no. 1, pp. 303–336,  
749 First 2014.
- 750 [25] Thomas Zimmerman *et al.*, “Stereo in-line holographic digital microscope,” presented at  
751 the Proc.SPIE, 2019, vol. 10883.
- 752 [26] B. Grindstaff, M. E. Mabry, P. D. Blischak, M. Quinn, and J. C. Pires, “Affordable Remote  
753 Monitoring of Plant Growth and Facilities using Raspberry Pi Computers,” *bioRxiv*, p.  
754 586776, Jan. 2019.
- 755 [27] C. Scherer *et al.*, *The development of UK pelagic plankton indicators and targets for the*  
756 *MSFD*. 2015.
- 757 [28] Z. Huang and J. Leng, “Analysis of Hu’s moment invariants on image scaling and rotation,”  
758 *2010 2nd Int. Conf. Comput. Eng. Technol.*, vol. 7, pp. V7-476-V7-480, 2010.
- 759 [29] Z. Yang and T. Fang, “On the Accuracy of Image Normalization by Zernike Moments,”  
760 *Image Vis. Comput*, vol. 28, no. 3, pp. 403–413, Mar. 2010.
- 761 [30] T. K. Ho, “Random decision forests,” in *Document analysis and recognition, 1995.,*  
762 *proceedings of the third international conference on*, 1995, vol. 1, pp. 278–282.
- 763 [31] R. Genuer, J.-M. Poggi, and C. Tuleau, “Random Forests: some methodological insights,”  
764 *ArXiv08113619 Stat*, Nov. 2008.
- 765 [32] L. Breiman, “Random Forests,” *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, Oct. 2001.

- 766 [33] “Random forest algorithm for classification of multiwavelength data - IOPscience.”  
767 [Online]. Available: <http://iopscience.iop.org/article/10.1088/1674-4527/9/2/011>.  
768 [Accessed: 11-Nov-2018].
- 769 [34] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, “Estimating  
770 the support of a high-dimensional distribution,” *Neural Comput.*, vol. 13, no. 7, pp. 1443–  
771 1471, Jul. 2001.
- 772  
773