# Towards Automatic Cetacean Photo-Identification: A Framework for Fine-Grain, Few-Shot Learning in Marine Ecology

Cameron Trotter*, Nick Wright*, A. Stephen McGough†, Matt Sharpe‡,
Barbara Cheney§, Mònica Arso Civil¶, Reny Tyson Moore‖, Jason Allen‖, and Per Berggren‡

*School of Engineering, †School of Computing, ‡School of Natural & Environmental Sciences
Newcastle University, UK
{c.trotter2, nick.wright, stephen.mcgough, m.j.sharpe, per.berggren}@ncl.ac.uk

§School of Biological Sciences
University of Aberdeen, UK
b.cheney@abdn.ac.uk

¶School of Biology
University of St Andrews, UK
mac64@st-andrews.ac.uk

‖Chicago Zoological Society's Sarasota Dolphin Research Program
Sarasota, FL, USA
renytysonmoore@gmail.com, allenjb@mote.org

*Abstract*—Photo-identification (photo-id) is one of the main non-invasive capture-recapture methods utilised by marine researchers for monitoring cetacean (dolphin, whale, and porpoise) populations. This method has historically been performed manually resulting in high workload and cost due to the vast number of images collected. Recently automated aids have been developed to help speed-up photo-id, although they are often disjoint in their processing and do not utilise all available identifying information. Work presented in this paper aims to create a fully automatic photo-id aid capable of providing most likely matches based on all available information without the need for data pre-processing such as cropping. This is achieved through a pipeline of computer vision models and post-processing techniques aimed at detecting cetaceans in unedited field imagery before passing them downstream for individual level catalogue matching. The system is capable of handling previously uncatalogued individuals and flagging these for investigation thanks to catalogue similarity comparison. We evaluate the system against multiple real-life photo-id catalogues, achieving mAP@IOU[0.5] = 0.91, 0.96 for the task of dorsal fin detection on catalogues from Tanzania and the UK respectively and 83.1, 97.5% top-10 accuracy for the task of individual classification on catalogues from the UK and USA.

*Index Terms*—Few-Shot, Fine-Grain Classification, Detection

## I. INTRODUCTION

In recent years there has been a concerted effort to apply computer vision techniques to challenging big data problems which can have a positive societal impact. A highly important area where computer vision can help is ecology [1]. One of the main goals of ecological research is to monitor animal populations in their distribution area, undertaking abundance estimates to inform policy change. This is most commonly performed using capture-recapture surveys where researchers identify the presence of individuals and estimate abundance of animals in an area to produce population estimates [2–4]. These surveys can be classified as invasive where animals are physically trapped, tagged, and released, or non-invasive where monitoring is performed passively such as via the collection of images – referred to as photo-id.

Photo-id is one of the main non-invasive capture-recapture methods utilised by cetacean researchers [5]. Surveys are usually undertaken at sea, although monitoring from coastlines or aircraft may also be utilised [6, 7]. The methodology is employed for the monitoring of multiple cetacean species, with a range of studies demonstrating its efficacy [8, 9].

All non-invasive capture-recapture methodologies rely on the target species having some form of individually identifiable markings. Depending on the species, different parts of the body are the primary identifying feature; for dolphins this is usually the dorsal fin as this body part is most likely to be visible above the waterline. During photo-id surveys, researchers often focus on long lasting stable markers such as dorsal fin shape, notches, scarring, and pigmentation. These markings can be difficult to capture in detail due to the free roaming nature of the animals causing high variances in angles of approach, direction of travel, distance from camera, and surfacing elevation, as seen in Figure 1. This is exacerbated when dealing with cetacean species which travel in pods, making it difficult to distinguish the individuals present.

Marine photo-id can be extremely labour and cost intensive compared to on-land surveys, which rely on the use of camera traps placed in stationary locations to capture images when they detect movement. This setup is not possible at sea due to a lack of stationary objects to attach devices to and rapid movement in the observed scene due to waves causing the

Figure 1. Two images of the same individual taken from different angles of approach, directions of travel, surfacing elevation, and distances from the camera. This changes the make-up of the dorsal fin but retains identifying information. Images from [10].

camera to trigger – producing a high false positive rate.

Upon survey completion, photo-id data must be analysed and individuals identified to produce a catalogue. Images collected during surveys are large in size and contain significant amounts of background noise. Historically curation of this data has been a manual process that often takes longer than the entire data collection period [11], further increasing labour and costs. As such, any techniques to speed up the curation process would be welcomed by both researchers and their funding bodies. As photo-id surveys are not guaranteed to capture all individuals in a given geographic area, naive approaches such as training a simple image classifier on existing catalogue examples would not suffice as they are incapable of flagging previously uncatalogued individuals.

This work details a framework for fully automatic catalogue matching based on unprocessed photo-id imagery. This is achieved by a pipeline of trained computer vision models and robust post-processing techniques capable of automatic fin detection and most likely catalogue matching based on latent space similarity. Images are passed through a Mask R-CNN [12] dorsal fin detector, removing the need for manual data pre-processing. Detections are post-processed ready for fine-grain, few-shot catalogue matching via a Siamese Neural Network (SNN) trained using triplet loss [13] and online semi-hard triplet mining to create a latent space based on the provided catalogue. Matches are obtained using the Euclidean distances between an input and generated class prototypes, allowing for the flagging of potentially uncatalogued individuals. This reduction in data processing affords cetacean researchers more time to work on application of their data, for example to inform mitigation and policy change, rather than curation.

## II. RELATED WORK

Due to the time and labour requirements of manual photo-id, multiple aids have been developed. Descriptions of these related works is provided, with an overview in Table I.

**Catalogue Management Systems** are widely used to aid manual photo-id, especially in geographic locations with large resident populations. FinScan [14] allows users to upload pre-processed fin images which they then trace around. This trace is checked against a database to determine most likely matches. Likewise, DARWIN [16] also provides automated photo-id based on traces.

FinBase [15] instead manages catalogues based on user defined attributes which can then be used for matching. Fins



Figure 2. An example showing the result of SURF [26] feature extraction on the dorsal fin of an Indo-Pacific bottlenose dolphin.

are matched based on querying a database for entries with matching attributes. Additionally, CatRlog [17] provides likely matches based on manually entered marking information.

**Non-Deep Learning Approaches** to aid cetacean photo-id are available. Karnowski *et al.* [19] used PCA to subtract background from underwater images of common bottlenose dolphins (*Tursiops truncatus*). Further, Weideman *et al.* [18] proposed CurvRank, an algorithm which automatically identifies a fin's trailing edge for likely matching.

**Deep Learning Approaches** have also been explored in recent years, inching work towards a fully automatic photo-id system. Data pre-processing such as cropping is one of the main labour costs in catalogue creation. Photo-ID Ninja[1] aims to speed up this processing by automatically cropping images to the dorsal fin.

Quiñonez *et al.* [20] proposed a detection system capable of distinguishing between four distinct classes: `dolphin`, `dolphin_pod`, `open_sea`, and `seabirds`. Morteo *et al.* [21] performed semi-automatic matching by projecting lines from the base of the fin's leading edge.

Bouma *et al.* [22] provided a system focusing on metric learning to photo-id individual New Zealand common dolphins (*Delphinus spp*), utilising Photo-ID Ninja to crop fins prior to matching. Lee *et al.* [23] proposed an architecture for cetacean identification via normalised segmentation and significant feature extraction.

DolFin [25] uses a SURF-based [26] approach to identify Risso's dolphins (*Grampus griseus*), a species prone to long term scarring which makes them ideal for matching via feature extractors. This approach fails with species where identifying markings are more subtle, such as in Figure 2 where SURF has failed to extract the top left notch from the dorsal fin of an Indo-Pacific bottlenose dolphin (*T. aduncus*).

FinFindR [24] allows for the identification of common bottlenose dolphins based on the trailing edge of the animal's dorsal fin, utilising this to cluster individuals.

### A. Comparison Against Related Work

Table I provides a summary of related works. The vast majority are standalone algorithms, requiring researchers to set up their own data pipelines. Only DolFin and finFindR propose fully automated, self-contained photo-id aids which require no pre-processing and are capable of handling previously uncatalogued individuals.

As noted previously, DolFin may struggle to handle species other than Risso's dolphin as SURF fails to extract subtle features such as notches. Further, finFindR fails to make use of

---

[1]Photo-ID Ninja: photoid.ninja

Table I
A COMPARISON OF AVAILABLE PHOTO-ID AIDS.

| System | Requires Data Pre-processing | Dorsal Fin Detection | Full Background Removal | Individual Photo-ID | Can Flag Individuals Not Currently In Catalogue | Uses All Information Found on Dorsal |
|---|---|---|---|---|---|---|
| FinScan [14] | ✓ | ✗ | ✗ | ✓ | ✓ | — |
| FinBase [15] | ✓ | ✗ | ✗ | ✓ | ✓ | — |
| DARWIN [16] | ✓ | ✗ | ✗ | ✓ | ✓ | — |
| catRlog [17] | ✓ | ✗ | ✗ | ✓ | ✓ | — |
| CurvRank [18] | ✓* | ✗* | ✗ | ✓ | ? | ✗ |
| Karnowski *et al.* [19] | ✗ | ✓ | ✓ | ✗ | — | — |
| Photo-ID Ninja[1] | ✗ | ✓ | ✗ | ✗ | — | — |
| Quiñonez *et al.* [20] | ✗ | ✓ | ✗ | ✗ | — | — |
| Morteo *et al.* [21] | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Bouma *et al.* [22] | ✓ | ✓† | ✗ | ✓ | ✓ | ✓ |
| Lee *et al* [23] | ✗ | ✓ | ✓ | ✓ | ? | — |
| finFindR [24] | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ |
| DolFin [25] | ✗ | ✓ | ✓ | ✓§ | ✓ | ✓ |
| Ours | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ |

\* CurvRank is included in Flukebook (flukebook.org), a website for cross-catalogue matching. In this scenario, Flukebook performs automatic data pre-processing and fin detection before passing to CurvRank, but the algorithm itself does not facilitate this.
? It is unclear whether the system is capable of flagging previously uncatalogued individuals to the researcher based on the literature available.
† Utilises Photo-ID Ninja for detection.
§ Utilises SURF for photo-id, thus unsuitable for cetacean species without well-defined markings.

all information on the dorsal, utilising only the trailing edge for matching, and does not fully remove background noise which may influence the matching process.

Work proposed in this paper aims to overcome these limitations. The outlined methodology performs full background removal, negating the effect of noise on matching, and is capable of operating on a range of species rather than just those prone to specific markings by extracting all available information. Images can also be operated over without the need for manual data pre-processing.

## III. METHODOLOGY

The proposed framework allows for the detection and individual identification of various cetacean species. This is achieved through a pipeline of models (see Figure 3) to sequentially process input images.

**Dorsal Fin Detection** is achieved using a Mask R-CNN model to locate regions of interest (RoIs) in images, defined as areas where the animal's dorsal fin is visible above the waterline. This model is trained on large scale images (3456x5184px) from DSLR cameras, typical of those utilised during photo-id surveys. By detecting RoIs in input images automatically, the requirement for data pre-processing is removed. Model outputs are in the form of masks which precisely differentiate between a coarse-grain dorsal fin or background. If an image contains a pod ($> 1$ dorsal fin), each detected RoI is processed sequentially downstream.

**Morphological Transformations** are performed based on *a priori* knowledge of cetaceans. It can be reasoned that holes in a mask are likely unintentional and a product of surrounding noise. In this instance, the model may have failed to capture all available information. Any holes present in masks are filled using dilation and erosion morphological transformations to ensure no identifiable information is lost. Note that any holes present in the dorsal fin from natural or anthropogenic activity such as from sting ray barbs are not transformed to retain identifying information. A *bitwise-and* operation is

then applied between the clean mask and the input image, segmenting the RoI to reduce the amount of noise passed downstream.

**Colour Thresholding** is executed over masks with multiple disjoint components. This may occur if, for example, an area of splash has been erroneously included as part of a detection. As a single cetacean cannot be made up of multiple disjoint components, it is known that some of these are erroneous and should be removed.

As the outer layer of cetaceans' skin is often a consistent grey colouring, minus any prominent identifiable markings, this can be used to filter mask components. By comparing the colour composition of each component against a calculated threshold, it is possible to discard those which have been erroneously detected. As these components are often areas of water, they will likely be much lighter in composition than the cetacean component. An example of noise removal via colour thresholding is shown in Figure 4. If multiple mask components pass noise removal and colour thresholding, each component is treated as a distinct mask downstream. Masks with only a single component are not colour thresholded to ensure no detections are ignored due to post-processing, preventing the discarding of an RoI which contains no disjoint components but is above the threshold, such as in the event of extreme over-exposure.

**Cropping** is then undertaken, reducing the input image down to the processed RoI. This vastly reduces the image file size and computational expense of downstream operations, as well as centres the RoI in the output image.

**Most Likely Catalogue Matching** is applied on a per-catalogue basis using a trained SNN with a triplet loss function [13]. SNNs have demonstrated usefulness for few-shot, fine-grain problems in research domains such as species identification [27, 28]. Our work compliments this through extension to individual-level identification.

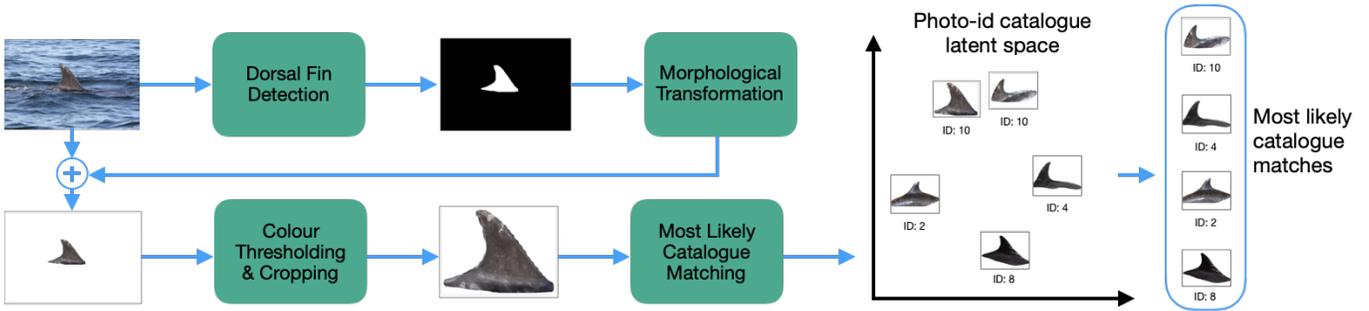At train time, the photo-id catalogue is processed as per the

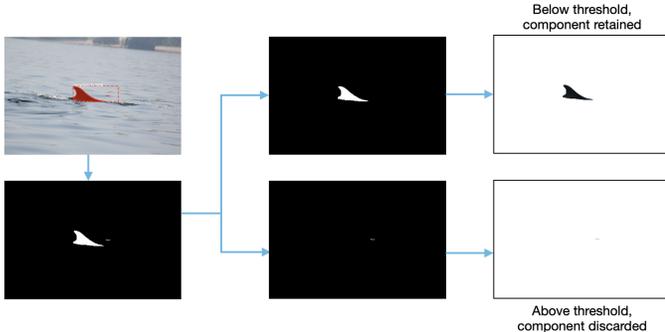Figure 3. A high level overview of data flow through the proposed system.



Figure 4. Workflow detailing colour thresholding to remove disjoint splash. The mask is split into individual components. The background subtracted images are colour thresholded, and erroneous splash is discarded.

methodology above. Each post-processed RoI for an individual is used as a class example for SNN training. A `noise` class is also included, containing examples of retained erroneous detections, to afford the model the ability to rule these out during inference without human intervention.

The trained SNN is capable of generating low-dimensional embeddings from fin images in such a way that those of the same individual generate embeddings which are close together in the latent space, creating individual class clusters which allow for most likely catalogue matching. During inference, most likely matching is performed via Euclidean distance measurement between the input and prototypes representing the median example of each class.

By clustering individuals together in the latent space based on similarity and comparing new images to the class prototypes, the system can flag potentially previously uncatalogued individuals. A new individual who enters the catalogue's survey area (e.g. through migration or birth) is placed in a distinct latent space location, resulting in large distances between it and the class prototypes. If a new individual is flagged it can then be verified by a human and, if correct, added to the catalogue. When a new entry is appended, a class prototype for the newly introduced individual can be generated, allowing the system to perform catalogue matching without the need for model re-training.

## IV. EXPERIMENTATION

To evaluate the proposed methodology, experimentation was performed using multiple real-life photo-id catalogues collected by various institutions from a range of geographic locations, times, and encompassing multiple cetacean species.

### A. Dorsal Fin Detection

**Precision** The dorsal fin detector's ability to precisely detect RoIs was evaluated using mean average precision at differing intersection over union thresholds (mAP@IOU). This was performed using 1021 photo-id survey images provided by Newcastle University's Marine MEGAfauna Lab of Indo-Pacific bottlenose dolphins resident in the coastal waters of Zanzibar, Tanzania in 2015 [2]. RoIs were labelled with a coarse-grain `dolphin` class, resulting in 616 total examples split into an 80-20 train-test split.

A ResNet50 architecture [29] was utilised as a model backbone, with a minimum confidence threshold of 0.9. Hyperparameter optimisation was performed using a grid search, with the following possibilities examined: (1) *Weight Decay*: 0.01, 0.001, 0.0001, or 0.0001. (2) *RPN Anchor Scale*: (8, 16, 32, 64, 128), (16, 32, 64, 128, 256), or (32, 64, 128, 256, 512). (3) *Optimiser*: Adam [30], or SGD with Warm Restarts [31]. (4) *Pre-trained on MSCOCO* [32]: yes or no. (5) *Data Augmentation Strategy*: `aug1`, `aug2`, or None.

The first strategy, `aug1`, selected up to three of the following: (1) *Horizontal Flip*: (p = 0.5). (2) *Vertical Flip*: (p = 0.5). (3) *Rotation*: 90, 180, or 270 degrees (p = 0.33). (4) *Scaling*: 80% to 120% on both axes independently. (5) *Brightness*: multiply all pixels in the image with a random value between 0.8 and 1.5. (6) *Gaussian Blur*: using a kernel with radius randomly assigned between 0 and 5.

The second strategy, `aug2`, was more complex, performing the following perturbations in a sequentially random order on 67% of the images: (1) *Horizontal Flip*: (p = 0.5). (2) *Cropping*: each side of the image randomly between 0% and 10% of the total side length. (3) *Gaussian Blur*: using a kernel with radius randomly assigned between 0 and 2.5 (p = 0.5). (4) *Contrast*: increase or decrease by a random factor between 0.75 and 1.5. (5) *Additive Gaussian Noise*: sample the noise per channel, adding noise to the colour of the pixels. (6) *Brightness*: multiply all pixels with a random value between 0.8 and 1.2. (7) *Scaling*: 80% to 120% on both axes independently. (8) *Rotation*: randomly between -180 and 180 degrees. The use of these augmentation strategies allowed

Figure 5. Example of dorsal fin detections. Left: An image showing cetaceans travelling in a pod. Each detection is post-processed and identified in isolation. Right: An image showing the effect of eco-tourism on the false positive rate of the detector.

TABLE II
INSTANCE SEGMENTATION RESULTS OF THE MASK R-CNN DETECTOR, TRAINED USING THE ZANZIBAR TRAINING DATA.

| Dataset | mAP@IOU[x] | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0.50 | 0.55 | 0.60 | 0.65 | 0.70 | 0.75 | 0.80 | 0.85 | 0.90 | 0.95 |
| **Zanzibar** | 0.91 | 0.91 | 0.89 | 0.86 | 0.85 | 0.79 | 0.69 | 0.50 | 0.15 | 0.00 |
| **NDD** | 0.96 | 0.95 | 0.93 | 0.91 | 0.88 | 0.83 | 0.71 | 0.51 | 0.16 | 0.00 |

for evaluation of whether a simple or more complex strategy would be more appropriate for this use case.

The search determined that training a model optimised using SGD with Warm Restarts, alongside an initial learning rate of 0.001 with 0.01 weight decay, RPN anchor scales of (16, 32, 64, 128, 256), pretrained on MSCOCO and using the *aug1* strategy produced the highest mAP@IOU scores.

Experimental results for the trained model on the Zanzibar data can be seen in Table II (Top). These images contain a wide variety of background noise. Furthermore, some dorsal fins in the images look similar in shape and structure to background objects – especially in choppy waters captured from a distance. The animal's bodies are also similarly coloured to their surroundings. These adaptations allow the animals to be better camouflaged in their environment, but can cause issues for detection systems. The model is capable of precisely detecting multiple fins when the animals are captured travelling in a pod, as seen in Figure 5 (Left) where the pod has been separated into its constituent parts ready for identification downstream. **Generalisability** The detector is also required to produce detections with high mAP scores when operating on data from different geographic regions, time intervals, and species, as this would negate the need for detector re-training when working in a new survey area. To evaluate this, the model was used to generate mask predictions for the above water image set found in the Northumberland Dolphin Dataset (NDD) [10]. This open-source dataset contains images of both common bottlenose dolphins and white-beaked dolphins (*Lagenorhynchus albirostris*) collected during a 2019 photo-id survey off the coast of Northumberland, UK.

To test this generalisability, the model trained using the Zanzibar data was evaluated on the whole above water set of NDD without re-training or fine-tuning. Results for this can be seen in Table II (Bottom). Interestingly, the model achieves a higher mAP at the given thresholds on NDD than the Zanzibar dataset on which it was trained. This is hypothesised to be due to the lack of other objects in NDD in comparison to the Zanzibar dataset. For example, some images in the Zanzibar
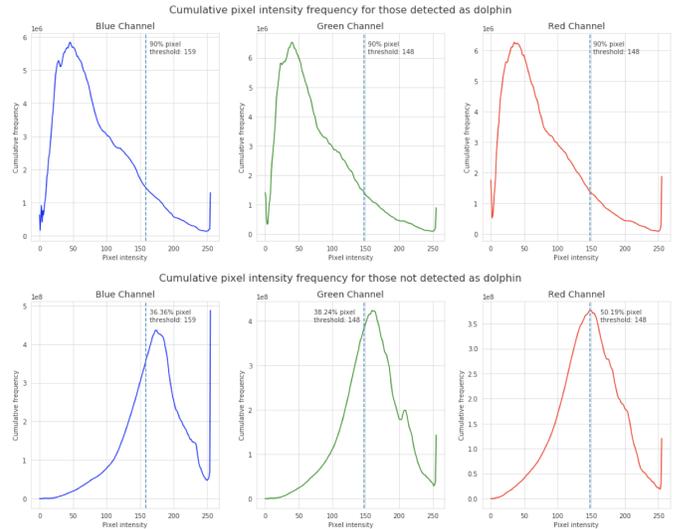


Figure 6. The global range of pixel intensities for each RGB colour channel of images in the Zanzibar dataset, split by pixel classification.

dataset contain vessels as well as humans as a result of high levels of eco-tourism operating in the survey area [33]. This is not the case for the data collection area of NDD, which may lead to a reduction in the false positive rate of the model when evaluated on this dataset. Figure 5 (Right) shows the effect of eco-tourism on the false positive rate, where the model believes a section of the boat's hull and the leg of a human to be a dolphin RoI. Regardless, this evaluation presents evidence that the model is robust enough to deal with data from a different geographic area, time, and cetacean species without the need for re-training or fine-tuning.

### B. Colour Thresholding Mask Components

Experimentation was undertaken to determine the optimal colour threshold value for mask post-processing. Histograms of the RGB colour channel pixel intensities for each object classification in the Zanzibar data were recorded, giving a total of six histograms per image. The histogram groups were then combined to give six global pixel intensity distributions, which can be seen in Figure 6. Regardless of colour channel, there is a near inversion in the distribution of pixel intensities between those detected as dolphin and those not, strongly suggesting it is possible to determine if a component is erroneous based on its colour composition.

Globally, for all masks detected as dolphin 90% of the RGB pixels are below intensities (148, 148, 159). As noise components in the mask are often areas of water or splash, these will be much lighter in composition than cetaceans, and thus can be removed from the mask with confidence by checking the percentage of pixels in the mask component below the threshold.

It was found however that using a 90% threshold when checking mask components at an individual image level was too restrictive, sometimes rejecting valid detections which may have been over-exposed due to lighting conditions. As such, whilst the colour threshold was kept the same, the percentage check was reduced to 50% – providing enough
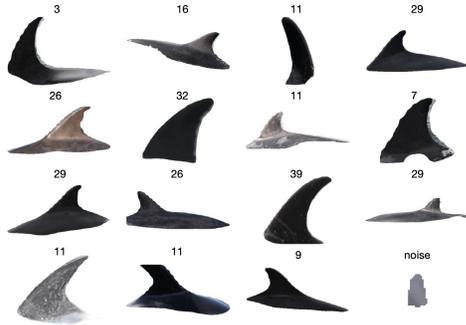
Figure 7. Example post-processed images from NDD used to train the SNN for the task of few-shot, fine-grain most likely photo-id catalogue matching. The class label is displayed above each image. Note the low inter-class and high intra-class differences between non-`noise` classes, such as between examples of class `11`.

leeway for over-exposed but valid detections to be kept whilst still rejecting a large portion of erroneous ones. This was confirmed with further threshold testing on the NDD dataset.

### C. Most Likely Catalogue Matching

**Top-N Accuracy** The SNN model was evaluated against its ability to clearly differentiate between classes in the latent space. To test this, an SNN was trained on the above water set of the NDD dataset after detection and post-processing.

Due to sparse amounts of example images for some individuals, additional photo-id data was provided from catalogues collected in waters around Eastern Scotland maintained by the University of Aberdeen and the University of St. Andrews Sea Mammal Research Unit [3, 4]. Due to the large home range of cetaceans, a 23 individual overlap between this data and the NDD dataset was determined. As a result, 1827 additional images of the overlapping individuals were processed and included in the NDD dataset used to train the SNN. This consisted of 2626 images representing 44 classes including `noise`. Non-`noise` classes (median = 22) contain low inter-class but high intra-class differences between them (example images in Figure 7). This provides a difficult few-shot, fine-grain dataset with which to evaluate the SNN's clustering ability. The dataset was divided using an 80-20 train-test split.

The SNN's backbone followed the structure outlined by Vetrova *et al.* [27], with triplets selected via online semi-hard triplet mining. Hyperparameter optimisation of the SNN was performed using Optuna [34]. During the search, the following possibilities were examined: (1) *Learning Rate* $\in \mathbb{R} \cap [1 \times 10^{-6}, 1 \times 10^{-3}]$, log uniform. (2) *Dropout* $\in \mathbb{R} \cap [0.1, 0.7]$, log uniform. (3) *Kernel Size* $\in \{5, 8\}$. (4) *Triplet Loss Margin* $\in \mathbb{R} \cap [0.1, 1.0]$. (5) *Weight Decay* $\in \mathbb{R} \cap [1 \times 10^{-6}, 1 \times 10^{-1}]$. (6) *Optimiser* $\in \{SGD, Adam\}$. (7) *Step Size* $\in \mathbb{Z} \cap [5, 10]$. (8) $\gamma \in \mathbb{R} \cap [0.1, 0.7]$, log uniform. (9) *Embedding Size* $\in \mathbb{Z} \cap [16, 128]$.

Optimisation of the number of network blocks was also examined. Each block consisted of a Convolutional layer and a MaxPool layer (stride = 2), separated by a ReLU layer and a Dropout layer. During searching, the number of blocks was treated as a hyperparameter optimising for an `int` between 1 and 5 blocks. The initial Convolutional layer size was also

Table III
EFFECT OF DIFFERENT DATA AUGMENTATION STRATEGIES ON SNN TOP-$N$ ACCURACY.

| Data Augmentation Strategy | NDD | | | Naples | | |
|---|---|---|---|---|---|---|
| | Top-10 | Top-5 | Top-1 | Top-10 | Top-5 | Top-1 |
| *Colour Jitter* | 76.83 | 61.18 | 38.82 | 96.25 | **91.25** | 70.00 |
| *Perspective Shift* | 73.58 | 51.22 | 23.17 | 96.25 | **91.25** | **73.75** |
| *Both* | 78.46 | 61.18 | 40.04 | **97.50** | 88.75 | 63.75 |
| *None* | **83.13** | **68.90** | **40.85** | 92.50 | 85.00 | 72.50 |

tuned, searching for an optimal `int` value between 16 and 100. Subsequent layers were double the size of the previous.

The best performing hyperparameters were found to be an initial learning rate of $7.25 \times 10^{-6}$ with weight decay of 0.043 optimised using Adam, a kernel size of 6, triplet loss margin of 0.8, $\gamma$ of 0.012, dropout of 0.17, and an embedding size of 106 (giving a 106-dimensional latent space for most likely matching). The optimal number of network blocks was determined as 2, with an initial Convolutional layer size of 59.

Various data augmentation strategies were examined. These were: (1) *Colour Jitter*: randomly perturb the input images' brightness by a factor between 0.8 and 1.2, contrast by a factor between 0.8 and 1.2, saturation by a factor between 0.9 and 1.1, and hue by a factor between -0.1 and 0.1. (2) *Perspective Shift*: randomly distort the input image's perspective by a factor of 0.5. (3) *Both*: perform both *Colour Jitter* and *Perspective Shift*. (4) *None*: no augmentations.

Prototypes were generated based on the median of the embeddings generated for each class example. Test images were then processed by the SNN and their embedding plotted into the latent space, which was compared using Euclidean distance against the prototypes to generate a list of most likely matches utilised for top-$N$ accuracy evaluation.

Results for the NDD dataset can be seen in Table III (Left). Best model performance was achieved without any data augmentation. These results confirm an SNN trained on automatically processed photo-id catalogue data is capable of accurately providing a list of most likely matches to cetacean researchers, even when trained only using a relatively small number of class examples.

**Generalisability** The generalisability of the SNN approach to most likely catalogue matching was evaluated using a photo-id catalogue subset provided by the Chicago Zoological Society's Sarasota Dolphin Research Program. The subset consisted of 250 images of 23 individual common bottlenose dolphins captured in the waters around Naples, FL, USA [35]. Images were passed through the detector, post-processed, and the generated RoIs used to create a dataset capable of training an SNN for most likely catalogue matching. For consistency, the same architecture and hyperparameters were utilised as with the NDD dataset.

Results for this dataset can be seen in Table III (Right). Unlike training on the NDD data where best results were achieved without any augmentation, here the results are more mixed. Whilst the best top-10 results are obtained using both *Colour Jitter* and *Perspective Shift* augmentations, the best top-5 results were obtained using only one strategy. Using

Figure 8. Example data used to examine the effect of retained background on most likely catalogue matching. Left: Bounding box detection containing both a dorsal fin and background. Centre: Corresponding dorsal fin mask. Right: Corresponding background mask.

*Perspective Shift* only provided the best top-1 accuracy. Whilst this suggests that data augmentation strategy is catalogue dependent, the results confirm that accurate individual level most likely catalogue matching of cetaceans can be performed using SNNs trained on automatically pre-processed data.

### D. Effect of Background on Most Likely Catalogue Matching

Of the four works in Table I which perform dorsal fin detection and downstream individual identification, only half remove all background beforehand. To examine the effect that background removal has on downstream identification, an SNN was trained using the NDD dataset processed into bounding box class examples. Due to the free roaming nature of the individuals, those in the NDD dataset were often photographed only during a single encounter leading to data with small intra-class but high inter-class background variation.

All variables, except the presence of background, were kept consistent with those used when training the best performing SNN on masked data, including model architecture, hyperparameters, and data augmentation strategy. Data was generated using the same Mask R-CNN detector as for previous experiments, modified to output bounding boxes rather than masks.

Analysis of the Euclidean distances between bounding boxed dorsal fins and corresponding fin and background masks show embedding generation is likely to be influenced more by features in the background than the fin. For example, the Euclidean distance between the bounding box data in Figure 8 (Left) and its corresponding dorsal fin mask (Centre) is 0.36, compared to a distance of 0.30 between the bounding box and the background mask (Right) and a mean distance of 0.97 between the bounding box and generated class prototypes. This suggests the SNN is performing likely matching based on features found in the background rather than on the dorsal fins, reflected in increased model performance whereby using bounding box data to train the SNN sees an increase of 22.94% top-1, 15.58% top-5, and 6.53% top-10 accuracies over using masked data.

By removing all background, the masked SNN is prevented from utilising environmental conditions to aid matching. This finding raises important questions regarding the performance of photo-id aids which do not remove all background before performing matching. If the photo-id catalogue utilised for system evaluation has been collected over a small temporal scale, then results obtained in this experiment suggest that performance may be artificially inflated by the retention of feature heavy background. Further studies will examine the effect of background retention on most likely matching to

catalogues gathered over a large temporal scale, as well as the use of out of distribution negative samples [36] to train networks robust to the issue of consistent background.

## V. LIMITATIONS

One limitation of the system currently is the need to retrain the SNN for each photo-id catalogue. As a result, initial manual curation must be performed before the methodology can be applied. The feasibility of a more general SNN capable of catalogue-agnostic photo-id will be examined in future work. This limitation does not apply to the Mask R-CNN however, which has been found not to require re-training when applied to a new photo-id catalogue regardless of changes in species, geography, or time.

Further, whilst the system has been shown to be robust enough to deal with multiple cetacean species, these have all been dolphins. It is not clear how well the pipeline would perform with cetacean species such as whales or porpoises, or with body parts like flukes. Further studies with catalogues of other species will be explored.

## VI. CONCLUSION

This work examines the use of a pipeline of computer vision models to aid researchers through the curation of photo-id data. The system is capable of operating on raw field images, no pre-processing required, thanks to the use of a dorsal fin detection model. Evaluation of this model shows it is capable of achieving high mAP on photo-id catalogues containing data from different species, collected in different geographical locations, and at different times to the catalogue on which it is trained. Drops in performance are observed when operating over data containing examples of eco-tourism, and further work will examine how best to mitigate this.

Detections are outputted as pixel wise masks, post-processed to improve the chance of catalogue matching. Experimentation to locate the optimal colour threshold confirms erroneous detections can be filtered out with confidence whilst over-exposed fins are retained.

Outputs are passed to an SNN, trained for the task of most likely catalogue matching, to generate an embedding which is plotted into a latent space. Embeddings are compared to class prototypes using Euclidean distance to generate matches. Evaluation of SNNs trained on processed photo-id data suggests that they are a viable approach to the problem of most likely catalogue matching, and are capable of flagging detections which may be of previously uncatalogued individuals. Experimental results studying the effect of retained background suggest this can negatively impact embedding generation, especially for catalogues collected over a small temporal scale. The use of detection masks negates this effect, preventing embedding generation interference.

## REFERENCES

[1] B. G. Weinstein, "A computer vision for animal ecology," *Journal of Animal Ecology*, vol. 87, no. 3, pp. 533–545, 2018.

[2] M. Sharpe and P. Berggren, "Indian Ocean humpback dolphin in the Menai Bay off the south coast of Zanzibar, East Africa

is Critically Endangered," *Aquatic Conservation: Marine and Freshwater Ecosystems*, vol. 29, no. 12, pp. 2133–2146, 2019.

[3] M. Arso Civil, N. J. Quick, B. Cheney, E. Pirotta, P. M. Thompson, and P. S. Hammond, "Changing distribution of the east coast of Scotland bottlenose dolphin population and the challenges of area-based management," *Aquatic Conservation: Marine and Freshwater Ecosystems*, vol. 29, pp. 178–196, 2019.

[4] B. Cheney, R. Corkrey, J. W. Durban, K. Grellier, P. S. Hammond, V. Islas-Villanueva, V. M. Janik, S. M. Lusseau, K. M. Parsons, N. J. Quick, B. Wilson, and P. M. Thompson, "Longterm trends in the use of a protected area by small cetaceans in relation to changes in population status," *Global Ecology and Conservation*, vol. 2, pp. 118–128, Dec. 2014.

[5] P. S. Hammond, S. A. Mizroch, and G. P. Donovan, "Individual recognition of cetaceans: use of photo-identification and other techniques to estimate population parameters: incorporating the proceedings of the symposium and workshop on individual recognition and the estimation of cetacean population parameters," *International Whaling Commission*, vol. 12, 1990.

[6] R. Payne, "Long term behavioral studies of the southern right whale (Eubalaena australis)," International Whaling Commission, Tech. Rep. 10, 1986.

[7] K. A. Forney and J. Barlow, "Seasonal Patterns in the Abundance and Distribution of California Cetaceans, 1991–1992," *Marine Mammal Science*, vol. 14, no. 3, pp. 460–489, 1998.

[8] L. J. Feyrer, M. Stewart, J. Yeung, C. Soulier, and H. Whitehead, "Origin and Persistence of Markings in a Long-Term Photo-Identification Dataset Reveal the Threat of Entanglement for Endangered Northern Bottlenose Whales (Hyperoodon ampullatus)," *Frontiers in Marine Science*, vol. 8, 2021.

[9] M. Bigg, "An Assessment of Killer Whale (Orcinus orca) Stocks off Vancouver Island, British Columbia," *Report of the International Whaling Commission*, vol. 32, no. 65, p. 12, 1982.

[10] C. Trotter, G. Atkinson, M. Sharpe, K. Richardson, A. S. McGough, N. Wright, B. Burville, and P. Berggren, "NDD20: A large-scale few-shot dolphin dataset for coarse and fine-grained categorisation," *arXiv:2005.13359 [cs]*, May 2020.

[11] R. B. Tyson Moore, K. W. Urian, J. B. Allen, C. Cush, J. R. Parham, D. Blount, J. Holmberg, J. W. Thompson, and R. S. Wells, "Rise of the Machines: Best Practices and Experimental Evaluation of Computer-Assisted Dorsal Fin Image Matching Systems for Bottlenose Dolphins," *Frontiers in Marine Science*, vol. 9, p. 849813, Apr. 2022.

[12] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," *arXiv:1703.06870 [cs]*, 2017.

[13] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 815–823.

[14] G. Hillman, N. Kehtarnavaz, B. Wursig, B. Araabi, G. Gailey, D. Weller, S. Mandava, and H. Tagare, ""Finscan", a computer system for photographic identification of marine animals," in *Engineering in Medicine and Biology*, vol. 2, 2002.

[15] N. Fisheries, "FinBase Photo-Identification Database System | NOAA Fisheries," Feb. 2018.

[16] S. A. Hale, "Unsupervised Threshold for Automatic Extraction of Dolphin Dorsal Fin Outlines from Digital Photographs in DARWIN (Digital Analysis and Recognition of Whale Images on a Network)," *arXiv:1202.4107 [cs]*, Feb. 2012.

[17] E. M. Keen, J. Wren, É. O'Mahony, and J. Wray, "catRlog: a photo-identification project management system based in R," *Mammalian Biology*, Aug. 2021.

[18] H. J. Weideman, Z. M. Jablons, J. Holmberg, K. Flynn, J. Calambokidis, R. B. Tyson, J. B. Allen, R. S. Wells, K. Hupman, K. Urian, and C. V. Stewart, "Integral Curvature Representation and Matching Algorithms for Identification of Dolphins and Whales," *arXiv:1708.07785 [cs]*, 2017.

[19] J. Karnowski, E. Hutchins, and C. Johnson, "Dolphin Detection and Tracking," in *IEEE Winter Applications and Computer Vision Workshops*, Waikoloa, HI, USA, 2015, pp. 51–56.

[20] Y. Quiñonez, O. Zatarain, C. Lizarraga, and J. Peraza, "Using Convolutional Neural Networks to Recognition of Dolphin Images," in *Trends and Applications in Software Engineering*, 2019, pp. 236–245.

[21] E. Morteo, A. Rocha-Olivares, R. Morteo, and D. W. Weller, "Phenotypic variation in dorsal fin morphology of coastal bottlenose dolphins (Tursiops truncatus) off Mexico," *PeerJ*, 2017.

[22] S. Bouma, M. D. Pawley, K. Hupman, and A. Gilman, "Individual Common Dolphin Identification Via Metric Embedding Learning," in *International Conference on Image and Vision Computing New Zealand (IVCNZ)*, 2018, pp. 1–6.

[23] Y.-C. Lee, H.-W. Hsu, J.-J. Ding, W. Hou, L.-S. Chou, and R. Y. Chang, "Backbone Alignment and Cascade Tiny Object Detecting Techniques for Dolphin Detection and Classification," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, 2020.

[24] J. W. Thompson, V. H. Zero, L. H. Schwacke, T. R. Speakman, B. M. Quigley, J. S. Morey, and T. L. McDonald, "finFindR: Automated recognition and identification of marine mammal dorsal fins using residual convolutional neural networks," *Marine Mammal Science*, vol. 38, no. 1, pp. 139–150, 2022.

[25] R. Maglietta, V. Renò, G. Cipriano, C. Fanizza, A. Milella, E. Stella, and R. Carlucci, "DolFin: an innovative digital platform for studying Risso's dolphins in the Northern Ionian Sea (North-eastern Central Mediterranean)," *Scientific Reports*, vol. 8, no. 1, p. 17185, Nov. 2018.

[26] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-Up Robust Features (SURF)," *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, Jun. 2008.

[27] V. Vetrova, S. Coup, E. Frank, and M. J. Cree, "Hidden Features: Experiments with Feature Transfer for Fine-Grained Multi-Class and One-Class Image Categorization," in *International Conference on Image and Vision Computing New Zealand (IVCNZ)*. Auckland, New Zealand: IEEE, Nov. 2018, pp. 1–6.

[28] V. M. Araújo, A. S. Britto Jr., L. S. Oliveira, and A. L. Koerich, "Two-view fine-grained classification of plant species," *Neurocomputing*, vol. 467, pp. 427–441, Jan. 2022.

[29] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *arXiv:1512.03385 [cs]*, Dec. 2015.

[30] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980 [cs]*, 2014.

[31] I. Loshchilov and F. Hutter, "SGDR: Stochastic Gradient Descent with Warm Restarts," *arXiv:1608.03983 [cs, math]*, 2016.

[32] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft COCO: Common Objects in Context," *arXiv:1405.0312 [cs]*, 2014.

[33] F. Christiansen, D. Lusseau, E. Stensland, and P. Berggren, "Effects of tourist boats on the behaviour of Indo-Pacific bottlenose dolphins off the south coast of Zanzibar," *Endangered Species Research*, vol. 11, pp. 91–99, Mar. 2010.

[34] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A Next-generation Hyperparameter Optimization Framework," Jul. 2019, arXiv:1907.10902 [cs, stat].

[35] R. B. Tyson Moore, A. Barleycorn, C. Cush, A. Honaker, S. McBride, C. Toms, and R. Wells, "Final Report: Abundance and distribution of common bottlenose dolphins (Tursiops truncatus) near Naples and Marco Island, Florida, USA, 2018-2019," The Batchelor Foundation, Tech. Rep., 2020.

[36] J. Lee, S. J. Oh, S. Yun, J. Choe, E. Kim, and S. Yoon, "Weakly Supervised Semantic Segmentation Using Out-of-Distribution Data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, New Orleans, LA, USA, 2022.