

# Automatic underwater fish species classification with limited data using few-shot learning

Sebastien Villon, Corina Iovan, Morgan Mangeas, Thomas Claverie, David Mouillot, Sébastien Villéger, Laurent Vigliola

### ▶ To cite this version:

Sebastien Villon, Corina Iovan, Morgan Mangeas, Thomas Claverie, David Mouillot, et al.. Automatic underwater fish species classification with limited data using few-shot learning. Ecological Informatics, 2021, 63, pp.101320. 10.1016/j.ecoinf.2021.101320. hal-03415715

# HAL Id: hal-03415715 https://hal.umontpellier.fr/hal-03415715

Submitted on 22 Nov 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Automatic underwater fish species 1 classification with limited data using 2 few-shot learning 3 4 Sébastien Villon<sup>a,\*</sup>, Corina Iovan<sup>a</sup>, Morgan Mangeas<sup>a</sup>, Thomas Claverie<sup>b,c</sup>, David 5 Mouillot<sup>c,d</sup>, Sébastien Villéger<sup>c</sup>, Laurent Vigliola<sup>a</sup> 6 7 8 <sup>a</sup> ENTROPIE, IRD, University of New-Caledonia, University of La Reunion, CNRS, Ifremer, Labex Corail, Noumea, New-Caledonia, France 9 <sup>b</sup> CUFR Mayotte, France 10 <sup>c</sup> MARBEC, University of Montpellier, CNRS, IRD, Ifremer, Montpellier, 11 France 12 <sup>d</sup> Institut Universitaire de France, Paris, France 13 14 \* Corresponding author: sebastien.villon@ird.fr 15 16 17

21 Abstract:

Underwater cameras are widely used to monitor marine biodiversity, and 22 the trend is increasing due to the availability of cheap action cameras. 23 The main bottleneck of video methods now resides in the manual 24 processing of images, a time-consuming task requiring trained experts. 25 Recently, several solutions based on Deep Learning (DL) have been 26 proposed to automatically process underwater videos. The main limitation 27 of such algorithms is that they require thousands of annotated images in 28 order to learn to discriminate classes (here species). This limitation 29 implies two issues: 1) the annotation of hundreds of common species 30 requires a lot of efforts 2) many species are too rare to gather enough 31 data to train a classic DL algorithm. Here, we propose to explore how 32 few-shot learning (FSL), an emerging research field, could overcome DL 33 limitations. Few-shot learning is based on the principle of training a Deep 34 Learning algorithm on "how to learn a new classification problem with 35 only few images". In our case-study, we assess the robustness of FSL to 36 discriminate 20 coral reef fish species with a range of training databases 37 from 1 image per class to 30 images per class, and compare FSL to a 38 classic DL approach with thousands of images per class. We found that 39 FSL outperform classic DL approach in situations where annotated images 40 are limited, yet still providing good classification accuracy. 41

42

43 Keywords: few-shot learning, Deep learning, video, marine biodiversity

44

19

## 45 Introduction

46

The world's ecosystems have entered an era of anthropogenic 47 48 defaunation where human activities have triggered global decline in animal abundance, species range contraction and a new wave of species 49 extinction [1]. This global change is threatening ecosystem services 50 worldwide hence the stability of our food systems, economies, and health. 51 Defaunation is more advanced in terrestrial and freshwater ecosystems 52 than in the marine environment where it started centuries later. However, 53 the pace of defaunation is accelerating in oceans mostly due to the 54 advent of industrial fishing since a century ago [2]. Given this context of 55 global changes rapidly affecting fish communities, it is imperative to 56 monitor fish biodiversity over time, on a large scale and using non-57 destructive methods. 58

Fish biodiversity surveys in the marine environment are typically 59 performed by divers. Although dive visual censuses provide a great deal 60 of information on some shallow habitats, there are many limitations. First, 61 divers are limited by depth and can hardly perform long dives to count 62 fish below 30 m, ignoring mesophotic habitats and deeper ecosystems. 63 Second, divers are limited by time and generally focus their 2-4 dives per 64 day in the most speciose hard-substrate habitats, and ignore less rich and 65 often immense adjacent soft-bottom habitats. Third, dive surveys provide 66 data at a slow rate so that the compilation of global fish biodiversity 67 database takes decades of efforts by multiple teams of highly skilled 68 taxonomic divers (e.g. [3], [4]). This is a major restriction to the 69 necessary temporal monitoring of global marine ecosystems, although a 70 few time series exist in some countries<sup>1</sup> [5]. 71

<sup>&</sup>lt;sup>1</sup> AIMS, Long-Term Monitoring Program: Visual Census Fish Data (Great Barrier Reef) https://apps.aims.gov.au/metadata/view/5be0b340-4ade-11dc-8f56-00008a07204e

Underwater videos (UV) are increasingly used [6] to overcome the 72 limitations of diver-based surveys to quickly collect large amounts of 73 data. For instance, more than 15,000 video stations were deployed in 58 74 countries in just three years for the first global assessment of the 75 conservation status of reef sharks [7]. Furthermore, underwater video 76 surveys can be performed in many habitats, with some example in 77 shallow reefs [8], sandy lagoons [9], deep sea [10], and even in the 78 pelagic ecosystem [11]. Deploying underwater video stations does not 79 require expert taxonomists and is now quite inexpensive with the 80 improvement of cheap action cameras since a few years. The bottleneck 81 to analyse this data now resides in the manual processing of the videos. 82 Indeed, manually extracting fish biodiversity and abundance data from 83 raw videos requires unsustainable workload by highly trained taxonomic 84 experts. Although this annotation work can be improved through citizen 85 science [12]-[14], such time-consuming and expensive task cannot 86 match the increasing size of datasets, up to 20,000 hours of videos for 87 global surveys [7] and the necessary monitoring of global oceans over 88 time. 89

90

As the demand for automatic methods to analyze underwater videos is 91 rising, the latest generation of deep learning algorithms (DL), and in 92 particular convolutional neural networks (CNNs) are increasingly used for 93 species identification [14]–[17] and fish detection [18]–[21]. However, 94 these algorithms require a large dataset of annotated images (thumbnails 95 hereafter) in order to train a robust model, able to provide satisfying 96 results. Therefore, this method still requires collecting an important image 97 dataset manually annotated by experts. This is especially problematic in 98 highly diverse faunas such as coral reef fish that encompass nearly 6500 99 species worldwide [22]. Furthermore, a universal pattern in species 100 distribution, including fish communities, is that both rare and common 101 species are found in every community, with the fraction of rare species 102

more important in rich ecosystems, such as coral reefs [23], [24]. It is
therefore almost impossible to gather enough thumbnails of rare species
to efficiently train a deep neural network in a "classic" way, which

requires thousands of images per species [25]–[27]

107

108

There are two ways to tackle this problem of lack of data. The first one 109 consists of directly addressing the data itself, through data augmentation 110 [28]–[30]. The second option is to change the classification algorithm. 111 Few-shots learning (FSL) algorithms [31], [32] are designed to compute a 112 classification task (query, noted Q) with only a few thumbnails to train 113 (Support Sets, noted SS), and it has been increasingly studied since 2017 114 [33]. Few-shots learning methods are divided into three main 115 approaches. Metric-based methods are embedding both queries (Q) and 116 support sets (Ss), before assigning to the query a class, according to 117 distances computed between Q and Ss ([34]–[36]). The second approach 118 consists of 1) training a model on a large database, and 2) adapt this 119 model to a new task with few examples, while not forgetting the concepts 120 learned previously [37], [38]. Finally, optimization-based methods are 121 designed to adapt quickly to new tasks, hence able to learn a 122 classification task with few examples [33], [39], [40]. Optimization-based 123 algorithms showed promising results in deep learning few-shot 124 classification [33], [41], [42]. Such methods propose to pre-train (or 125 "meta-train") a model with existing databases (e.g. MiniImageNet [43], 126 Ominglot [44]) on different tasks so it can adapt easily to a new one. For 127 object identification, a task is defined by the classes the model has to 128 discriminate. Once this model, called "meta-model" has been trained, it 129 can then be tuned to operate on a new task with a very limited dataset, 130 131 usually only 1-5 thumbnails per class.

In this study we propose to compare the efficiency of optimization-based 133 few-shot learning and standard large dataset deep-learning methods to 134 identify coral reef fish species on images. More specifically, we aim to 135 determine how well a classic deep learning architecture trained with 136 thousands of images and the benefit of data augmentation (hereafter DL) 137 and FSL algorithms perform in situations where training thumbnail 138 dataset is large or limited. To achieve this, we first trained a classic DL 139 architecture built for image classification [45] on a large dataset of 140 69,169 thumbnails, and on a more limited dataset of 6,320 thumbnails for 141 20 coral reef fish species. Then, we trained a few-shots, optimization-142 based learning algorithm [39] on the exact same training datasets while 143 varying the number of shots from 1 to 30. Finally, we compared the 144 capacity of DL and FSL models to correctly identify species on an 145 independent thumbnail dataset, and modelled the asymptotic relationship 146 between classification accuracy and the number of thumbnails in the 147 training datasets for both classic DL and FSL algorithms 148

149

#### 151 Material and methods

152

#### 153 Thumbnail datasets

We used three fish thumbnail datasets (*T0*, *T1*, and *T2*) extracted from 175 underwater videos recorded on reefs around Mayotte Island (Western Indian Ocean) using GoPro Hero 3+ and GoPro hero 4+ cameras with a resolution of 1920x1080 pixels. A thumbnail is defined as an image containing a single labelled fish belonging to one of the 20 most common fish species in the videos, and representing a broad range of sizes, colors, body orientations, and background (Supp. Fig. 1, Supp. Fig. 2).

*TO* is composed of 69,169 thumbnails extracted from 130 videos, with a

range of 1,134 to 7,345 thumbnails per species (Table 1). T1 is composed

of 6,320 thumbnails extracted from 20 videos with 40-1,436 images per
species whereas *T2* is composed of 13,232 thumbnails extracted from 25
videos with 55-3,896 images per species. Thumbnails size originally

- ranged from 55x55 pixels to 500x450 pixels, but were resized to 84x84
   pixels before being processed through FS and DL algorithms.
- 168 The datasets *T1* and *T2* correspond to two real scenarii where videos were 169 recorded during two trips in the field of a week each.

The three thumbnails datasets are fully independent, as they were
extracted from videos recorded at different sites, with different conditions
(weather, lighting, depth, time of the day, seascape) and on different
days.

174

175 To train our DL architecture, we applied data augmentation to *TO* and *T1*.

176 For each natural thumbnails in *TO* and *T1*, we created 9 thumbnails

177 through contrast augmentation or diminution, and horizontal flip. We then

obtained augmented datasets composed of 691,690 (AT0) and 63,200

- (AT1) images respectively Supp. Table 1. Further details on thumbnaildatasets and data augmentation are given in [46].
- 181
- 182 Table 1: Number of natural thumbnails extracted from the videos to build
- 183 our three datasets

Family	<sup>-</sup> amily Species		Training dataset <i>T1</i>	Test dataset <i>T2</i>
Acanthuridae	Acanthurus leucosternon	3,259	235	491
Acanthuridae	Acanthurus lineatus	1,008	114	864
Acanthuridae	Naso brevirostris	1,134	539	1932
Acanthuridae	Naso elegans	7,345	1,435	3,896
Acanthuridae	Zebrasoma scopas	4,970	48	579
Chaetodontidae	Chaetodon auriga	2,134	737	502
Chaetodontidae	Chaetodon guttatissimus	1,182	221	68
Chaetodontidae	Chaetodon trifascialis	5,234	41	630
Chaetodontidae	Chaetodon trifasciatus	4,421	71	82
Labridae	Gomphosus caeruleus	3,131	57	173
Labridae	Halichoeres	3,192	40	287

#### hortulanus

Labridae	Thalassoma hardwicke	4,951	181	275
Lethrinidae	Monotaxis grandoculis	3,893	797	1,422
Monacanthidae	Oxymonacanthus Iongirostris	2,553	54	55
Pomacentridae	Abudefduf vaigiensis	5,124	376	216
Pomacentridae	Amblyglyphidodon indicus	1,188	636	1,310
Pomacentridae	Chromis opercularis	1,525	81	93
Pomacentridae	Chromis ternatensis	3,640	300	156
Pomacentridae	Pomacentrus sulfureus	5,409	270	142
Zanclidae	Zanclus cornutus	3,876	86	59
TOTAL		69,169	6,320	13,232

184

185

186

187

188 Experimental design

189 To compare classic deep-learning and few-shot algorithms in situations of 190 large or small thumbnail datasets, we led five experiments using datasets

191 T0, T1, T2, AT0 and AT1 described in Supp. Table 1 :

- We trained a classic DL algorithms architecture with our biggest
   dataset *ATO* as a baseline for the DL accuracy;
- 194 2) We trained the same DL architecture with the same hyper-
- 195 parameters (e.g. model architecture and training process) but on a
- 196 much more limited dataset (*AT1*). Hyper-parameters are the
- 197 parameters defining the architecture (number of layers, number
- and size of convolutions, connections between layers) and the
  training process of a Deep Model (learning rate, neurone activation,
  back-propagation compotation).;
- 3) We trained the same DL architecture with limited datasets obtained
  by subsampling T0 to 250 and 500 images per class (here after
  "species" when we are referring to our experiments), corresponding
  to 2500 and 5000 thumbnails in *ATO*;
- 4) We pre-trained a FSL architecture on the 64 training classes of
  MiniImageNet (Supp. Fig. 3) and used T0 to build support sets (*SS*)
  with 1, 5, 15 and 30 thumbnails for each fish species;
- 5) We pre-trained the same FSL architecture on MiniImageNet and
  used the more limited T1 dataset to build support sets with 1, 5, 15
  and 30 images per species.

We used ResNet 100 [45] as our classic deep-learning algorithm. Resnet 212 is a convolutional neural network (CNN), a DL architecture which is able 213 to both extract features from images and classify these images thanks to 214 those features [47]. In order for a CNN to build an image classification 215 model, the architecture is fed a large dataset, composed of pairs of labels 216 and images. Using this dataset, the algorithms change their inner 217 parameters in order to minimize the classification error, through a 218 process called back-propagation. The ResNet architecture achieved the 219 best results on ImageNet Large Scale Visual Recognition Competition 220 (ILSVRC [43]) in 2015, considered the most challenging image 221

classification competition. It is still one of the best classificationalgorithms, while being easy to use and implement.

224

For the few-shot implementation, we used the Reptile algorithm [39]. 225 Few-shot learning algorithms are specific DL algorithms, whose goal is to 226 be able to fit a model with very few training images. The Reptile algorithm 227 is based on the well-known MAML architecture [33], and more precisely 228 on the first-order version of MAML [48]. The Reptile algorithm is based on 229 the division of the training dataset into a number of tasks  $T_{i}$ , a task being 230 a learning problem. Through repetitively changing the task during the first 231 232 training phase (known as meta-training), this algorithm produces a quick learner, i.e. a learner than can quickly adapt to a new task with a small 233 number of examples. 234

Here, the few-shot algorithms were tested on a classic *n*-ways *k*-shots 235 procedure, *n* being the number of classes per support set, and *k* the 236 number of images per class in the support set. For instance, a 5-ways 1-237 shots consists of training 5 classes with supports sets composed of 1 238 image per classes (e.g. species). We set n=5 [34], [36], [40], [42], [49] 239 and allowed k to vary between 1 shot and 30 shots for both experiments 240 4 and 5. We did not use data augmentation for FSL experiments for 241 several reasons. First, the goal of FSL is to adapt quickly with a very 242 limited number of images. Second, to have similar settings for method 243 comparison. There were no data-augmentation in the original paper, so 244 we reproduced that. It also allowed us to compare our results with those 245 obtained on benchmarks. Third, the reason behind the use of raw data 246 instead of augmented data in few-shot learning paper is that with very 247 few training samples and few conditions, the risk of overfitting by using 248 the same image modified multiple times is far greater than in classic 249 approaches with important datasets with many conditions. 250

253 Model comparison

All the DL and FSL models were tested on the independent *T2* dataset.

<sup>255</sup> First, we compared the results of experiments 1 and 2 in order to

estimate the decrease in performance of a classic ResNet DL architecture

when trained on a large dataset *ATO* (i.e. between 11,340 and 73,450

images per species after data augmentation, with an average of 3458

natural thumbnails per species) or trained on a more limited dataset AT1

260 (i.e. between 400 and 14,360 images per species after data

augmentation, with an average of 315 natural thumbnails per species).

Second, we compared the results of experiments 1 and 4 in order to
evaluate if the ResNet architecture outperforms the Reptile architecture in
a real-case situation where thumbnail dataset is large (*TO and ATO*).

Finally, we compared the results of experiments 2 and 5 to determine whether and to which extent a Reptile model performs better than a ResNet model in a real-case situation where thumbnail dataset is limited (*T1 and AT1*).

269

270

In order to better evaluate the performance of ResNet and Reptile
algorithms, we also modelled the relationship between model accuracy
and the number of thumbnails used to train the models. To achieve this,
we fitted the following asymptotic function to the results of experiments
1, 3 and 4 (obtained through training DL and FSL architectures on
datasets of various size obtained from ATO and TO):

277

278  $Accuracy = Accuracy_{\infty} \cdot (1 - exp(-R \cdot N_{image}))$  (eq.1)

where  $Accuracy_{\infty}$  is the asymptotic model accuracy when the number of thumbnails  $N_{image}$  is infinite, and R is the rate at which the asymptote is reached.

Equation 1 was fitted by non-linear mixed-effect modelling (NLME [50]) 282 using species as a random effect. This method is widely used for fitting 283 asymptotic processes. It allows estimating and comparing asymptotic 284 accuracies of both FSL and DL algorithms, and the number of image to 285 reach these asymptotic accuracies. The number of images required to 286 reach the asymptotic accuracy was calculated as the number of images 287 corresponding to an accuracy of 0.99 times the asymptotic value, 288 meaning the asymptote was reached within 1%. 289

- 290
- 291

#### 292 Results

The deep ResNet model trained on the large ATO dataset (3458 natural 293 thumbnails in average per species) during the first experiment obtained a 294 mean accuracy (i.e. percentage of correct classification) of 78.00% 295 (standard deviation (SD) of 15.16%) on T2 test-dataset (Table 2). With 296 this model, accuracy varied among species between 54.14% (Naso 297 brevirostris) and 99.07% (Abudefduf vagiensis). The same ResNet DL 298 model trained on smaller AT1 (315 natural thumbnails in average per 299 species) during the second experiment showed highly degraded 300 performance with a mean accuracy of only 42.21% (SD=24.95%). Among 301 species variation ranged with this model from only 3.49% (Chaetodon 302 trifascialis) to 85.86% (Chaetodon auriga). 303

The few-shot Reptile architecture trained on limited T1 dataset during our fifth experiment obtained a mean accuracy of 32.04% for the 1-shot learning (SD= 12.70%) and 51.77% mean accuracy for the 30-shots learning (SD= 18.96%) (Table 2,). In this scenario of limited *T1* training dataset, the few-shot Reptile algorithm nearly equalled the ResNet DL model with only 5 shots (41.47% accuracy for 5-shots learning on T1 vs 42.21% for DL on A*T1*), and performed better beyond 10 shots (45.92%

of accuracy on *T2* with 10-shots learning). A pairwise proportion test 311 showed a p-value <0.0001, assessing that FSL was significantly better 312 than DL in this scenario beyond 10 shots (Supp. Table 3) accuracy of 313 Reptile models had a standard deviation from 12.70% with one-shot 314 learning, to 18.96% with 30-shots learning, indicating important variation 315 in accuracy among species. However, this standard deviation was smaller 316 than that of the ResNet algorithm trained on the same AT1 limited 317 dataset (24.95%). 318

319 Figure 1: Relationship between the number of natural thumbnails per species for training and the accuracy of deep-learning and few-shot learning models. Non-320 linear mixed effects asymptotic model fit for (a) DL architecture at the fixed-321 322 effect level, and for FSL architecture at (b) the fixed-effect level and (c) the random-effect level. Grey areas represent 95% CI in fixed-effects estimates. 323 324 Dotted lines represent the NLME estimate of the number of images per species required to reach the 99% asymptote value. We obtained similar magnitude with 325 326 the 95% asymptote value, reached with 750 images for DL and 5 with FSL.



Number of images

- 329
- 330

The same few-shot Reptile architecture trained on subsets of T0 during the fourth experiment obtained even better results than when trained on T1, with a mean accuracy on T2 of 34.57% for 1-shot, 50.23% for 5shots, and up to 64.92% for 30-shots (Table 2).

Mixed-effects modelling (NLME) of T0 and AT0 experimental data showed a clear pattern of asymptotic increase of accuracy with the number of natural thumbnails for both Resnet and Reptile architectures (Figure 1).

NLME models included significant species random effect for both DL and
FSL (Log-likelihood tests, P<0.0001).</li>

The fixed-effect asymptotic value of accuracy was higher for ResNet 340 model ( $Accuracy_{\infty}$  = 77.34%, 95% CI: 71.26-83.41%) than for Reptile 341 model (*Accuracy*<sub>∞</sub> =60.87%, 95% CI: 54.48–67.26%), illustrating higher 342 classification power of ResNet over Reptile when large numbers of 343 thumbnails are available. However, the slope of the asymptotic model 344 was two-orders of magnitude higher for Reptile (0.707, 95% CI: 0.559-345 0.854) than for ResNet architecture (0.0040, 95% CI: 0.0032-0.0048), 346 illustrating the high capacity of Reptile FSL algorithm to learn from only a 347 few images. NLME modelling further showed that average asymptotic 348 accuracy was reached with only 7 natural thumbnails per species for 349 Reptile architecture, compared to 1153 natural thumbnails per species for 350 ResNet, confirming the strong power of Reptile method in situation of 351 limited thumbnail training dataset. However, model random effects 352 showed that some variation existed among species. For Reptile 353 architecture, asymptotic accuracy values ranged from 38.09% 354 (Amblyglyphidodon indicus) to 89.78% (Pomacentrus sulfureus), and was 355 reached with 4 to 16 training images per species. For DL architecture, 356 species asymptotes varied from 62.72% (Monotaxis grandoculis) to 357

96.81% (*Abudefduf vaigiensis*), and could be reached with 786-1776
thumbnails per species (*Supp. Table 4*).

- Table 2: Accuracy of our ResNet deep-learning (DL) and Reptile few-shots learning (FSL) models trained on T0 or
- *T1 thumbnails datasets for different number of shots. Accuracy is the % correct classification of models on T2 test*
- *dataset. DL models were trained from TO and T1 after data augmentation (ATO and AT1).*

	D	ρL			FS	SL		
	то	T1		T1			то	
Image per species (on average)	3458	315	1 shot	5 shots	30 shots	1 shot	5 shots	30 shots
Abudefduf vaigiensis	99.07	69.91	16.08	11.39	11.38	47.67	70.9	86.35
Acanthurus leucosternon	86.15	44.67	25.51	30.71	38.74	19.23	28.8	42.66
Acanthurus lineatus	59.72	20.37	39.86	56.04	72.50	32.93	61.01	72.02
Amblyglyphidodon indicus	58.78	60.78	25.75	26.74	32.86	28.26	32.55	40.64
Chaetodon auriga	87.05	85.86	18.16	25.68	36.56	27.8	35.18	53.20
Chaetodon guttatissimus	85.50	44.12	33.58	44.21	58.26	29.61	51.18	79.29
Chaetodon trifascialis	90.00	3.49	29.02	25.48	28.44	27.14	43.17	63.51
Chaetodon trifasciatus	87.80	28.05	38.73	50.63	66.72	32.41	51.07	70.63

Chromis opercularis	61.29	9.68	44.01	61.81	62.94	45.34	68.28	81.50
Chromis ternatensis	59.61	55.77	18.91	24.94	35.07	35.4	55.44	67.22
Gomphosus caeruleus	75.72	20.81	26.01	38.99	58.74	28.96	39.16	54.22
Halichoeres hortulanus	82.93	17.07	31.82	44.81	57.01	28.94	41.35	58.87
Monotaxis grandoculis	57.10	53.37	32.03	41.52	50.64	32.8	45.85	59.13
Naso brevirostris	54.14	68.60	47.06	54.26	64.66	54.47	58.08	61.00
Naso elegans	93.24	79.43	34.54	43.11	52.17	28.47	33.71	44.36
Oxymonacanthus	96.43	14.54	39.29	53.48		42.15	65.44	
longirostris					66.26			84.86
Pomacentrus sulfureus	90.14	61.97	70.90	88.18	93.93	65.53	86.21	90.00
Thalassoma hardwicke	90.90	51.64	25.64	44.4	67.96	26.72	45.7	70.60
Zanclus cornutus	81.36	40.68	18.33	31.28	44.44	26.08	40.82	62.72
Zebrasoma scopas	63.04	13.30	25.56	31.81	36.05	31.42	50.69	55.70
MEAN	78.00	42.21	32.04	41.47	51.77	34.57	50.23	64.92
SD	15.16	24.95	12.70	16.93	18.96	11.14	14.75	14.55

## 366

#### 367 Discussion

Our experiments demonstrated that few-shot learning methods based on Reptile 368 architecture can be effectively used to drastically reduce the number of annotated 369 370 images for underwater fish identification. Accuracy levels obtained with few-shot learning algorithm trained with only five training images are close to those of a 371 372 standard Deep Learning architecture such as ResNet trained with 400-14350 images per species. Further, FSL architecture trained with 10 images outperformed a 373 ResNet 100 architecture trained with at least 400 images per species. This is a very 374 375 promising result in situations where many species need to be identified from models trained with a few images, a typical characteristic in marine biodiversity 376 377 applications.

However, the important standard deviation among the different trained species 378 379 (18.96 SD on 30-shots) showed that few-shot algorithms may not be robust enough 380 to discriminate among similar species showing only subtle differences. Nevertheless, in our 2<sup>nd</sup> experiment, our ResNet model achieved an accuracy under 40% for all 381 382 the species with fewer training images than 1140 (after data augmentation, i.e. 114 383 natural images), and only 7 species were identified with an accuracy greater than 384 45%. These species were represented with a range of 2700-14350 images during the training phase. We also show better results with the model trained on T0 than 385 386 the model trained on T1. As expected, increasing the number of images per shot 387 rely on better performances as well as increasing the per species images variability. However, in real conditions, few-shot learning is to be used in a context where very 388 389 few images per classes are at disposal. Therefore, the dataset T1 corresponded 390 more to a real use case scenario.

Thus, there is a trade-off to make between accuracy and robustness on one hand, and the cost of video annotation by experts on the other.

393

Modeling the accuracy of neural networks using NLME allowed to understand the number of images per species required for the Few-shot and Deep architecture to reach 99% of their maximum potential accuracy. In our case study, there was a 150-fold factor between the average number of images required for a Deep Learning architecture (1153 images) and for a Few-shot architecture (7) to reach asymptotic accuracy. However, it is important to note that these numbers could vary according to the number and complexity of classes fed to the deep classifier.

In this work we used a Reptile FSL architecture. As the field of few-shot learning is quickly improving, new methods are proposed at a fast rate. While Reptile obtained a mean accuracy of 61.98% on the MiniImageNet dataset (the most used benchmark for few-shot learning methods) through a 5-shots learning, [51] recently achieved 80.51% of accuracy on the same dataset. Although further studies are required, we can reasonably assume that the improvements of FSL algorithms will further expand the possible use of few-shot learning for real-life use cases.

408

Applied to marine and coral reef ecology, such methods requiring few examples to 409 410 fit a model on an identification task could be used for studies on species rarely seen on screen. A key characteristic of highly diverse ecosystems is that they are 411 composed of few very common species and a large proportion of less-common and 412 rare species. Hence, the important effort required to build databases with a 413 414 sufficient number of images of all these rare species is the main bottleneck preventing the use of Deep Learning on a large number of species. The 415 improvement of few-shot learning algorithms offers promises to build efficient 416 identification models to automatically process images and videos to localise and 417 identify rare fish species. Such models could then be paired with more classic deep 418 419 architectures, more efficient to identify abundant species with the leverage of 420 important datasets.

421

#### 422 Acknowledgements:

- This study was funded by the French National Research Agency project ANR 18-CE02-0016 SEAMOUNTS.

#### 430 References

- 431 [1] R. Dirzo, H. S. Young, M. Galetti, G. Ceballos, N. J. B. Isaac, and B. Collen,
  432 "Defaunation in the Anthropocene," *Science (80-. ).*, vol. 345, no. 6195, pp.
  433 401–406, 2014.
- H. S. Young, D. J. Mccauley, M. Galetti, and R. Dirzo, "Patterns, Causes, and
  Consequences of Anthropocene Defaunation," *Annu. Rev. Ecol. Evol. Syst.*, no.
  August, pp. 333–358, 2016.
- J. E. Cinner *et al.*, "Meeting fisheries, ecosystem function, and biodiversity
  goals in a human-dominated world," *Science (80-. ).*, vol. 311, no. April, pp.
  307–311, 2020.
- R. D. Stuart-smith *et al.*, "Integrating abundance and functional traits reveals
  new global hotspots of fish diversity," *Nature*, vol. 501, no. 7468, pp. 539–
  542, 2013.
- 443 [5] A. Heenan *et al.*, "Long-term monitoring of coral reef fish assemblages in the
  444 Western central pacific," *Sci. Data*, vol. 4, pp. 1–12, 2017.
- [6] S. K. Whitmarsh, P. G. Fairweather, and C. Huveneers, "What is Big BRUVver
  up to? Methods and uses of baited underwater video," *Rev. Fish Biol. Fish.*,
  vol. 27, no. 1, pp. 53–73, 2017.
- 448 [7] M. Aaron MacNeil *et al.*, "Global status and conservation potential of reef
  449 sharks," *Nature*, no. July 2019, 2020.
- J. Juhel, L. Vigliola, L. Wantiez, T. B. Letessie, J. J. Meeuwig, and D. Mouillot,
  "Isolation and no-entry marine reserves mitigate anthropogenic impacts on
  grey reef shark behavior," *Sci. Rep.*, vol. 9, no. November 2018, pp. 1–11,
  2019.
- M. Cappo, G. De, and P. Speare, "Inter-reef vertebrate communities of the
  Great Barrier Reef Marine Park determined by baited remote underwater video
  stations," *Mar. Ecol. Prog. Ser.*, vol. 350, pp. 209–221, 2007.

- [10] V. Zintzen, M. J. Anderson, C. D. Roberts, E. S. Harvey, and L. Andrew,
  "Effects of latitude and depth on the beta diversity of New Zealand fish
  communities," *Sci. Rep.*, vol. 7, no. July, pp. 1–10, 2017.
- [11] Tom B Letessier *et al.*, "Remote reefs and seamounts are the last refuges for
  marine predators across the Indo- Pacific," *PLoS Biol.*, vol. 17, pp. 1–20,
  2019.
- 463 [12] C. J. Torney *et al.*, "A comparison of deep learning and citizen science
  464 techniques for counting wildlife in aerial survey images," *Methods Ecol. Evol.*,
  465 vol. 10, no. October 2018, pp. 779–787, 2019.
- 466 [13] E. C. Mcclure *et al.*, "Artificial Intelligence Meets Citizen Science to
  467 Supercharge Ecological Monitoring," *Patterns*, vol. 1, no. 7, p. 100109, 2020.
- 468 [14] M. Willi *et al.*, "Identifying animal species in camera trap images using deep
  469 learning and citizen science," *Methods Ecol. Evol.*, vol. 10, no. 1, pp. 80–91,
  470 2019.
- 471 [15] Z. Miao *et al.*, "Insights and approaches using deep learning to classify
  472 wildlife," *Sci. Rep.*, no. May, pp. 1–9, 2019.

[16] M. Lasseck, "Audio-based Bird Species Identification with Deep Convolutional
 Neural Networks Audio-based Bird Species Identification with Deep
 Convolutional Neural Networks," no. January, 2020.

- 476 [17] Y. Shiu *et al.*, "Deep neural networks for automated detection of marine
  477 mammal species," pp. 1–12, 2020.
- [18] D. Rathi, S. Jain, and S. Indu, "Underwater Fish Species Classification using
  Convolutional Neural Network and Deep Learning. (arXiv:1805.10106v1
  [cs.CV])," no. June, 2018.
- 481 [19] A. Salman *et al.*, "OCEANOGRAPHY : METHODS Fish species classification in
  482 unconstrained underwater environments based on deep learning," pp. 570–
  483 585, 2016.
- 484 [20] S. Villon *et al.*, "A Deep Learning algorithm for accurate and fast identification

- 485 of coral reef fishes in underwater videos," *PeerJ Prepr.*, vol. 6, p. e26818v1,
  486 2018.
- 487 [21] H. Qin, X. Li, J. Liang, Y. Peng, and C. Zhang, "DeepFish: Accurate underwater
  488 live fish recognition with a deep architecture," *Neurocomputing*, vol. 187, pp.
  489 49–58, 2016.
- 490 [22] P. Chabanet, S. R. Floeter, A. Friedlander, J. Mcpherson, and R. E. Myers,
  491 "Global Biogeography of Reef Fishes : A Hierarchical Quantitative Delineation
  492 of Regions," *PLoS One*, vol. 8, no. 12, 2013.
- [23] A. P. Hercos, M. Sobansky, H. L. Queiroz, A. E. Magurran, and A. Andre, "Local
  and regional rarity in a diverse tropical fish assemblage," *Biol. Sci.*, vol. 280,
  pp. 81–101, 2013.
- 496 [24] G. E. Jones, M. J. Caley, and P. L. Munday, "Rarity in Coral Reef Fish
  497 Communities," *Coral reef fishes Dyn. Divers. a complex Ecosyst.*, pp. 88–101,
  498 2002.
- L. Liu, T. Zhou, G. Long, J. Jiang, and C. Zhang, "Many-Class Few-Shot
  Learning on Multi-Granularity Class Hierarchy," *IEEE Trans. Knowl. Data Eng.*,
  pp. 1–14, 2020.
- [26] A. Li, T. Luo, Z. Lu, T. Xiang, and L. Wang, "Large-Scale Few-Shot Learning:
  Knowledge Transfer With Class Hierarchy," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7212–
  7220.
- 506 [27] P. Zhuang, Y. Wang, and Y. Qiao, "WildFish : A Large Benchmark for Fish
  507 Recognition in the Wild," in *Proceedings of the 26th ACM international*508 *conference on Multimedia*, 2018, vol. 2, pp. 1301–1309.
- 509 [28] J. Wang and L. Perez, "The Effectiveness of Data Augmentation in Image
  510 Classification using Deep Learning," *arXiv Prepr. arXiv1712.04621.*, 2017.
- [29] D. A. Van Dyk and X. Meng, "The Art of Data Augmentation," *J. Comput. Graph. Stat.*, vol. 8600, no. 2001, pp. 1–50, 2012.

513 514 515 516	[30]	S. C. Wong, M. D. Mcdonnell, G. Adam, and S. Victor, "Understanding data augmentation for classification : when to warp?," in <i>2016 international conference on digital image computing: techniques and applications (DICTA)</i> , 2016, pp. 1–6.
517 518 519	[31]	L. Fei-fei, R. Fergus, S. Member, and P. Perona, "One-Shot Learning of Object Categories," <i>EEE Trans. pattern Anal. Mach. Intell.</i> , vol. 28, no. 4, pp. 594–611, 2006.
520 521 522	[32]	M. Fink, "Object Classification from a Single Example Utilizing Class Relevance Metrics," in <i>Advances in neural information processing systems</i> , 2005, pp. 449–456.
523 524	[33]	C. Finn, P. Abbeel, and S. Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks," <i>arXiv Prepr. arXiv1703.03400</i> , 2017.
525 526 527 528	[34]	F. Sung, Y. Yang, and L. Zhang, "Learning to Compare : Relation Network for Few-Shot Learning Queen Mary University of London," in <i>Proceedings of the</i> <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , 2018, pp. 1199–1208.
529 530 531	[35]	Liu Yanbin <i>et al.</i> , "Learning to proagate labels: Transductive propagation network for few-shot learning," in <i>arXiv preprint arXiv:1805.10002</i> , 2019, pp. 1–14.
532 533	[36]	J. Victor, Garcia Bruna, "FEW-SHOT LEARNING WITH GRAPH NEURAL NETWORKS," in <i>arXiv preprint arXiv:1711.04043, 2017.</i> , 2018, pp. 1–13.
534 535 536	[37]	S. Gidaris, P. Paristech, N. Komodakis, and P. Paristech, "Dynamic Few-Shot Visual Learning without Forgetting," in <i>Proceedings of the IEEE Conference on</i> <i>Computer Vision and Pattern Recognitio</i> , 2018, pp. 4367–4375.
537 538 539	[38]	B. Hariharan, R. Girshick, and F. Ai, "Low-shot Visual Recognition by Shrinking and Hallucinating Features," in <i>Proceedings of the IEEE International Conference on Computer Vision</i> , 2017, pp. 3018–3027.
540 541	[39]	A. Nichol and J. Schulman, "Reptile : a Scalable Metalearning Algorithm," <i>arXiv</i> <i>Prepr. arXiv1803.02999, 2018</i> , pp. 1–11, 2018.
		20

- [40] Q. Sun and Y. L. T. Chua, "Meta-Transfer Learning for Few-Shot Learning,"
   *Conf. Comput. Vis. Pattern Recognit.*, pp. 403–412, 2018.
- [41] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, "Generalizing from a Few
  Examples: A Survey on Few-Shot Learning arXiv : 1904 . 05046v2 [ cs . LG ]
  13 May 2019," 2019.
- 547 [42] M. A. Jamal and H. Cloud, "Task Agnostic Meta-Learning for Few-Shot
  548 Learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and*549 *Pattern Recognition (CVPR)*, 2019.
- [43] O. Russakovsky *et al.*, "ImageNet Large Scale Visual Recognition Challenge,"
   *Int. J. Comput. Vis.*, pp. 211–252, 2015.
- 552 [44] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "The Omniglot challenge :
  a 3-year progress report," *COBEHA*, vol. 29, pp. 97–104, 2019.
- [45] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image
  Recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [46] S. Villon, D. Mouillot, M. Chaumont, and G. Subsol, "A new method to control
  error rates in automated species identification with deep learning algorithms," *Sci. Rep.*, vol. 10, pp. 1–13, 2020.
- 560 [47] Y. Lecun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, pp. 436–444,
  561 2015.
- 562 [48] A. Nichol, J. Achiam, and J. Schulman, "On First-Order Meta-Learning
  563 Algorithms," *arXiv*, pp. 1–15, 2018.
- [49] Y. Wang, Q. Yao, and L. M. Ni, "Generalizing from a Few Examples : A Survey
  on Few-shot Generalizing from a Few Examples : A Survey on Few-shot," *ACM Comput. Surv.*, vol. 53, no. June, 2020.
- [50] J. Pinheiro and D. Bates, *Mixed-effects models in S and S-PLUS*. 2006.
- [51] H. Li, D. Eigen, S. Dodge, M. Zeiler, and X. Wang, "Finding Task-Relevant
   Features for Few-Shot Learning by Category Traversal," in *Proceedings of the*

- 570 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR),
- 571 2019, vol. 1.

# 576 Supplementary

		5	
Abudefduf	Acanthurus	Acanthurus	Amblyglyphidodon
vaigiensis	leucosternon	lineatus	indicus
Chaetodon auriga	Chaetodon	Chaetodon	Chaetodon
	guttatissimus	trifascialis	trifasciatus
Mit have			
Chromis	Chromis	Gomphosus	Halichoeres
opercularis	ternatensis	caeruleus	hortulanus

Monotaxis grandoculis	Naso brevirostris	Naso elegans	Oxymonacanthus Iongirostris
	THE STREET		
<i>Pomacentrus sulfureus</i>	Thalassoma hardwicke	Zanclus cornutus	Zebrasoma scopas

578 Supp. Fig. 1 : The 20 reef fish species considered in this study



- 585 Supp. Fig. 2: Diversity of individuals of the same species and of their
- *environments.*



588 Supp. Fig. 3: Examples of classes' images in MiniImageNet

## 589 Supp. Table 1: Dataset usage during our experiments

name         Building         Number of annotations         Usage           Building of support         Building of support	
Building of suppor	
	ts sets of 1,5,15 and 30
TO Human Annotation 69,169 images for our fo	urth experiment
Building of suppor	ts sets or 1, 5, 15 and
T1Human Annotation6,32030 images for our	fifth experiment
T2 Human Annotation 13,232 testing dataset	
ATO Data augmentation applied on TO 691,690 DL training for our	first experiment
AT1 Data augmentation applied on T1 63,200 DL training for our	second experiment

- 590
- 591
- 592

593 Supp. Table 2: Mean accuracy obtained with FSL models trained with

1,5,10,15,20,25 and 30 images per species. All the images used for the

supports set are from T1.

Number of thumbnails in the Support set								
Species	1	5	10	15	20	25	30	
Abudefduf vaigiensis	16.08	11.39	13.18	12.59	14.02	12.98	11.38	
Acanthurus								
leucosternon	25.51	30.71	34.90	36.68	37.16	36.95	38.74	
Acanthurus lineatus	39.86	56.04	63.52	66.21	70.02	70.93	72.50	

Amblyglyphidodon							
indicus	25.75	26.74	28.06	28.80	30.24	32.20	32.86
Chaetodo auriga	18.16	25.68	30.67	32.63	33.78	35.42	36.56
Chaetodon							
guttatissimus	33.58	44.21	47.40	49.26	54.12	55.91	58.26
Chaetodon trifascialis	29.02	25.48	27.44	28.02	28.18	29.65	28.44
Chaetodon							
trifasciatus	38.73	50.63	53.47	60.17	63.19	64.38	66.72
Chromis opercularis	44.01	61.81	61.85	63.85	64.67	61.98	62.94
Chromis ternatensis	18.91	24.94	26.27	31.21	33.87	33.49	35.07
Gomphosus							
caeruleus	26.01	38.99	46.84	52.88	53.96	58.69	58.74
Halichoeres							
hortulanus	31.82	44.81	50.73	53.52	54.11	55.85	57.01
Monotaxis							
grandoculis	32.03	41.52	45.68	47.59	48.34	49.19	50.64
Naso brevirostris	47.06	54.26	61.01	59.46	59.30	62.40	64.66
Naso elegans	34.54	43.11	48.38	50.19	51.31	50.55	52.17
Oxymonacanthus							
longirostris	39.29	53.48	59.71	58.84	62.70	64.37	66.26
Pomacentrus							
sulfureus	70.90	88.18	90.92	93.08	92.67	94.13	93.93
Thalassoma							
hardwicke	25.64	44.40	57.80	62.03	64.97	67.86	67.96
Zanclus cornutus	18.33	31.28	37.27	41.71	41.70	42.95	44.44

Zebrasoma scopas	25.56	31.81	33.32	34.62	37.02	37.22	36.05
Mean	32.04	41.47	45.92	48.17	49.77	50.86	51.77
SD	12.70	16.93	17.69	18.00	18.10	18.52	18.96

Supp. Table 3: Probability values of the pairwise proportional test used to
assess the significance of the difference between accuracies obtained
through few-shot models with 5, 10, 15, 20, 25 and 30 shots and deep
learning model trained on T1.

	DL	FS1	FS5	FS10	FS15	FS20	FS25
FS1	<2e-16						
FS5	0.28	<2e-16					
FS10	1.20E-08	<2e-16	3.80E-12				
FS15	<2e-16	<2e-16	<2e-16	0.0016			
FS20	<2e-16	<2e-16	<2e-16	4.10E-09	0.0379		
FS25	<2e-16	<2e-16	<2e-16	1.30E-14	8.90E-05	0.2362	
FS30	<2e-16	<2e-16	<2e-16	<2e-16	3.90E-08	0.0059	0.28

602

603 Supp. Table 4: Value of the asymptotic accuracy predicted by the NLME

604 *models, and number of natural images required for both Deep Learning* 

architecture and Few-shot Learning architecture to reach 99% of this

606 *asymptote.* 

	Deep Learning		Few-Shot Learning	
	Number of images		Number of images	
	required to reach	Accuracy	required to reach	
	99% of the	asymptote	99% of the	Accuracy
Species	asymptote.	value	asymptote.	asymptote value
Naso brevirostris	1,506.47	67.43	5.04	38.10
Monotaxis				
grandoculis	1,769.71	62.72	7.59	38.14
Amblyglyphidodon				
indicus	1,766.55	62.89	5.34	41.76
Chromis				
ternatensis	1,492.65	67.75	6.16	46.55
Acanthurus	1,425.04	69.30	6.42	50.59

Mean	1221.54	77.34	7.45	60.87
vaigiensis	785.95	96.82	4.12	89.78
Abudefduf				
longirostris	833.75	93.29	5.68	80.97
Oxymonacanthus				
Naso elegans	1,074.08	80.48	6.66	78.13
hardwicke	994.27	84.00	5.86	78.42
Thalassoma				
sulfureus	960.99	85.63	15.41	72.98
Pomacentrus				
trifascialis	1,087.93	79.94	7.42	68.66
Chaetodon				
trifasciatus	1,076.99	80.34	7.70	63.88
Chaetodon				
Chaetodon auriga	1,053.06	81.31	14.97	65.94
guttatissimus	868.64	90.93	6.46	64.95
Chaetodon				
leucosternon	1,184.97	76.31	10.21	58.28
Acanthurus				
hortulanus	1,250.29	74.18	11.60	57.14
Halichoeres				
Zanclus cornutus	1,069.22	80.57	3.66	60.08
opercularis	1,112.28	78.67	5.85	54.95
Chromis				
caeruleus	1,341.80	71.51	6.90	53.15
Gomphosus				
Zebrasoma scopas	1,776.17	62.74	5.92	54.94
lineatus				