

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

Sébastien Villon^{a,*}, Corina Iovan^a, Morgan Mangeas^a, Thomas Claverie^{b,c}, David Mouillot^{c,d}, Sébastien Villéger^c, Laurent Vigliola^a

^a ENTROPIE, IRD, University of New-Caledonia, University of La Reunion, CNRS, Ifremer, Labex Corail, Noumea, New-Caledonia, France

^b CUFR Mayotte, France

^c MARBEC, University of Montpellier, CNRS, IRD, Ifremer, Montpellier, France

^d Institut Universitaire de France, Paris, France

ARTICLE INFO

Keywords:

Few-shot learning
Deep learning
Video
Marine biodiversity

ABSTRACT

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

1. Introduction

The world's ecosystems have entered an era of anthropogenic defaunation where human activities have triggered global decline in animal abundance, species range contraction and a new wave of species extinction (Dirzo et al., 2014). This global change is threatening ecosystem services worldwide hence the stability of our food systems, economies, and health. Defaunation is more advanced in terrestrial and freshwater ecosystems than in the marine environment where it started centuries later. However, the pace of defaunation is accelerating in oceans mostly due to the advent of industrial fishing since a century ago (Young et al., 2016). Given this context of global changes rapidly affecting fish communities, it is imperative to monitor fish biodiversity over time, on a large scale and using non-destructive methods.

Fish biodiversity surveys in the marine environment are typically

performed by divers. Although dive visual censuses provide a great deal of information on some shallow habitats, there are many limitations. First, divers are limited by depth and can hardly perform long dives to count fish below 30 m, ignoring mesophotic habitats and deeper ecosystems. Second, divers are limited by time and generally focus their 2–4 dives per day in the most speciose hard-substrate habitats, and ignore less rich and often immense adjacent soft-bottom habitats. Third, dive surveys provide data at a slow rate so that the compilation of global fish biodiversity database takes decades of efforts by multiple teams of highly skilled taxonomic divers (e.g. (Cinner et al., 2020; Stuart-smith et al., 2013)). This is a major restriction to the necessary temporal monitoring of global marine ecosystems, although a few time series exist in some countries¹ (Heenan et al., 2017).

Underwater videos (UV) are increasingly used (Whitmarsh et al., 2017) to overcome the limitations of diver-based surveys to quickly

* Corresponding author.

E-mail address: sebastien.villon@ird.fr (S. Villon).

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

<https://doi.org/10.1016/j.ecoinf.2021.101320>

Received 10 December 2020; Received in revised form 11 May 2021; Accepted 11 May 2021

Available online 14 May 2021

1574-9541/© 2021 Elsevier B.V. All rights reserved.

collect large amounts of data. For instance, more than 15,000 video stations were deployed in 58 countries in just three years for the first global assessment of the conservation status of reef sharks (Aaron MacNeil et al., 2020). Furthermore, underwater video surveys can be performed in many habitats, with some example in shallow reefs (Juhel et al., 2019), sandy lagoons (Cappo et al., 2007), deep sea (Zintzen et al., 2017), and even in the pelagic ecosystem (Letessier et al., 2019). Deploying underwater video stations does not require expert taxonomists and is now quite inexpensive with the improvement of cheap action cameras since a few years. The bottleneck to analyze this data now resides in the manual processing of the videos. Indeed, manually extracting fish biodiversity and abundance data from raw videos requires unsustainable workload by highly trained taxonomic experts. Although this annotation work can be improved through citizen science (McClure et al., 2020; Torney et al., 2019; Willi et al., 2019), such time-consuming and expensive task cannot match the increasing size of datasets, up to 20,000 h of videos for global surveys (Aaron MacNeil et al., 2020) and the necessary monitoring of global oceans over time.

As the demand for automatic methods to analyze underwater videos is rising, the latest generation of deep learning algorithms (DL), and in particular convolutional neural networks (CNNs) are increasingly used for species identification (Lasseck, 2020; Miao et al., 2019; Shiu et al., 2020; Willi et al., 2019) and fish detection (Qin et al., 2016; Rathi et al., 2018; Salman et al., 2016; Villon et al., 2018). However, these algorithms require a large dataset of annotated images (thumbnails hereafter) in order to train a robust model, able to provide satisfying results. Therefore, this method still requires collecting an important image dataset manually annotated by experts. This is especially problematic in highly diverse faunas such as coral reef fish that encompass nearly 6500 species worldwide (Chabanet et al., 2013). Furthermore, a universal pattern in species distribution, including fish communities, is that both rare and common species are found in every community, with the fraction of rare species more important in rich ecosystems, such as coral reefs (Hercos et al., 2013; Jones et al., 2002). 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 (Li et al., 2019a; Liu et al., 2020; Zhuang et al., 2018).

There are two ways to tackle this problem of lack of data. The first one consists of directly addressing the data itself, through data augmentation (Van Dyk and Meng, 2012; Wang and Perez, 2017; Wong et al., 2016). The second option is to change the classification algorithm. Few-shots learning (FSL) algorithms (Fei-fei et al., 2006; Fink, 2005) are designed to compute a classification task (query, noted *Q*) with only a few thumbnails to train (Support Sets, noted *SS*), and it has been increasingly studied since 2017 (Finn et al., 2017). Few-shots learning methods are divided into three main approaches. Metric-based methods are embedding both queries (*Q*) and support sets (*Ss*), before assigning to the query a class, according to distances computed between *Q* and *Ss* (Sung et al., 2018; Victor and Bruna, 2018; Yanbin et al., 2019). The second approach consists of 1) training a model on a large database, and 2) adapt this model to a new task with few examples, while not forgetting the concepts learned previously (Gidaris et al., 2018; Har-iharan et al., 2017). Finally, optimization-based methods are designed to adapt quickly to new tasks, hence able to learn a classification task with few examples (Finn et al., 2017; Nichol and Schulman, 2018; Sun and Chua, 2018). Optimization-based algorithms showed promising results in deep learning few-shot classification (Finn et al., 2017; Jamal and Cloud, 2019; Wang et al., 2019). Such methods propose to pre-train (or “meta-train”) a model with existing databases (e.g. MiniImageNet (Russakovsky et al., 2015), Ominglot (Lake et al., 2019)) on different tasks so it can adapt easily to a new one. For object identification, a task is defined by the classes the model has to discriminate. Once this model, called “meta-model” has been trained, it can then be tuned to operate on a new task with a very limited dataset, usually only 1–5 thumbnails per class.

In this study we propose to compare the efficiency of optimization-based few-shot learning and standard large dataset deep-learning methods to identify coral reef fish species on images. More specifically, we aim to determine how well a classic deep learning architecture trained with thousands of images and the benefit of data augmentation (hereafter DL) and FSL algorithms perform in situations where training thumbnail dataset is large or limited. To achieve this, we first trained a classic DL architecture built for image classification (He et al., 2016) on a large dataset of 69,169 thumbnails, and on a more limited dataset of 6320 thumbnails for 20 coral reef fish species. Then, we trained a few-shots, optimization-based learning algorithm (Nichol and Schulman, 2018) on the exact same training datasets while varying the number of shots from 1 to 30. Finally, we compared the capacity of DL and FSL models to correctly identify species on an independent thumbnail dataset, and modelled the asymptotic relationship between classification accuracy and the number of thumbnails in the training datasets for both classic DL and FSL algorithms.

2. Material and methods

2.1. 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 1920×1080 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. Figs. 1 and 2).

T0 is composed of 69,169 thumbnails extracted from 130 videos, with a range of 1134 to 7345 thumbnails per species (Table 1). *T1* is composed of 6320 thumbnails extracted from 20 videos with 40–1436 images per species whereas *T2* is composed of 13,232 thumbnails extracted from 25 videos with 55–3896 images per species. Thumbnails size originally ranged from 55×55 pixels to 500×450 pixels, but were

Table 1

Number of natural thumbnails extracted from the videos to build our three datasets.

Family	Species	Training dataset	Training dataset	Test dataset
		<i>T0</i>	<i>T1</i>	<i>T2</i>
Acanthuridae	<i>Acanthurus leucosternon</i>	3259	235	491
Acanthuridae	<i>Acanthurus lineatus</i>	1008	114	864
Acanthuridae	<i>Naso brevirostris</i>	1134	539	1932
Acanthuridae	<i>Naso elegans</i>	7345	1435	3896
Acanthuridae	<i>Zebrasoma scopas</i>	4970	48	579
Chaetodontidae	<i>Chaetodon auriga</i>	2134	737	502
Chaetodontidae	<i>Chaetodon</i>	1182	221	68
	<i>guttatissimus</i>			
Chaetodontidae	<i>Chaetodon trifascialis</i>	5234	41	630
Chaetodontidae	<i>Chaetodon trifasciatus</i>	4421	71	82
Labridae	<i>Gomphosus caeruleus</i>	3131	57	173
Labridae	<i>Halichoeres</i>	3192	40	287
	<i>hortulanus</i>			
Labridae	<i>Thalassoma</i>	4951	181	275
	<i>hardwicke</i>			
Lethrinidae	<i>Monotaxis grandoculis</i>	3893	797	1422
Monacanthidae	<i>Oxymonacanthus</i>	2553	54	55
	<i>longirostris</i>			
Pomacentridae	<i>Abudefduf vaigiensis</i>	5124	376	216
Pomacentridae	<i>Amblyglyphidodon</i>	1188	636	1310
	<i>indicus</i>			
Pomacentridae	<i>Chromis opercularis</i>	1525	81	93
Pomacentridae	<i>Chromis ternatensis</i>	3640	300	156
Pomacentridae	<i>Pomacentrus sulfureus</i>	5409	270	142
Zanclidae	<i>Zanclus cornutus</i>	3876	86	59
Total		69169	6320	13232

resized to 84×84 pixels before being processed through FS and DL algorithms.

The datasets *T1* and *T2* correspond to two real scenarios where videos were 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.

To train our DL architecture, we applied data augmentation to *T0* and *T1*. For each natural thumbnails in *T0* and *T1*, we created 9 thumbnails 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 thumbnail datasets and data augmentation are given in (Villon et al., 2020).

2.2. Experimental design

To compare classic deep-learning and few-shot algorithms in situations of large or small thumbnail datasets, we led five experiments using datasets *T0*, *T1*, *T2*, *AT0* and *AT1* described in Supp. Table 1:

- 1) We trained a classic DL algorithms architecture with our biggest dataset *AT0* as a baseline for the DL accuracy;
- 2) We trained the same DL architecture with the same hyper-parameters (e.g. model architecture and training process) but on a much more limited dataset (*AT1*). Hyper-parameters are the 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 computation);
- 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 *AT0*;
- 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 (He et al., 2016) as our classic deep-learning algorithm. Resnet is a convolutional neural network (CNN), a DL architecture which is able to both extract features from images and classify these images thanks to those features (Lecun et al., 2015). In order for a CNN to build an image classification model, the architecture is fed a large dataset, composed of pairs of labels and images. Using this dataset, the algorithms change their inner parameters in order to minimize the classification error, through a process called back-propagation. The ResNet architecture achieved the best results on ImageNet Large Scale Visual Recognition Competition (ILSVRC (Russakovsky et al., 2015)) in 2015, considered the most challenging image classification competition. It is still one of the best classification algorithms, while being easy to use and implement.

For the few-shot implementation, we used the Reptile algorithm (Nichol and Schulman, 2018). Few-shot learning algorithms are specific DL algorithms, whose goal is to be able to fit a model with very few training images. The Reptile algorithm is based on the well-known MAML architecture (Finn et al., 2017), and more precisely on the first-order version of MAML (Nichol et al., 2018). The Reptile algorithm is based on the division of the training dataset into a number of tasks T_i , a task being a learning problem. Through repetitively changing the task during the first 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 number of examples.

Here, the few-shot algorithms were tested on a classic n -ways k -shots procedure, n being the number of classes per support set, and k the number of images per class in the support set. For instance, a 5-ways 1-shots consists of training 5 classes with supports sets composed of 1 image per classes (e.g. species). We set $n = 5$ (Jamal and Cloud, 2019; Sun and Chua, 2018; Sung et al., 2018; Victor and Bruna, 2018; Wang et al., 2020) and allowed k to vary between 1 shot and 30 shots for both experiments 4 and 5. We did not use data augmentation for FSL experiments for several reasons. First, the goal of FSL is to adapt quickly with a very limited number of images. Second, to have similar settings for method comparison. There was no data-augmentation in the original paper, so we reproduced that. It also allowed us to compare our results with those obtained on benchmarks. Third, the reason behind the use of raw data instead of augmented data in few-shot learning paper is that with very few training samples and few conditions, the risk of overfitting by using the same image modified multiple times is far greater than in classic approaches with important datasets with many conditions.

2.3. Model comparison

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

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 *AT0* (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* (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 (*T0* and *AT0*).

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*).

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 *AT0* and *T0*):

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

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.

Eq. (1) was fitted by non-linear mixed-effect modelling (NLME (Pinheiro and Bates, 2006)) using species as a random effect. This method is widely used for fitting asymptotic processes. It allows estimating and comparing asymptotic accuracies of both FSL and DL algorithms, and the number of image to reach these asymptotic accuracies. The number of images required to reach the asymptotic accuracy was calculated as the number of images corresponding to an accuracy of 0.99 times the asymptotic value, meaning the asymptote was reached within 1%.

3. Results

The deep ResNet model trained on the large *AT0* dataset (3458 natural thumbnails in average per species) during the first experiment obtained a mean accuracy (i.e. percentage of correct classification) of 78.00% (standard deviation (SD) of 15.16%) on *T2* test-dataset (Table 2). With this model, accuracy varied among species between 54.14% (*Naso brevirostris*) and 99.07% (*Abudefduf vagiensis*). The same

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 T0 and T1 after data augmentation (AT0 and AT1).

Image per species (on average)	DL		FSL					
	T0	T1	T1			T0		
	3458	315	1 shot	5 shots	30 shots	1 shot	5 shots	30 shots
<i>Abudefduf vaigiensis</i>	99.07	69.91	16.08	11.39	11.38	47.67	70.9	86.35
<i>Acanthurus leucosternon</i>	86.15	44.67	25.51	30.71	38.74	19.23	28.8	42.66
<i>Acanthurus lineatus</i>	59.72	20.37	39.86	56.04	72.50	32.93	61.01	72.02
<i>Amblyglyphidodon indicus</i>	58.78	60.78	25.75	26.74	32.86	28.26	32.55	40.64
<i>Chaetodon auriga</i>	87.05	85.86	18.16	25.68	36.56	27.8	35.18	53.20
<i>Chaetodon guttatissimus</i>	85.50	44.12	33.58	44.21	58.26	29.61	51.18	79.29
<i>Chaetodon trifascialis</i>	90.00	3.49	29.02	25.48	28.44	27.14	43.17	63.51
<i>Chaetodon trifasciatus</i>	87.80	28.05	38.73	50.63	66.72	32.41	51.07	70.63
<i>Chromis opercularis</i>	61.29	9.68	44.01	61.81	62.94	45.34	68.28	81.50
<i>Chromis ternatensis</i>	59.61	55.77	18.91	24.94	35.07	35.4	55.44	67.22
<i>Gomphosus caeruleus</i>	75.72	20.81	26.01	38.99	58.74	28.96	39.16	54.22
<i>Halichoeres hortulanus</i>	82.93	17.07	31.82	44.81	57.01	28.94	41.35	58.87
<i>Monotaxis grandoculis</i>	57.10	53.37	32.03	41.52	50.64	32.8	45.85	59.13
<i>Naso brevirostris</i>	54.14	68.60	47.06	54.26	64.66	54.47	58.08	61.00
<i>Naso elegans</i>	93.24	79.43	34.54	43.11	52.17	28.47	33.71	44.36
<i>Oxymonacanthus longirostris</i>	96.43	14.54	39.29	53.48	66.26	42.15	65.44	84.86
<i>Pomacentrus sulfureus</i>	90.14	61.97	70.90	88.18	93.93	65.53	86.21	90.00
<i>Thalassoma hardwicke</i>	90.90	51.64	25.64	44.4	67.96	26.72	45.7	70.60
<i>Zanclus cornutus</i>	81.36	40.68	18.33	31.28	44.44	26.08	40.82	62.72
<i>Zebrasoma scopas</i>	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

ResNet DL model trained on smaller AT1 (315 natural thumbnails in average per species) during the second experiment showed highly degraded performance with a mean accuracy of only 42.21% (SD = 24.95%). Among species variation ranged with this model from only 3.49% (*Chaetodon trifascialis*) to 85.86% (*Chaetodon auriga*).

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 AT1), and performed better beyond 10 shots (45.92% of accuracy on T2 with 10-shots learning). A pairwise proportion test showed a p -value < 0.0001 , assessing that FSL was significantly better than DL in this scenario beyond 10 shots (Supp. Table 3) accuracy of Reptile models had a standard deviation from 12.70% with one-shot learning, to 18.96% with 30-shots learning, indicating important variation in accuracy among species. However, this standard deviation was smaller than that of the ResNet algorithm trained on the same AT1 limited dataset (24.95%).

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 5-shots, 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 (Fig. 1).

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

The fixed-effect asymptotic value of accuracy was higher for ResNet model ($Accuracy_{\infty} = 77.34\%$, 95% CI: 71.26–83.41%) than for Reptile model ($Accuracy_{\infty} = 60.87\%$, 95% CI: 54.48–67.26%), illustrating higher classification power of ResNet over Reptile when large numbers of thumbnails are available. However, the slope of the asymptotic model was two-orders of magnitude higher for Reptile (0.707, 95% CI: 0.559–0.854) than for ResNet architecture (0.0040, 95% CI: 0.0032–0.0048), illustrating the high capacity of Reptile FSL algorithm to learn from only a few images. NLME modelling further showed that

average asymptotic accuracy was reached with only 7 natural thumbnails per species for Reptile architecture, compared to 1153 natural thumbnails per species for ResNet, confirming the strong power of Reptile method in situation of limited thumbnail training dataset. However, model random effects showed that some variation existed among species. For Reptile architecture, asymptotic accuracy values ranged from 38.09% (*Amblyglyphidodon indicus*) to 89.78% (*Pomacentrus sulfureus*), and was reached with 4 to 16 training images per species. For DL architecture, species asymptotes varied from 62.72% (*Monotaxis grandoculis*) to 96.81% (*Abudefduf vaigiensis*), and could be reached with 786–1776 thumbnails per species (Supp. Table 4).

4. Discussion

Our experiments demonstrated that few-shot learning methods based on Reptile architecture can be effectively used to drastically reduce the number of annotated 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 standard Deep Learning architecture such as ResNet trained with 400–14,350 images per species. Further, FSL architecture trained with 10 images outperformed a ResNet 100 architecture trained with at least 400 images per species. This is a very promising result in situations where many species need to be identified from models trained with a few images, a typical characteristic in marine biodiversity applications.

However, the important standard deviation among the different trained species (18.96 SD on 30-shots) showed that few-shot algorithms may not be robust enough to discriminate among similar species showing only subtle differences. Nevertheless, in our 2nd experiment, our ResNet model achieved an accuracy under 40% for all the species with fewer training images than 1140 (after data augmentation, i.e. 114 natural images), and only 7 species were identified with an accuracy greater than 45%. These species were represented with a range of 2700–14,350 images during the training phase. We also show better results with the model trained on T0 than the model trained on T1. As expected, increasing the number of images per shot 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 few images per classes are at disposal. Therefore, the dataset

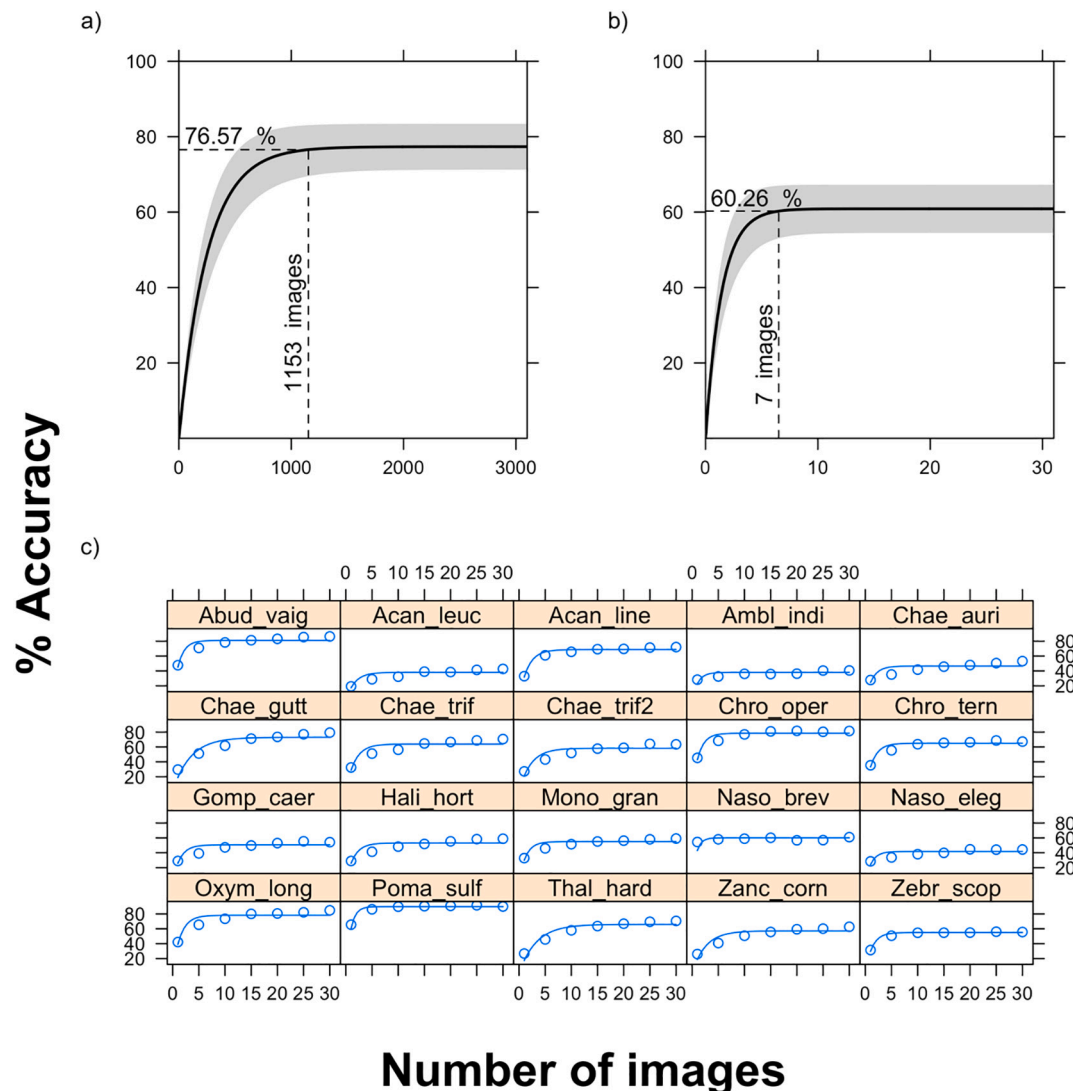


Fig. 1. Relationship between the number of natural thumbnails per species for training and the accuracy of deep-learning and few-shot learning models. Non-linear mixed effects asymptotic model fit for (a) DL architecture at the fixed-effect level, and for FSL architecture at (b) the fixed-effect level and (c) the random-effect level. Grey areas represent 95% CI in fixed-effects estimates. 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 the 95% asymptote value, reached with 750 images for DL and 5 with FSL.

T1 corresponded 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.

Modelling 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 Mini-ImageNet dataset (the most used benchmark for few-shot learning methods) through a 5-shots learning, (Li et al., 2019b) 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.

Applied to marine and coral reef ecology, such methods requiring few examples to 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 composed of few very common species and a large proportion of less-common and rare species. Hence, the important effort required to build databases with a 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 improvement of few-shot learning algorithms offers promises to build efficient identification models to automatically process images and videos to localise and identify rare fish species. Such models could then be paired with more classic deep architectures, more efficient to identify abundant species with the leverage of important datasets.

Declaration of Competing Interest

None.

Acknowledgements

This study was funded by the French National Research Agency

project ANR 18-CE02-0016 SEAMOUNTS.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.ecoinf.2021.101320>.

References

- Aaron MacNeil, M., et al., 2020. Global status and conservation potential of reef sharks. *Nature* (July), 2019.
- Cappo, M., De, G., Speare, P., 2007. Inter-reef vertebrate communities of the Great Barrier Reef Marine Park determined by baited remote underwater video stations. *Mar. Ecol. Prog. Ser.* 350, 209–221.
- Chabanet, P., Floeter, S.R., Friedlander, A., McPherson, J., Myers, R.E., 2013. Global biogeography of reef fishes : a hierarchical quantitative delineation of regions. *PLoS One* 8 (12).
- Cinner, J.E., et al., 2020. Meeting fisheries, ecosystem function, and biodiversity goals in a human-dominated world. *Science* (80-.) 311 (April), 307–311.
- Dirzo, R., Young, H.S., Galetti, M., Ceballos, G., Isaac, N.J.B., Collen, B., 2014. Defaunation in the anthropocene. *Science* (80-.). 345 (6195), 401–406.
- Fei-fei, L., Fergus, R., Member, S., Perona, P., 2006. One-shot learning of object categories. *IEEE Trans. pattern Anal. Mach. Intell.* 28 (4), 594–611.
- Fink, M., 2005. Object classification from a single example utilizing class relevance metrics. In: *Advances in Neural Information Processing Systems*, pp. 449–456.
- Finn, C., Abbeel, P., Levine, S., 2017. Model-agnostic Meta-learning for Fast Adaptation of Deep Networks. *arXiv Prepr. arXiv1703.03400*.
- Gidaris, S., Paristech, P., Komodakis, N., Paristech, P., 2018. Dynamic few-shot visual learning without forgetting. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4367–4375.
- Hariharan, B., Girshick, R., Ai, F., 2017. Low-shot visual recognition by shrinking and hallucinating features. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3018–3027.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778.
- Heenan, A., et al., 2017. Long-term monitoring of coral reef fish assemblages in the Western central pacific. *Sci. Data* 4, 1–12.
- Hercos, A.P., Sobansky, M., Queiroz, H.L., Magurran, A.E., Andre, A., 2013. Local and regional rarity in a diverse tropical fish assemblage. *Biol. Sci.* 280, 81–101.
- Jamal, M.A., Cloud, H., 2019. Task agnostic meta-learning for few-shot learning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jones, G.E., Caley, M.J., Munday, P.L., 2002. “Rarity in Coral Reef Fish Communities,” *Coral Reef Fishes Dyn. Divers. A Complex Ecosyst.* pp. 88–101.
- Juhel, J., Vigliola, L., Wantiez, L., Letessier, T.B., Meeuwig, J.J., Mouillot, D., 2019. Isolation and no-entry marine reserves mitigate anthropogenic impacts on grey reef shark behavior. *Sci. Rep.* 9 (November 2018), 1–11.
- Lake, B.M., Salakhutdinov, R., Tenenbaum, J.B., 2019. The Omniglot challenge : a 3-year progress report. *COBEHA* 29, 97–104.
- Lasseck, M., 2020. Audio-based Bird Species Identification With Deep Convolutional Neural Networks Audio-based Bird Species Identification With Deep Convolutional Neural Networks. January.
- Lecun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 436–444.
- Letessier, Tom B., et al., 2019. Remote reefs and seamounts are the last refuges for marine predators across the Indo- Pacific. *PLoS Biol.* 17, 1–20.
- Li, A., Luo, T., Lu, Z., Xiang, T., Wang, L., 2019a. Large-scale few-shot learning : knowledge transfer with class hierarchy. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7212–7220.
- Li, H., Eigen, D., Dodge, S., Zeiler, M., Wang, X., 2019b. Finding task-relevant features for few-shot learning by category traversal. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 1.
- Liu, L., Zhou, T., Long, G., Jiang, J., Zhang, C., 2020. Many-class few-shot learning on multi-granularity class hierarchy. *IEEE Trans. Knowl. Data Eng.* 1–14.
- McClure, E.C., et al., 2020. Artificial intelligence meets citizen science to supercharge ecological monitoring. *Patterns* 1 (7), 100109.
- Miao, Z., et al., 2019. Insights and approaches using deep learning to classify wildlife. *Sci. Rep.* (May), 1–9.
- Nichol, A., Schulman, J., 2018. Reptile : A Scalable Metalearning Algorithm. *arXiv Prepr. arXiv1803.02999*, pp. 1–11 (2018).
- Nichol, A., Achiam, J., Schulman, J., 2018. On First-order Meta-learning Algorithms. *arXiv*, pp. 1–15.
- Pinheiro, J., Bates, D., 2006. Mixed-effects Models in S and S-PLUS.
- Qin, H., Li, X., Liang, J., Peng, Y., Zhang, C., 2016. DeepFish: accurate underwater live fish recognition with a deep architecture. *Neurocomputing* 187, 49–58.
- Rathi, D., Jain, S., Indu, S., 2018. Underwater Fish Species Classification using Convolutional Neural Network and Deep Learning. (*arXiv:1805.10106v1 [cs.CV]*). June.
- Russakovsky, O., et al., 2015. ImageNet large scale visual recognition challenge. *Int. J. Comput. Vis.* 211–252.
- Salman, A., et al., 2016. Oceanography : Methods Fish Species Classification in Unconstrained Underwater Environments Based on Deep Learning, pp. 570–585.
- Shiu, Y., et al., 2020. Deep Neural Networks For Automated Detection of Marine Mammal Species, pp. 1–12.
- Stuart-smith, R.D., et al., 2013. Integrating abundance and functional traits reveals new global hotspots of fish diversity. *Nature* 501 (7468), 539–542.
- Sun, Q., Chua, Y.L.T., 2018. Meta-transfer learning for few-shot learning. *Conf. Comput. Vis. Pattern Recognit.* 403–412.
- Sung, F., Yang, Y., Zhang, L., 2018. Learning to compare : relation network for few-shot learning Queen Mary University of London. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1199–1208.
- Torney, C.J., et al., 2019. A comparison of deep learning and citizen science techniques for counting wildlife in aerial survey images. *Methods Ecol. Evol.* 10 (October 2018), 779–787.
- Van Dyk, D.A., Meng, X., 2012. The art of data augmentation. *J. Comput. Graph. Stat.* 8600 (2001), 1–50.
- Victor, J., Bruna, Garcia, 2018. Few-shot Learning With Graph Neural Networks. *arXiv preprint arXiv:1711.04043*, 2017, pp. 1–13.
- Villon, S., et al., 2018. A Deep Learning algorithm for accurate and fast identification of coral reef fishes in underwater videos. *PeerJ Prepr.* 6, e26818v1.
- Villon, S., Mouillot, D., Chaumont, M., Subsol, G., 2020. A new method to control error rates in automated species identification with deep learning algorithms. *Sci. Rep.* 10, 1–13.
- Wang, J., Perez, L., 2017. The Effectiveness of Data Augmentation in Image Classification Using Deep Learning. *arXiv Prepr. arXiv1712.04621*.
- Wang, Y., Yao, Q., Kwok, J.T., Ni, L.M., 2019. Generalizing from a Few Examples: A Survey on Few-shot Learning *arXiv* : 1904 . 05046v2 [cs. LG] 13 May 2019.
- Wang, Y., Yao, Q., Ni, L.M., 2020. Generalizing from a few examples : a survey on few-shot generalizing from a few examples : a survey on few-shot. In: *ACM Comput. Surv.* 53. June.
- Whitmarsh, S.K., Fairweather, P.G., Huveneers, C., 2017. What is Big BRUVver up to ? Methods and uses of baited underwater video. *Rev. Fish Biol. Fish.* 27 (1), 53–73.
- Willi, M., et al., 2019. Identifying animal species in camera trap images using deep learning and citizen science. *Methods Ecol. Evol.* 10 (1), 80–91.
- Wong, S.C., McDonnell, M.D., Adam, G., Victor, S., 2016. Understanding data augmentation for classification : when to warp ? In: *2016 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pp. 1–6.
- Yanbin, Liu, et al., 2019. Learning to Proagate Labels: Transductive Propagation Network for Few-shot Learning. *arXiv preprint arXiv:1805.10002*, pp. 1–14.
- Young, H.S., Mccauley, D.J., Galetti, M., Dirzo, R., 2016. Patterns, causes, and consequences of anthropocene defaunation. *Annu. Rev. Ecol. Syst.* (August), 333–358.
- Zhuang, P., Wang, Y., Qiao, Y., 2018. Wildfish : a large benchmark for fish recognition in the wild. In: *Proceedings of the 26th ACM international conference on Multimedia*, 2, pp. 1301–1309.
- Zintzen, V., Anderson, M.J., Roberts, C.D., Harvey, E.S., Andrew, L., 2017. Effects of latitude and depth on the beta diversity of New Zealand fish communities. *Sci. Rep.* 7 (July), 1–10.