

**IDENTIFYING HUMPBACK WHALE FLUKES BY
SEQUENCE MATCHING OF TRAILING EDGE
CURVATURE**

By

Zachary Jablons

A Thesis Submitted to the Graduate
Faculty of Rensselaer Polytechnic Institute
in Partial Fulfillment of the
Requirements for the Degree of
MASTER OF SCIENCE
Major Subject: COMPUTER SCIENCE

Examining Committee:

Dr. Charles Stewart, Thesis Adviser

Dr. Barbara Cutler, Member

Dr. Bülent Yener, Member

Rensselaer Polytechnic Institute
Troy, New York

April 2016
(For Graduation May 2016)

© Copyright 2016
by
Zachary Jablons
All Rights Reserved

CONTENTS

LIST OF TABLES	vi
LIST OF FIGURES	vii
ACKNOWLEDGMENTS	x
ABSTRACT	x
1. Introduction	1
1.1 Humpback Whales	1
1.1.1 Distinguishing Individual Flukes	1
1.2 Current Identification Methods	1
1.2.1 Based on Trailing Edge	2
1.2.2 Based on General Fluke Appearance	4
1.3 Method Outline	5
1.4 The Dataset	6
1.5 Thesis Outline	6
2. Background	8
2.1 Convolutional Networks	8
2.1.1 Facial Keypoint Prediction	9
2.1.2 Fully Convolutional Networks	9
2.2 Contour Extraction	10
2.2.1 Active Contour	10
2.2.2 Seam Carving	10
2.3 Curvature Measures	10
2.4 Dynamic Time Warping	11
3. Methods	13
3.1 Trailing Edge Extraction	13
3.1.1 Fluke keypoint prediction	13
3.1.1.1 Network Design	14
3.1.1.2 Training Details	14
3.1.1.3 Evaluation	15
3.1.2 Basic Trailing Edge Extraction Algorithm	16

3.1.3	Trailing Edge Scoring	19
3.1.3.1	Trailing Edge Scoring Architectures	20
3.1.3.2	Using the trailing edge scores	23
3.1.3.3	Training Details	24
3.1.3.4	Evaluation	25
3.2	Trailing Edge Matching	26
3.2.1	Curvature Measurement	26
3.2.2	Sequence Matching	28
3.3	Alternative Approaches	29
3.3.1	Aligning Trailing Edges	29
3.3.1.1	Keypoint Alignment	29
3.3.1.2	Dynamic Time Warping Alignment	31
3.3.2	Histogram Matching	31
3.3.3	Embedding via Convolutional Networks	31
4.	Results	33
4.1	Main method	33
4.2	Configuration Options	33
4.2.1	Variability in Matching Score	33
4.2.2	Effectiveness of Keypoint Extractor	35
4.2.3	Cropping Width	35
4.2.4	Trailing Edge Extraction	37
4.2.4.1	Trailing Edge Scorer variations	37
4.2.4.2	Using a Trailing Edge Scorer	38
4.2.4.3	Number of neighbors in the extraction	40
4.2.4.4	Using the notch	41
4.2.5	Curvature Scales	41
4.2.6	Sakoe-Chiba bound	43
4.3	In Combination with Hotspotter	44
4.3.1	Characterization of when to use which method	45
5.	Discussion	47
5.1	Issues with the Proposed Method	47
5.2	Future work	48
5.3	Conclusion	49

REFERENCES 50

APPENDIX

LIST OF TABLES

- 3.1 Table showing the precision, recall, and IoU of each of the evaluated trailing edge scorers on each section of the trailing edge dataset. For the purposes of this analysis, we use the `argmax` over the classes to determine a positive (i.e. trailing edge) or negative pixel. 26

LIST OF FIGURES

1.1	Example Flukes. Example images of humpback whale flukes from the SPLASH [9] dataset. These flukes both have distinctive internal textures (more so on the left). However, the trailing edge on the left is far more distinctive than the trailing edge on the right.	2
1.2	Uniform Internal Texture. This image of a humpback fluke shows no clear internal texture, but a distinctive trailing edge.	3
1.3	Change in Trailing Edge. The above images show that out of plane rotations of the fluke can obscure it or otherwise make it hard to match. These images are both of the same individual, however in the top image the fluke is rotated slightly towards the camera.	4
1.4	Example Keypoint and Trailing Edge Annotation. This image shows a typical set of fluke keypoints and trailing edge extracted by our algorithm.	5
3.1	Example Keypoint Prediction. Example image showing the left tip, bottom of the notch, and right tip located by the keypoint extractor convolutional network.	14
3.2	Example Keypoint Failure. Example image showing a keypoint extraction failure case from its testing set. Note the difference in pose of the fluke from the success case shown in Figure 3.1. This is an example of a fluke image that violates our assumptions.	16
3.3	Histogram of Keypoint Distances. This is a histogram of the average distance from predicted keypoints to annotated keypoints on the testing set for the keypoint extraction network. The vast majority of keypoints are predicted within 10 pixels of the true keypoints.	17
3.4	Example Trailing Edge Extraction. Example of the baseline trailing edge extraction with $n = 2$. Note that the gradient image has a significant black area where the trailing edge is, making this an easy case. 18	18
3.5	Example Trailing Edge Score. Bottom image is the Residual scorer's classification of the top image. Trailing edge is class is colored black. . .	21
3.6	Trailing Edge Scores. These are the trailing edge scores given by each of the networks described in section 3.1.3.1 on the image used in Figure 3.5. A darker pixel is predicted to be part of the trailing edge by the network.	23

3.7	Trailing Edge Scoring Failure. Unfortunately even the Residual architecture is still imperfect, and can lead to catastrophic trailing edge failures like this one. However this is a rare case.	25
3.8	Trailing Edge Curvature. The top images are of the same individual, and the bottom images visualize the corresponding curvatures for the trailing edges that were extracted. Each row in the visualization is a curvature scale, increasing from top to bottom. Note that darker blue implies a “valley” in the trailing edge, whereas lighter blue implies a “peak”.	28
3.9	Example Matches. The left side shows a success case, and the right side shows a failure case.	30
4.1	Score Separability Histogram. The blue bars in this figure represent true matches, and the red bars represent false matches. The line at score = 0.88 represents the optimal threshold at which to accept a match, although we can see it is not perfect.	34
4.2	Varying Manual Extraction. There is a small difference in matching accuracy between using the manually annotated points (red) provided for this dataset versus the keypoint extractor’s predicted points (cyan). The bottom of the notch keypoint is not used in either of these evaluations.	35
4.3	Distribution of Unresized Trailing Edge Lengths. This shows a significant distribution of trailing edges centered around a width of 800 pixels.	36
4.4	Varying w. Note that we use the manually annotated points in this analysis to control for any issues with keypoint extraction.	37
4.5	Trailing Edge Scorer Architectures. The highest performing trailing edge scorer (Residual) is shown in red, followed by Simple, Jet, and Upsample (in descending order of accuracy).	38
4.6	Trailing Edge Scorer Architectures at $\beta = 1$. Upsample (cyan) performs significantly worse than the other networks, which all perform comparably.	39
4.7	Varying β.	39
4.8	Example Use of Trailing Edge Scorer. In (a), we have the trailing edge extracted with the Residual scorer. Compare to (b), which did not use any scorer at all, resulting in a match failure.	40
4.9	Varying n. This shows that the optimal neighborhood constraint is $n = 3$, despite qualitatively producing worse-looking trailing edges. Beyond $n = 5$, the trailing edges can become very noisy.	41

4.10	Using the notch as a control point.	42
4.11	Curvature Diversity. Left panel (a) shows the average standard deviation of the (fixed length) curvature at different scales. Right panel (b) shows the average Euclidean distance between successive scales of curvature	42
4.12	Varying Sakoe-Chiba bound.	43
4.13	Example Disagreements Between Hotspotter and our method. On the left side, (a) was matched correctly to (c) by our method, whereas Hotspotter could not find any matches for (a). On the right hand side however, Hotspotter was able to match two flukes with a large variance in pose and lighting, while our method did not rank (d) in the top 5 matches for (b).	44
4.14	Comparison between Hotspotter and our method.	45
5.1	Histogram of Ground-Truth Ranks. Note that the histogram ranges are uneven to better show the lower end of the range. In order to have all matches found within the top- k matches we would have to set $k = 414$	48

ACKNOWLEDGMENTS

This project could not have been completed so quickly and thoroughly without the invaluable help from the rest of the IBEIS team. I would like to thank Jon Crall for helping to write the experimentation framework which has saved a lot of time and effort. I would also like to thank Jason Parham and Hendrik Weidemann for the great discussions and advice, as well as the handy L^AT_EX template. Additionally, without the development and annotation efforts of Andrew Batbouda, a lot of models could not have been so successfully trained. Importantly, I would like to thank the Wildbook team (Jason Holmberg, Jon van Oast) for providing the dataset and corresponding annotations with which a lot of models were trained and experiments run. Of course, this work could never have been undertaken without the help, advice, and direction of my advisor, Charles Stewart.

ABSTRACT

Photographic identification of humpback whale (*Megaptera novaeangliae*) flukes (i.e. their tail) is an important task in marine ecology, and is used in tracking migration patterns and estimating populations [7, 9]. In this thesis, we lay out a method that automates the photo-identification of humpback flukes using the “trailing edge” of the fluke. The method uses convolutional networks to identify keypoints on the fluke and possible trailing edge locations. It then uses this information to extract a detailed trailing edge and its curvature, which is then matched to other trailing edges in the database via dynamic time warping. Using this method, we achieve nearly 80% top-1 ranking accuracy on a large subset of the SPLASH [9] dataset consisting of about 400 identified individuals. We also show that in combination with Hotspotter [12], a general pattern based matching algorithm, we can achieve 93% accuracy.

To our knowledge, this is the first method that can achieve this level of accuracy on humpback fluke identification without extensive manual effort at inference time.

CHAPTER 1

Introduction

1.1 Humpback Whales

Since the international ban on commercial hunting of Humpback whales in 1966, humpback whales have grown from a population of only 5000 [4] to over 50000 [8]. As the population grows, it becomes more and more important to be able to automatically identify individual whales in order to accurately monitor their population growth and follow their migration patterns, among other ecological conservation endeavours. One of the most reliable methods for photo-identifying humpback whales is by taking pictures of their flukes as they dive after breaching the surface of the water.

1.1.1 Distinguishing Individual Flukes

The primary distinguishing features of these flukes are — for the purposes of this work — separated into two main areas; the “trailing edge” of the fluke, and the internal texture (see Figure 1.1). The internal texture is a more obvious choice for identification, as it is distinctive even from a distance even when the image is blurred. Unfortunately, some humpback flukes have indistinct (e.g. all-black) internal textures, which make matching based on texture impossible (see Figure 1.2). Additionally, the work of Blackmer et al. [7] finds that the trailing edge changes less with age than the internal texture of the fluke, which means that it can (potentially) be a more reliable identifier over time. That said, the requirements for getting a good photograph of the trailing edge can be impractical, as the trailing edge is vulnerable to being obscured by out of plane rotations (see Figure 1.3), and as such are optimally photographed when the fluke is “flat” to the photographer.

1.2 Current Identification Methods

Computer-assisted photo-identification of humpback whale flukes has been attempted since the early 90s [36]. While early efforts mostly relied on a manual de-

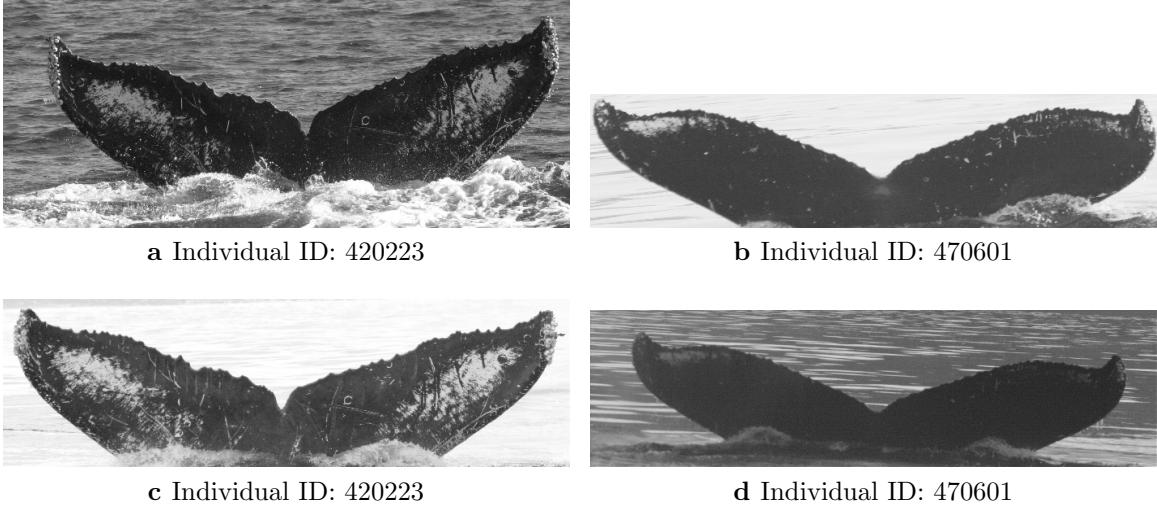


Figure 1.1: Example Flukes. Example images of humpback whale flukes from the SPLASH [9] dataset. These flukes both have distinctive internal textures (more so on the left). However, the trailing edge on the left is far more distinctive than the trailing edge on the right.

scription of the fluke that would then be matched against other stored descriptions [36, 55], later efforts have involved matching flukes based on automated analysis of both the internal texture and trailing edge [25, 30, 20].

Existing computer-assisted photo-identification methods can be broadly separated into three categories. There are manual methods, in which humans both identify and describe distinguishing features, semi-automated methods, in which humans identify distinguishing features that are then described and matched automatically, and automated methods — which can match based solely on raw images with no human involvement.

1.2.1 Based on Trailing Edge

In the I³S contour system [20], the user must input start and end points on the query trailing edge, after which its contour is extracted. This trailing edge is then resized and aligned so that it can be compared with absolute difference against the database trailing edges. It also compares a set of possible shifts, rotations, and scales of the query trailing edge to account for these differences. At the time of writing no published results on this system applied to humpback whales could be



Figure 1.2: Uniform Internal Texture. This image of a humpback fluke shows no clear internal texture, but a distinctive trailing edge.

found.

Automatically identifying humpback whales by their entire trailing edge contour is done experimentally in Hughes et al. [25], using a technique that is originally designed for great white sharks. This technique segments the trailing edge into a set of possible contours and matches them combinatorially using Difference of Gaussians. The authors achieve a comparable accuracy to our method, however for a much smaller dataset of humpback flukes than the one evaluated here.

While trailing edge matching has seen limited use in humpback whale identification, it is a much more common technique in sperm whale (*P. macrocephalus*) identification [24, 5, 55] — with varying levels of manual effort. In Whitehead’s work [55], points of interest on the trailing edge are entered and catalogued manually along with their positions. In order to match these trailing edges, all of the points are compared against points on annotated trailing edges in the database, using a distance threshold to ensure locality. This is a manual method, requiring extensive annotation for matching.

The method proposed by Huele et al. [24] uses a semi-automatic extraction of sperm whale trailing edge, and then applies wavelet transformations which are cross-correlated to determine a similarity measurement.

An investigation of the above methods for sperm whale identification is carried out in the work of Beekmans et al. [5], showing that combining these two methods yields an 80% top-1 accuracy (meaning that the correct identity is ranked at the top of the potential matches 82% of the time) on a slightly smaller dataset than ours.

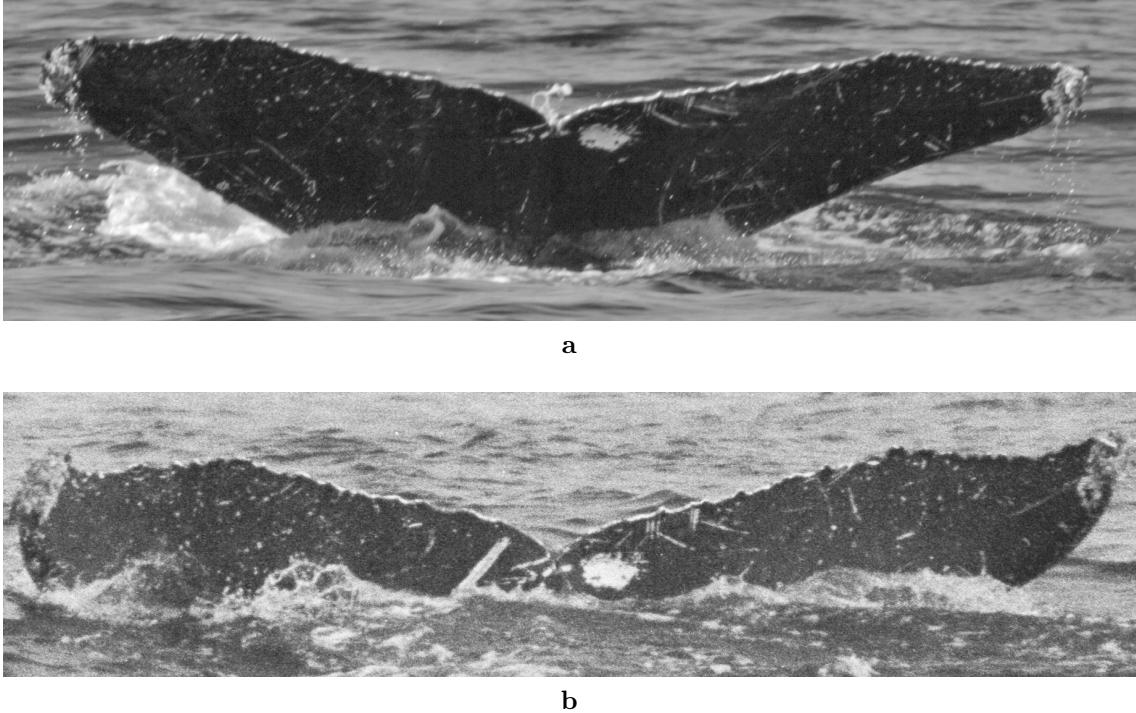


Figure 1.3: Change in Trailing Edge. The above images show that out of plane rotations of the fluke can obscure it or otherwise make it hard to match. These images are both of the same individual, however in the top image the fluke is rotated slightly towards the camera.

However on their own each method only achieves about 65% top-1 accuracy.

1.2.2 Based on General Fluke Appearance

The primary paradigm for computer-assisted photo-identification of humpback whale flukes is to use the internal fluke pattern, as seen in the work of Mizroch et al. [36] and the Flukematcher [30] of Kniest et al¹.

In Mizroch et al., information about the fluke is manually catalogued and used to match individual whales. The fields that are catalogued contain information primarily about the overall coloration patterns of the fluke, as well as the shape of the central notch. The matching algorithm generates potential fluke matches by looking at how similar the annotated patterns are. This requires significant manual effort to identify individuals.

¹Flukematcher also uses information about the trailing edge, but is focused on the internal fluke texture

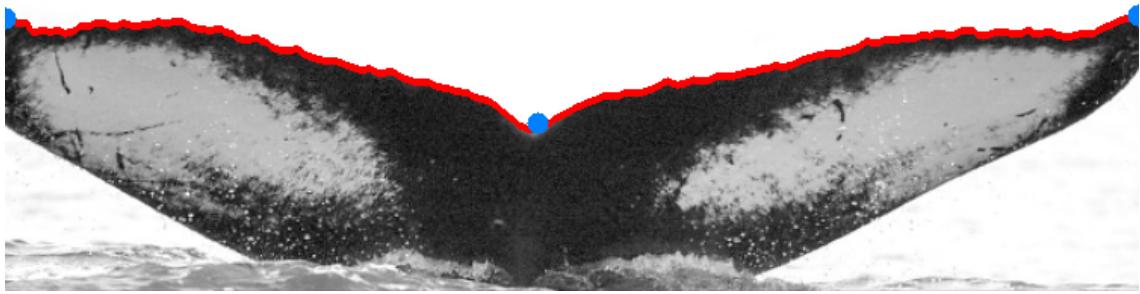


Figure 1.4: Example Keypoint and Trailing Edge Annotation. This image shows a typical set of fluke keypoints and trailing edge extracted by our algorithm.

In Flukematcher [30], control points are manually annotated which allow the program to automatically find pigmentation patterns in the fluke and align accordingly. Optionally, distinctive fluke patterns can also be selected and annotated by the user. A variety of heuristic features are then extracted, which are matched using a variety of similarity measures. This method has achieved a 82% top-1 accuracy on a smaller dataset than ours. However, it requires significant manual effort on the order of five minutes per fluke photograph.

1.3 Method Outline

In this thesis, we develop and present an efficient fully-automated² algorithm that identifies humpback flukes based on their trailing edges. For each fluke, the algorithm first determines the left and right tip points of the fluke (as well as the bottom of the central notch). The trailing edge contour is then extracted between the left and right tip points (see Figure 1.4), and for each point on this contour curvature is measured at multiple scales. Once these curvatures are extracted for each fluke photograph, the algorithm computes the distance between a query fluke and possible identities by aligning its contour curvature against contour curvatures in the database and computing distance with dynamic time warping. The query fluke is then given the identity of the database fluke with the lowest distance. This method is based on both traditional edge detection and curvature techniques, modern machine learning, and classical fast sequence comparison.

We also show we can greatly improve the accuracy of this method by combining

²With the caveat that manual annotation is needed to train parts of the algorithm

it with matches found by Hotspotter, a generalized pattern based identification method that is a powerful matching algorithm for several species [12]. Hotspotter extracts keypoints with SIFT descriptors in an image that are then spatially verified and matched with other keypoints in the database to produce a ranking over possible identities. Despite the general nature of this method, we find that it does not identify keypoints on the trailing edge and thus struggles with flukes that have no significant internal texture. This work is the first to our knowledge that details Hotspotter's efficacy when applied to humpback whale flukes, and the results are presented in Chapter 4.

1.4 The Dataset

The main dataset that is used and evaluated in this work is a subset of the dataset collected by the SPLASH project [9]. It consists of about 1400 identified photographs spread over about 860 identified individuals. Of these, only 433 individuals have more than one image associated with them, giving 942 images that can be used in a one-to-one comparison. We refer to this dataset as the Flukebook dataset, which is the team from which it originates. All images shown in this thesis come from this dataset.

Additionally, an external dataset of unidentified (but annotated) humpback flukes is used for training individual components of the method.

1.5 Thesis Outline

The rest of this thesis starts by giving a background on the algorithms that our method is based on in Chapter 2. Afterwards, the main method is detailed in Chapter 3, along with a description of some alternate methods that failed to work as well. Chapter 4 goes into the results of the primary algorithm as well as various configuration changes and how they affect the results. We also describe the effectiveness of this identification method in combination with Hotspotter. This thesis concludes with a discussion on the failings of the primary method, as well as ways to improve it and generalize it to extracting and matching edge characteristics in other animals.

The primary contributions of this work are the individual components of the main method for extracting trailing edges and fluke keypoints, as well as the combination of all the components into a coherent identification pipeline. We also contribute an evaluation of dynamic time warping on curvature measures as a method for matching humpback flukes, as well as an evaluation of Hotspotter for this task.

CHAPTER 2

Background

In this chapter, we provide a series of sections detailing background information on the algorithms on top of which our trailing edge identifier was developed. We describe some of the applications and variants of deep convolutional networks on top of which we build our fluke keypoint and trailing edge extractors. We briefly introduce the concept of seam carving as it relates to our trailing edge extraction algorithm, and provide a small overview of contour curvature measures. Additionally, we describe briefly the concept of dynamic time warping for sequence matching.

2.1 Convolutional Networks

In recent years, convolutional neural networks have provided state of the art results in several challenging computer vision tasks, including general image classification [31, 53], image segmentation [35, 10] and individual identification (specifically for human faces) [15, 48].

The essential idea of a convolutional network is that if we can use the gradient of an error signal to learn series of convolution kernels separated by nonlinear activation functions. These networks are a type of neural network, and are used in image data as convolutions provide a meaningful prior when applied to data with spatial relationships. Convolutional networks were introduced nearly 40 years ago by Fukushima in [17]. The current incarnation of these networks can be traced back to the seminal work of LeCun et al. [33]. Modern convolutional networks follow a common framework of using Rectified Linear Units (ReLUs) as activation functions, and Dropout [22] layers after fully connected layers for regularization. Additionally, the convolutional kernels used are commonly small square kernels with “same” padding³ alternated with $2 \times$ downsampling layers (specifically max pooling layers) [50, 49, 31].

³Essentially zero-padding such that the output image has the same shape as the input image. For more information on this see [14]

The convolutional networks used in this thesis follow the above framework, and also use batch normalization [26] at every layer. Every network in this thesis uses orthogonal initialization [47] at all layers to help ensure gradient flow.

2.1.1 Facial Keypoint Prediction

Facial keypoints are essentially coordinates on an image of a (human) face that detail the locations of the nose, eyes, mouth, etc. They are commonly used in facial identification pipelines [54], as well as in motion capture [1] and expression recognition [6]. We draw on recent work for using convolutional networks for facial keypoint prediction [52, 41]. The essential idea behind these networks is that they predict points (rather than classifications) in the form of (x, y) coordinates, and are trained with a regression loss function. In this work, we adjust these methods to predict fluke keypoints marking the left and right tips of the trailing edge using essentially the same paradigm as the above works.

2.1.2 Fully Convolutional Networks

Classically, convolutional networks reduce an image to a single (spatially invariant) vector, which is then used for classification (or embedding, regression, etc.) To do this, these networks usually have fully connected (or dense) layers towards the end. This ensures that the receptive field of the network covers the entire image, which is practical for many applications where a scalar or fixed size prediction is required. However, when arbitrarily sized predictions are required (e.g. for semantic segmentation) from arbitrarily sized images, it is useful to use networks that are “fully convolutional”, in which case the entire network consists of convolution kernels. Convolutional networks that reduce to dense layers can be cast as fully convolutional networks by replacing the dense layers with 1×1 kernels, using the dense units as channels. This technique is especially applicable in segmentation tasks [40, 11, 19], as this process allows for (downsampled) predictions that adhere spatially to the input image. By then upsampling and combining different stages of prediction, the authors in [35] produce high quality image segmentations, a technique that we replicate for classifying pixels as being part of the trailing edge. In this

work, fully convolutional networks are used for predicting the “trailing edginess” of an image, which allows us to refine the trailing edge contour extraction.

2.2 Contour Extraction

2.2.1 Active Contour

Automatically extracting contours from edges in images is an old technique, and one of the primary methods for getting coherent contours from image information is the active contour (or snakes) method [2, 28]. This technique explicitly models the contour as a function that minimizes a combination of curvature, smoothness, and image edge-ness constraints, weighted appropriately. Often these contours are fairly smooth, as the constraints imposed require a continuous function to represent the contour. Additionally, this is an iterative process that is subject to local minima, meaning that it can produce different results based on the initialization in the image. For these reasons, we did not investigate using active contours for trailing edge extraction.

2.2.2 Seam Carving

Seam carving is a technique that tries to resize images without warping or distorting the objects shown in the image [3]. This technique uses a dynamic programming algorithm to find minimal salience paths through an image, where salience is often defined as the gradient. The motivation for this is that these minimal saliency paths are not important to the image, so they can be removed to reduce its size. While this method is not directly used in this work, the underlying dynamic programming algorithm for trailing edge extraction is essentially a single iteration of the seam carving algorithm, using gradient information.

2.3 Curvature Measures

Contour curvature measures are commonly used to characterize the overall shape of an contour. A significant amount of work has been done on using curvature information for detection [37] and classification [16, 32]. This curvature information

can be broadly broken down into either integral or differential curvature, and is usually computed at multiple scales.

Differential curvature can generally be seen as measuring the angle of the normal to the gradient at each point in an image [16]. The curvature of a contour is then expressed by using the points on that contour. While doing this directly can be fast to compute, it tends to be noise sensitive and we found that integral curvature (below) works better for our purposes.

Integral curvature works (conceptually) by sliding a circle of some radius r along the contour [43], and measuring how much of the circle is “inside” the contour. This measurement is usually taken at multiple scales, and has the appealing property of being invariant to rotation and translation (of the entire contour). In this work, we approximate the circular curvature with a square of size r , which appears to perform just as well but can be computed much faster.

2.4 Dynamic Time Warping

In deciding a sequence comparator, one criterion that is often important is ensuring that small shifts in the sequence do not balloon into large differences. Dynamic Time Warping (DTW) is a sequence comparison method that, roughly, finds the optimal matching between all sets of points in the two given sequences that minimizes the overall distance (for some defined distance function) between the matched points, while keeping the locality of the points intact [45]. This allows for shifts and some warps in the two sequences to be compensated for, and results in a nonlinear mapping of one sequence onto another. The algorithmic complexity of dynamic time warping can be limiting in large datasets, as it is quadratic in both space and time – making a linear scan of each database fluke to identify a query fluke a bit inefficient.

There are several variants on DTW that give faster speeds [46, 34], however we only use the Sakoe-Chiba bound [45], which constrains the neighborhood in which points can be matched. This bound gives a complexity of $O(nT)$, where T is a user set parameter constraining the neighborhood.

Sequences of curvature measures have been used with DTW for signature

verification [39], however this combination has not been used for matching trailing edges to our knowledge.

CHAPTER 3

Methods

In this chapter, we describe in detail the humpback fluke trailing edge identification algorithm pipeline. We also briefly discuss some alternative approaches that we found to be unsuccessful.

3.1 Trailing Edge Extraction

Extracting good, high quality trailing edges images is one of the primary challenges when matching Humpback whales by their trailing edge. In this section, we describe the steps that go into automating the extraction of high quality trailing edges.

One major assumption we make when extracting these trailing edges is that the humpback whale fluke is aligned such that its major axis is horizontal. Additionally, we generally assume that all parts of the fluke are present. The nature of the problem (and the Flukebook dataset) makes these assumptions reasonable. It would also be possible to add a detection and orientation prediction step beforehand to ensure this assumption, however we do not explore this in this work.

3.1.1 Fluke keypoint prediction

When extracting trailing edges, one of the first steps we take is to identify the starting and ending points of the trailing edge, as well as the bottom of the central notch (although we do not necessarily use it). To do this, we train a convolutional network to predict these three points.

This network can use arbitrarily resized images, so the first step of the keypoint extraction pipeline is to resize the image to 128×128 pixels. This size choice is somewhat arbitrary, but we find that it provides strong performance without using an unnecessary amount of memory. The network predicts the points as values between 0 and 1, essentially giving coordinates into the image as percentages of its height and width. These predictions are then rescaled back up to the original



Figure 3.1: Example Keypoint Prediction. Example image showing the left tip, bottom of the notch, and right tip located by the keypoint extractor convolutional network.

image size by multiplying with the original image height and width. An example prediction is shown in Figure 3.1.

3.1.1.1 Network Design

The overall design of the network follows the pattern of alternating small (3×3) convolutional filters with 2×2 max pooling layers, at each step doubling the number of kernels (starting with 8 kernels). This is somewhat similar to VGG-16 [50], although with half the trainable layers. After a $32 \times$ downsample has been achieved, we attach a decision layer which consists of a dense layer followed by three separate dense layers with separate predictions layers after (one for each point being predicted). While this is not a common approach in keypoint prediction, we found that it gave better performance than having the points predicted as a single vector. We theorize that this may be because shared units between each of the three predictions leads to stronger correlations between them, reducing overall prediction flexibility.

3.1.1.2 Training Details

Generating the training data for this is straightforward given a set of annotations with the associated points to learn. The dataset that we created for this purpose contains approximately 2100 training images, 700 validation images, and

900 test images.

First, each image and its corresponding points are resized to a fixed width while maintaining the aspect ratio. This is done to somewhat normalize the scale of objects in the image, so that large differences in original image size don't produce large differences in training loss magnitude. Each image is then further rescaled to the network size (128×128), and then the corresponding targets are rescaled to the range $[0 - 1]$. The size of the fixed width image is recorded as well as \vec{s} . The loss function that we use for each fluke keypoint is the Euclidean distance between the target and predicted point. We also include a scalar scaling factor α , which scales each point by a proportion of \vec{s} . Thus, we have the scaled Euclidean loss function SE

$$SE_\alpha(\vec{t}, \vec{p}, \vec{s}) = \|(\alpha * \vec{s}) \odot (\vec{t} - \vec{p})\| \quad (3.1)$$

Where \vec{t} and \vec{p} are the target and predicted coordinates respectively ⁴. We then average this loss function over the keypoints that the network outputs.

The networks are trained for 1000 epochs with α set to $2e-2$ using the Adam [29] optimizer (with recommended settings) and $l2$ regularization on the trainable parameters with a decay of $1e-4$. All of these hyper parameters were tuned using the validation set, although the possible parameter space was not fully explored due to time constraints.

3.1.1.3 Evaluation

On average, the best network achieved a 15 pixel distance error on the validation and testing sets (in the original image scale). While this may seem like a lot, the trailing edge extraction (and subsequently matching accuracy) was not severely affected in by switching from ground truth points to predicted points (see Figure 4.2).

We find that, for the vast majority of images, the network achieves a low pixel distance error, while there are a few that have a much higher error (see Figure 3.3). Qualitative inspection of these images shows that they are either of flukes which are

⁴ \odot denotes elementwise multiplication



Figure 3.2: Example Keypoint Failure. Example image showing a keypoint extraction failure case from its testing set. Note the difference in pose of the fluke from the success case shown in Figure 3.1. This is an example of a fluke image that violates our assumptions.

not the singular or major object in the image, or flukes that are significantly rotated out-of-plane relative to the camera. An example of this is shown in Figure 3.2.

We attempted to use a spatial transformer network [27] to try and handle these cases, but we were unable to get it to perform as well as the standard convolutional network, nor produce sensible transformations. Since these failure cases represent a small amount of the dataset, it is both difficult to train a network to handle them, and simultaneously not a large issue. Additionally, we find that keypoint extraction failures are only a small percentage of the matching failure cases that we encounter.

3.1.2 Basic Trailing Edge Extraction Algorithm

The base algorithm for extracting the trailing edge uses the vertical gradient information of the image (denoted as I_y). We extract I_y using a vertically oriented 5×5 Sobel kernel [51]. I_y is then normalized with min-max scaling, giving N_y as

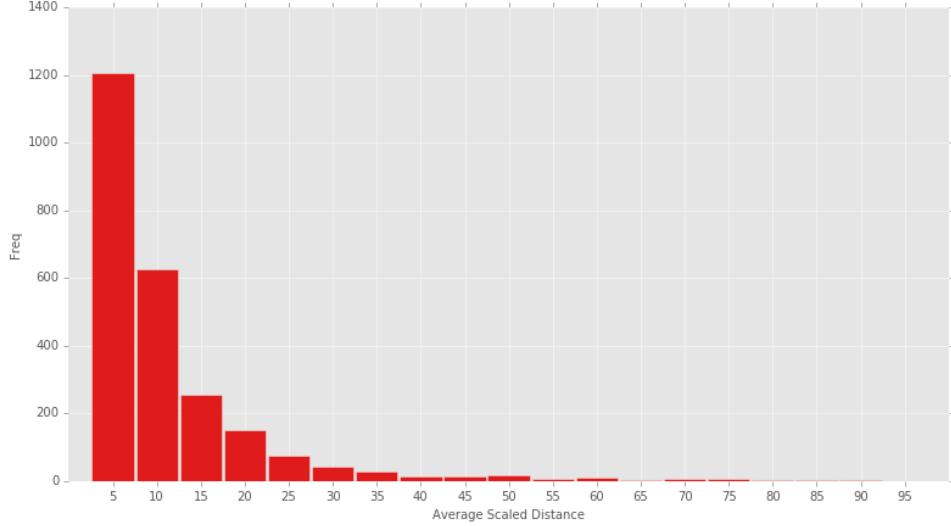


Figure 3.3: Histogram of Keypoint Distances. This is a histogram of the average distance from predicted keypoints to annotated keypoints on the testing set for the keypoint extraction network. The vast majority of keypoints are predicted within 10 pixels of the true keypoints.

$$N_y = \frac{I_y - \min(I_y)}{\max(I_y) - \min(I_y)} \quad (3.2)$$

Given N_y , we extract a trailing edge starting at the left tip of the fluke (a point denoted s) and ending at the right tip (a point denoted e). To do this, we scan each pixel (i, j) in N_y starting from the column $s_x + 1$ and ending at the column e_x , updating a cost matrix C cost with the following recursive update rule:

$$C(x, y) = \begin{cases} 0 & x < 0 \\ \infty & y < 0 \text{ or } y > h \\ \min_{y - \frac{n}{2} \leq y_c \leq y + \frac{n}{2}} C(x - 1, y_c) + N_y(x, y) & \text{else} \end{cases} \quad (3.3)$$

Where n is a neighborhood constraint and h is the height of N_y . We default n to 3, meaning that each pixel considers its immediate “neighbors” from the previous

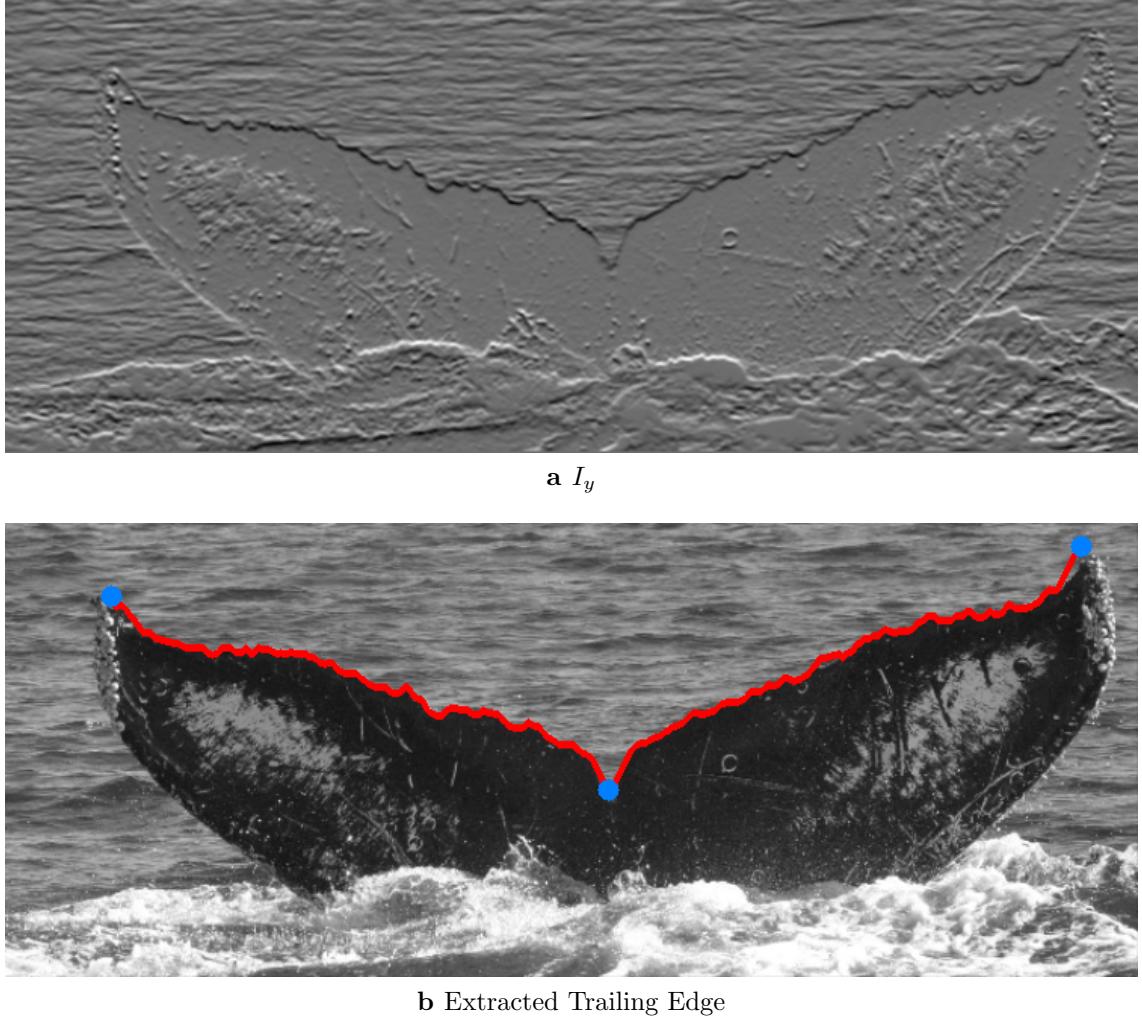


Figure 3.4: Example Trailing Edge Extraction. Example of the baseline trailing edge extraction with $n = 2$. Note that the gradient image has a significant black area where the trailing edge is, making this an easy case.

column.

As C is filled out, we also keep a backtrace matrix B , which keeps track of the index of the minimal candidate neighbor chosen in equation (3.3).

$$B(x, y) = \operatorname{argmin}_{-n \leq n_c \leq n} C(x - 1, y + n_c) \quad (3.4)$$

Once the end column is reached, we scan the columns in reverse order from e to construct the path by adding the chosen neighbor from B at each step. More formally, we start the trailing edge sequence TE^0 at (e_x, e_y) , and add elements to

TE as follows

$$TE^i = (TE_x^{i-1} - 1, (TE^{i-1})_y + B(TE_x^{i-1}, TE_y^{i-1})) \quad (3.5)$$

In order to enforce that the path begins at s , we apply the following process before running the above.

$$N_y(s_x, \cdot) = \infty \quad (3.6)$$

$$N_y(s_x, s_y) = 0 \quad (3.7)$$

This “forces” the path to start at s , as otherwise it would take on infinite cost. We initially did the same “forcing” for the point denoting the bottom of the notch, although we find that this can affect matching accuracy negatively if it is too far from the trailing edge. An example of the base algorithm’s output is given in Figure 3.4. While this algorithm has no understanding of humpback whale flukes, it generally finds high quality trailing edges in images that are constrained around the fluke with oceanic backgrounds (which is a large majority of the dataset at hand). In the next section, we outline our approach for trying to make this more robust.

3.1.3 Trailing Edge Scoring

As mentioned in the beginning of this section, using only the gradient information for extracting the trailing edge works in many cases, but is not a robust method.

If we had a score of each pixel’s “trailing edginess” in an image, the trailing edge extractor could make use of this information to make better choices in trailing edge extraction. An example of this at work is shown in Figure 4.8. To do this, we need a prediction of whether or not each pixel belongs to the trailing edge of a fluke — a task that is best suited to a fully convolutional network.

In the fully convolutional networks that we use, we learn “same” convolutional kernels at each layer so as to maintain spatial shape from one layer to the next, except at explicit downsampling layers. In order to construct a “same” kernel, we

use square kernels of size $k \times k$ such that k is odd, and apply them with a stride of 1 in each direction. Then, in order to ensure that the size of the output is the same as the size of the input, we add a zero-padding of size $\lfloor \frac{k}{2} \rfloor$ pixels to each side (filling in the corners). In the networks used for trailing edge scoring, all convolutions (aside from the downsampling, i.e. max pooling layers) are “same” convolutions. The four major variants on trailing edge scoring networks that we evaluated are detailed below. All of these networks function on the same paradigm of taking an arbitrarily sized image and producing an image of the same size but with a class score for each pixel.

The dataset used to train these networks is sampled from trailing edges extracted using the basic method detailed above, however annotated with manual adjustments to fix many of the more common issues we encountered with this algorithm. These manual adjustments often included manually adding “control points” along the edge, which often fixed major failures. Additionally, small adjustments were made that would not be found by the gradient finding algorithm. Due to the extensive manual effort required to annotate and correct these images, only about 500 images were annotated, many of which required little to no manual adjustments.

Given these manually verified trailing edges, we generated the training set for trailing edge scoring by extracting a 128×128 patch at 128 pixel intervals along the trailing edge (giving a positive patch), and a corresponding patch (with no trailing edge pixels) randomly sampled from the left over space above and below the trailing edge (giving a negative patch). As the images that this dataset were extracted from were generally 960 pixels in width (with the trailing edges being only slightly smaller), on average about 15 patches were extracted from each image. These patches are then randomly split into training, testing, and validation sets each with 3700, 1200, and 1600 patches respectively.

3.1.3.1 Trailing Edge Scoring Architectures

Several network architectures were tried, although here we only report the major variants. One major consideration that has to be made when selecting a fully

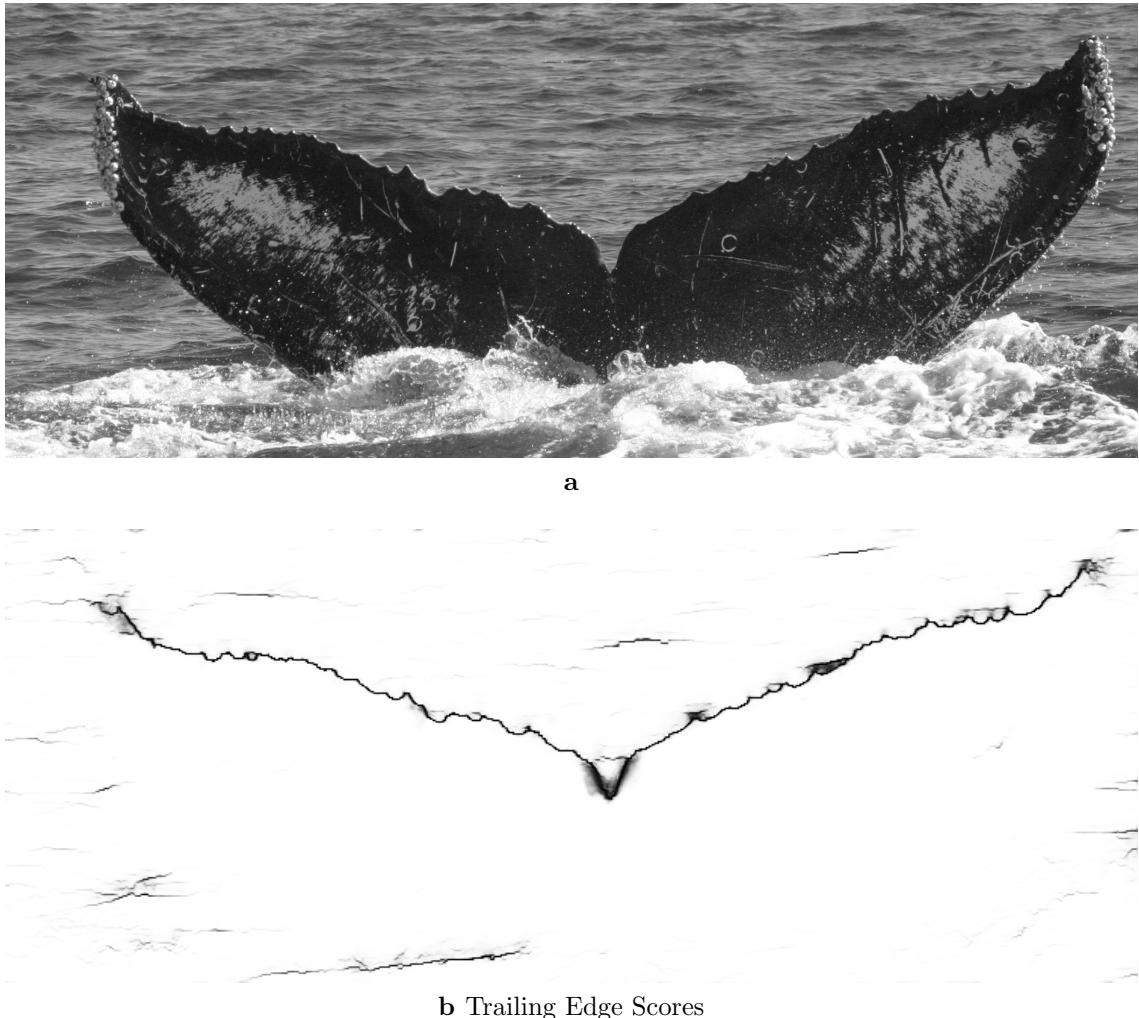


Figure 3.5: Example Trailing Edge Score. Bottom image is the Residual scorer’s classification of the top image. Trailing edge is class is colored black.

convolutional network architecture is its receptive field⁵. If the receptive field is too small, it may not have enough information to accurately determine if a given pixel is part of the trailing edge. However, in order to increase the receptive field without massively increasing the depth of the network, we must downsample, which makes the output less fine-grained.

Simple This network is simply a stack of 6 “same” convolutional layers (of decreasing spatial extent), with no downsampling regions. This has a small receptive

⁵The receptive field of a convolutional network is, at a high level, the region of input that affects the output at a given layer

field, but can produce detailed predictions. Due to the small receptive field however, at convergence it gives low precision predictions, and seems to (in many cases) produces many of the same mistakes that normal gradient based trailing edges do.

Upsample To deal with the small receptive fields, convolutional networks typically downsample at various stages, which doubles the receptive field. This down-sampling usually takes the form of max pooling (instead of e.g. convolutional kernels with a stride greater than one). The Upsample network downsamples the input $8\times$ through alternating “same” convolution layers and max pooling layers, makes a prediction, and then simply upsamples its output (with bilinear interpolation). The Upsample architecture is analogous to the FCN-32s architecture in [35]. The goal of this is to increase the receptive field of the network’s output layer, although the predictions it makes are very “blocky” (see Figure 3.6c).

Jet Following the deep-jet architecture from [35], we modify the Upsample network to combine prediction layers at intermediate levels of downsampling. The essential idea behind this is to take the prediction output from the last downsampling layer, and combine it with predictions made on the output of the layers before the downsample. This is repeated in a cascade, ending with a combination of the merged predictions and a prediction made on the last convolutional layer before the first downsampling layer. Ideally this allows the network to take better advantage of both a deep representation of its input and a lower level spatial layout that allows it to make finer predictions. This architecture gives far more fine-grained predictions than the Upsample network while taking advantage of the increased receptive field size, without increasing the number of parameters greatly.

Residual While the receptive field of the Simple network is small, it can produce very fine trailing edges, which we found to be necessary for good trailing edge extraction. We can still increase the receptive field by stacking small convolutions, but the rate of increase is lower, meaning that more layers are needed. We could simply create a very deep network of “same” convolutions, however training very deep networks like this can run into problems very quickly, as found in the work of

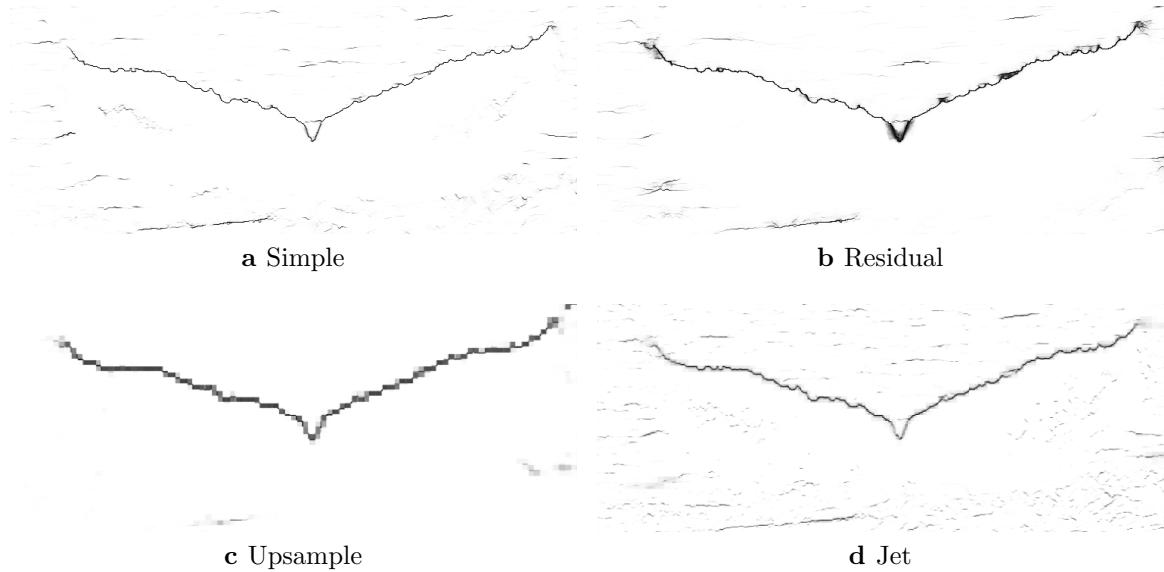


Figure 3.6: Trailing Edge Scores. These are the trailing edge scores given by each of the networks described in section 3.1.3.1 on the image used in Figure 3.5. A darker pixel is predicted to be part of the trailing edge by the network.

He et al. [21], which introduce residual connections to address this problem. These connections are simply shortcut connections from one layer to another, which allows the network to learn the “residual” of its input rather than the whole transformation. We create a residual network architecture by stacking 64 3×3 “same” convolution kernels, and adding a residual connection every other layer. This is our best performing network. We also find that it performs far better than an equivalent network without shortcut connections.

An example of all of the above networks’ outputs can be seen in Figure 3.6.

3.1.3.2 Using the trailing edge scores

Once we have a “trailing edginess” score for each pixel in an image, we need to combine this information with the normalized gradient N_y in a way that causes the trailing edge extraction algorithm to follow those pixels that the scoring map marks as trailing edge. More formally, the trailing edge scoring network gives us an image $T_p \in [0 - 1]^{w \times h}$ (where w and h are the width and height of the image) which denotes the network’s predicted probability of each pixel being part of the trailing edge. The most simple and obvious way to combine this information with N_y from

before is to combine them with a mixing parameter β .

$$S_{te} = (1 - \beta) * N_y + \beta * (1 - T_p) \quad (3.8)$$

We use $1 - T_p$ in this case because we are minimizing the path through S_{te} . Once done, the trailing edge extraction algorithm proceeds as described in the previous section.

Another variant on combining T_p and N_y that we tried was to dilate the trailing edge predictions and then forbid the trailing edge from going outside of those predictions. An inherent difficulty with this approach is that in order for it to work, we need to have a guarantee from the trailing edge scorer that it will not produce major gaps in its predictions. If there are such gaps, then the trailing edge will not be extractable with this method. It is difficult if not impossible to guarantee that this will not happen, although with more data the risk can be mitigated. In our experience, these breaks were common enough that this approach was abandoned.

3.1.3.3 Training Details

All networks were trained for 100 epochs (or until convergence) with a batch size of 32 and l_2 regularization using a decay of $1e-4$. We used the Adam optimizer [29] (with the recommended settings) for calculating weight updates.

One detail that turned out to be important is the class imbalance. The trailing edge pixels (necessarily) make up a small percentage of the total image, meaning that these networks could get a fairly high accuracy (and thus low loss) simply by predicting only background pixels. In order to prevent this, we only sample a negative patch once for every positive patch (as detailed above), and additionally we weight the loss for the trailing edge pixels $10\times$ higher than the loss for the background pixels. This provides a much better trailing edge extraction even if it can reduce precision.

As noted earlier, many of the manually annotated trailing edges did not require extra input. This can be an issue as it biases the network towards marking pixels as trailing edge when they would be marked as such based on just the image gradient. In order to mitigate this issue, we included some data augmentation, by

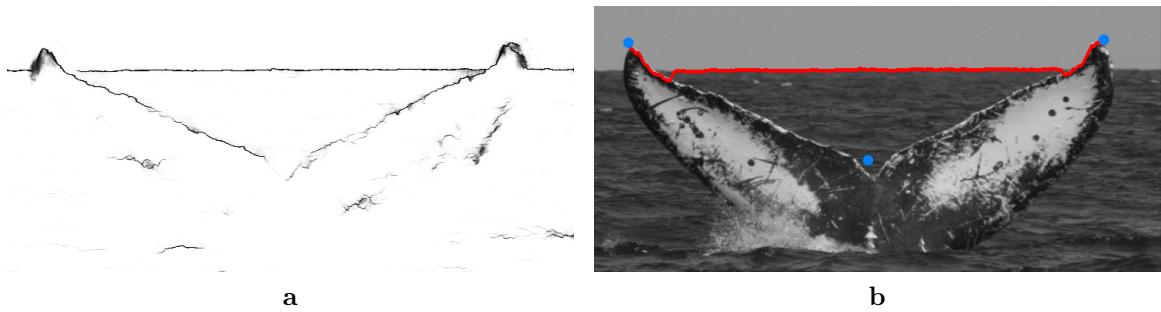


Figure 3.7: Trailing Edge Scoring Failure. Unfortunately even the Residual architecture is still imperfect, and can lead to catastrophic trailing edge failures like this one. However this is a rare case.

random application of Gaussian blur and randomly inverting the pixel intensities with probability of $P = 0.3$. The inverting of pixel intensities as data augmentation is meant to try and influence the network to handle cases where the trailing edge is white against a somewhat darker background, producing an inverted gradient compared to when the trailing edge is dark. It is difficult to measure how successful this was without further annotated data, however in terms of overall matching accuracy we observed a small improvement using a network trained with this augmentation scheme. All of these hyper parameters were tuned using the validation set, although the possible parameter space was not fully explored due to time constraints.

3.1.3.4 Evaluation

Due to the overwhelming class imbalance, the model accuracy is rather meaningless (e.g. the accuracy of an all-background prediction is 99%). Instead, we report the intersection-over-union (IoU) score between the two pixel label classes (i.e. background and trailing edge), which gives a much better idea of model performance. We also report the precision and recall of the model.

The poor performance of the Upsample architecture is also reflected in the match accuracy of trailing edges extracted from its scores (see Figure 4.6). With the exception of the Upsample architecture, the performance differences of the other architectures are insignificant, and we find that they all perform comparably.

Unfortunately it does not appear that these models learned to reliably find only the trailing edge, as can be seen in 3.7. However we observe that unless a false

Architecture	Training			Validation			Testing		
	Pr.	Re.	IoU	Pr.	Re.	IoU	Pr.	Re.	IoU
Simple	0.59	0.95	0.57	0.59	0.94	0.57	0.60	0.95	0.59
Upsample	0.17	0.88	0.17	0.17	0.86	0.17	0.17	0.87	0.17
Jet	0.62	0.89	0.57	0.62	0.88	0.57	0.63	0.89	0.58
Residual	0.57	0.93	0.54	0.57	0.92	0.54	0.58	0.93	0.56

Table 3.1: Table showing the precision, recall, and IoU of each of the evaluated trailing edge scorers on each section of the trailing edge dataset. For the purposes of this analysis, we use the `argmax` over the classes to determine a positive (i.e. trailing edge) or negative pixel.

trailing edge prediction allows for a continuous path from the left to right tip, it is rare for this failed prediction to cause the trailing edge extraction algorithm to go severely “off-course”.

3.2 Trailing Edge Matching

Given extracted trailing edges, we now define a method for doing a one-to-one comparison between a given query and database trailing edge. The simplest way to do this is to define a (potentially non-metric) distance function between any two trailing edges. Once this distance function is defined, we can identify an individual by its trailing edge (referred to as the query trailing edge) by looking at the identity of the closest trailing edge in the database. As a distance function we use dynamic time warping over block curvature measurements, using a weighted Euclidean distance as a local distance function between curvatures. This is essentially a ranking function over a database of identities, where each identity is assigned the rank of the closest image.

3.2.1 Curvature Measurement

In order to do this, we first extract the curvature from the trailing edge. Given the trailing edge as a sequence of coordinates into the original image, we construct a zero-image I_0 of shape similar to the original image. Each pixel corresponding to and below the trailing edge in I_0 is set to 1. This essentially means that everything “inside” the fluke is set to 1, and everything outside the fluke is set to 0 — with the

assumption of course that everything below the trailing edge is part of the fluke. Because the trailing edge extractor produces only one coordinate per column, we can do this safely, however this algorithm could be easily adapted to this not being the case. Once this is done, we calculate a summed area table [13] ST from I_0 as follows.

$$ST(x, y) = \sum_{i=0}^{i=y} \sum_{j=0}^{j=x} I_0 \quad (3.9)$$

Conceptually, the next step is to slide a square of size $m \times m$ pixels centered on each point, and measure the percentage of that square that is within the filled in trailing edge. The value m is the scale at which we measure curvature, and is computed as a percentage of the trailing edge length. We compute multiple scales M of curvature for each trailing edge. This curvature measurement is carried out by computing $BC_m(x, y)$ for each (x, y) coordinate in the trailing edge.

$$\begin{aligned} b(i) &= i - \frac{s}{2} \\ e(i) &= i + \frac{s}{2} \\ BC_m(x, y) &= \frac{(ST(b(x), b(y)) + ST(e(x), e(y))) - (ST(b(x), e(y)) + ST(e(x), b(y)))}{m^2} \end{aligned} \quad (3.10)$$

The numerator in equation (3.10) gives the total area within the square that is below the trailing edge, which we then normalize by dividing by the square's area. The set of scales to choose presents a large parameter space, however we have found that the scales $M = [2\%, 4\%, 6\%, 8\%]$ (as percentages of the trailing edge length) work well for our purposes. We then treat this curvature measurement as a $|M| \times l$ matrix BC , where l is the length of the trailing edge. Two example curvatures (with their corresponding trailing edges) are given in Figure 3.8. In this case, we can see that at a high level, these curvature patterns look very similar.

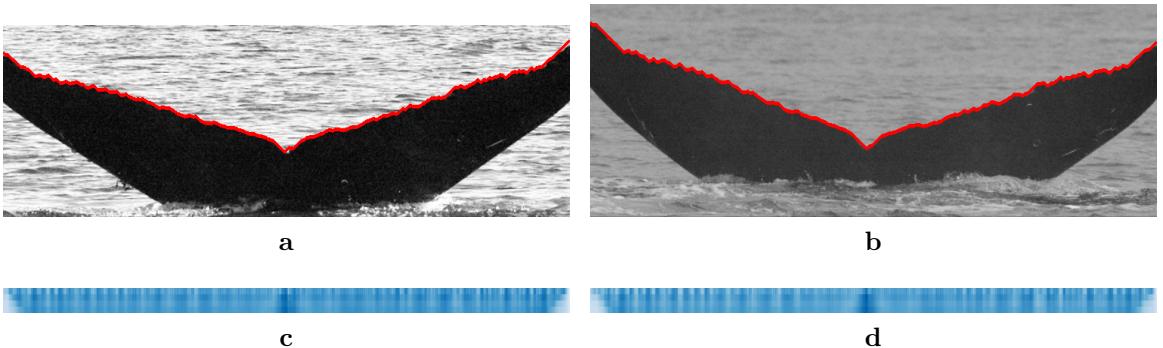


Figure 3.8: Trailing Edge Curvature. The top images are of the same individual, and the bottom images visualize the corresponding curvatures for the trailing edges that were extracted. Each row in the visualization is a curvature scale, increasing from top to bottom. Note that darker blue implies a “valley” in the trailing edge, whereas lighter blue implies a “peak”.

3.2.2 Sequence Matching

Given two sequences of curvatures BC_1 and BC_2 (referred to as query and database curvature respectively) we match them using dynamic time warping as follows. First, we create a cost matrix C of size $l_1 \times l_2$ (i.e. the length of the first and second trailing edge respectively). We initialize this cost matrix by setting the first column and row to ∞ , and then $C(0, 0) = 0$, intuitively forcing the optimal path to match align the beginning of BC_1 with BC_2 . Then, for each cell (i, j) in the cost matrix starting with $C(1, 1)$, we use the following update rule

$$D(c_1, c_2) = \|\vec{c}_1 - \vec{c}_2\|_2 \quad (3.11)$$

$$C(i, j) = D(BC_1(i, \cdot), BC_2(j, \cdot)) + \min(C(i - 1, j), C(i, j - 1), C(i - 1, j - 1)) \quad (3.12)$$

Where \vec{c} is a vector of the curvatures at different scales for a point. Initially there was also a weighting term that would weight each curvature scale differently, however we found that weighting each curvature equally was the best option. We then take the value of $C(l_1 - 1, l_2 - 1)$ as the distance between the two curvatures⁶.

Additionally, we impose the Sakoe-Chiba [45] locality constraint T so that

⁶Similarly to making $C(0, 0) = 0$, this enforces that the end of each trailing edge aligns

for each element i in BC_1 , we only consider the range over elements j of BC_2 $j \in [\min(i - T, 0), \max(i + T, l_2)]$. This essentially provides a bound over the possible matches between points on the trailing edge, preventing (for example) the first element of BC_1 from matching with the last element of BC_2 . If we do not believe that these matches would be reasonable, then this bound not only prevents these (likely) erroneous matches, but also greatly speeds up the computation of the distance.

We set T as a percentage of l_1 . For most of these experiments T is set to 10%, which appears to minimize the time taken for each comparison while preserving the overall accuracy of the algorithm (see Figure 4.12).

It's worth noting that while this distance measure is not a metric distance (i.e. it doesn't satisfy the triangle inequality), it is a symmetric distance as the local distance function (3.11) is symmetric [38].

3.3 Alternative Approaches

In this section, we list and briefly describe alternative approaches that were tried, although they did not prove accurate enough to make it into the final system.

3.3.1 Aligning Trailing Edges

One obvious pre-processing step that would make sense when comparing trailing edges is to make sure that they are aligned in image space. However, we found that doing so when comparing curvature was often unnecessary (due to the invariances to rotation and translation), and using the Euclidean distance between points on the aligned trailing edges (i.e. in place of curvature for (3.11)) did not give good results. There were two approaches to alignment that we evaluated, although neither achieved top-1 accuracies above 20%.

3.3.1.1 Keypoint Alignment

As noted in the section on fluke keypoints, there are three points to predict — left, notch and right. Originally the intention for recording (and predicting) all three, as opposed to just left and right, was to have three corresponding points with which

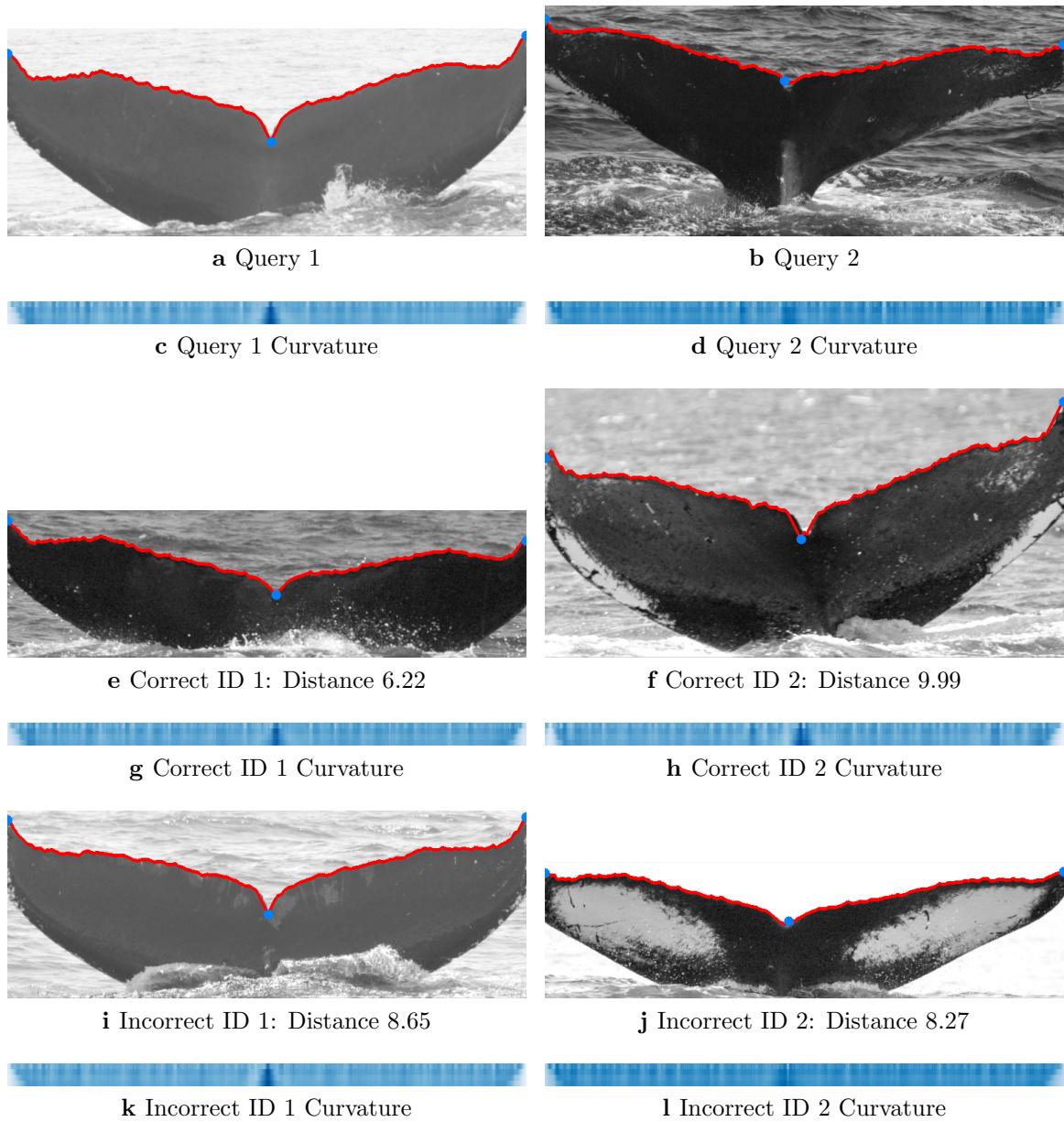


Figure 3.9: Example Matches. The left side shows a success case, and the right side shows a failure case.

to estimate an affine transformation from database image onto query image. This would be done prior to any computation of trailing edge or curvature. One major issue with aligning these images however is that if a non-affine transformation was required, the trailing edge itself would be warped in such a way that made matching difficult.

3.3.1.2 Dynamic Time Warping Alignment

We also attempted to align the trailing edges based on matches generated by the dynamic time warping, in similar fashion to the AI-DTW approach laid out in [44]. This approach used an iterative alignment process using the correspondences found by DTW (using either curvature distance or Euclidean distance as criteria). Essentially this method would find alignments using DTW, and then use these alignments to estimate an affine transformation of the database image onto the query image — and then repeat until convergence. However, we found that this process oftentimes wouldn’t converge, and when it did the alignments provided were of worse quality than those found by aligning the three fluke keypoints. Additionally, the extra time taken to carry out this process proved untenable.

3.3.2 Histogram Matching

One early curvature comparison method that we evaluated was to use histograms to match instead of a sequence-based method. This is similar to the approach taken in Leafsnap [32]. We found that for our task even high resolution histograms did not provide enough detail to match the trailing edges properly, although this could potentially be explored further.

3.3.3 Embedding via Convolutional Networks

We also made an attempt at training convolutional networks to directly embed the images of flukes into a k -dimensional vector, much like the work done for face recognition in Schroff et al. [48] and Parkhi et al. [42]. However, most of the previous literature on this technique is applied to larger datasets such as Labeled Faces in the Wild [23], which is significantly larger than the dataset that we had available. A major factor in this is that these larger datasets often have five to ten images per identity (if not more), whereas most of the identities in our dataset had one or two images associated.

Regardless, we attempted the embedding approach (from raw images), however even a severely overfit convolutional network could only achieve less than half the top-1 accuracy on its training set that the main method is able to achieve. We tried

both triplet loss (a modified version of the one detailed in [48]) and contrastive loss [18] to no avail. We believe that the small amount of images per identity is the main factor for the failure of these methods, and that a larger dataset would be necessary to properly train them.

CHAPTER 4

Results

In this chapter we present the results that our primary method achieves on the Flukebook dataset. The main results for the optimal method are given briefly, and then we discuss how different variations on the method affect accuracy.

4.1 Main method

The main method we settled on achieves an 80% top-1 accuracy — meaning that for 80% of the query images, the correct identity is ranked first — on the dataset that we evaluated. The figures that we show in this section give the accuracy up to top-5 cumulatively. In general we find that relative accuracies between configurations do not change significantly as we increase the rank at which we allow a match.

The optimal configuration that we used for this method is given below.

- Every image is cropped between the predicted left and right tips, and resized to the same width (750 pixels) while maintaining aspect ratio.
- We do not use the bottom of the notch as a control point for trailing edge extraction.
- We set the number of neighbors used in trailing edge extraction n to 3.
- We use the Residual architecture for scoring the trailing edge.
 - β is set to 0.5.
- We use [2%, 4%, 6%, 8%] for our curvature scales.

4.2 Configuration Options

4.2.1 Variability in Matching Score

A major issue with our method is that the correlation of the distance between two trailing edges and whether or not it is a true match is not very strong. We can

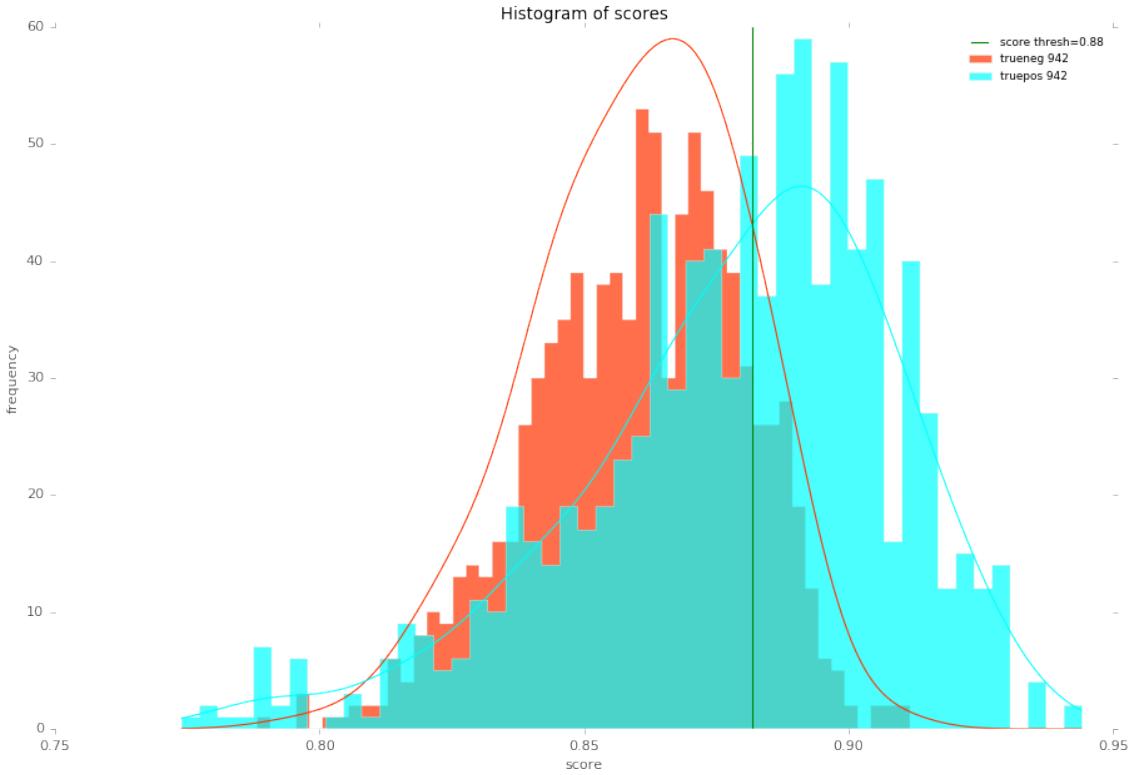


Figure 4.1: Score Separability Histogram. The blue bars in this figure represent true matches, and the red bars represent false matches. The line at score = 0.88 represents the optimal threshold at which to accept a match, although we can see it is not perfect.

see the effect of this in Figure 4.1, which shows that the scores⁷ between query trailing edges and the corresponding ground-truth trailing edge overlaps heavily with the distance between queries and unrelated trailing edges. The result of this is that small changes in the query and database trailing edges can (in some cases) cause the ranking of a ground-truth trailing edge to change significantly. While this effect is somewhat rare, we find that it is the cause of many discrepancies in matching accuracy between different configurations. With this in mind, we discuss the difference in matching accuracy only when it is a result of significant discrepancies — meaning that the score between query and ground-truth trailing edge changes by more than $\epsilon = 0.02$.

⁷The score is computed as $e^{\frac{-\text{distance}}{50}}$ for implementation reasons

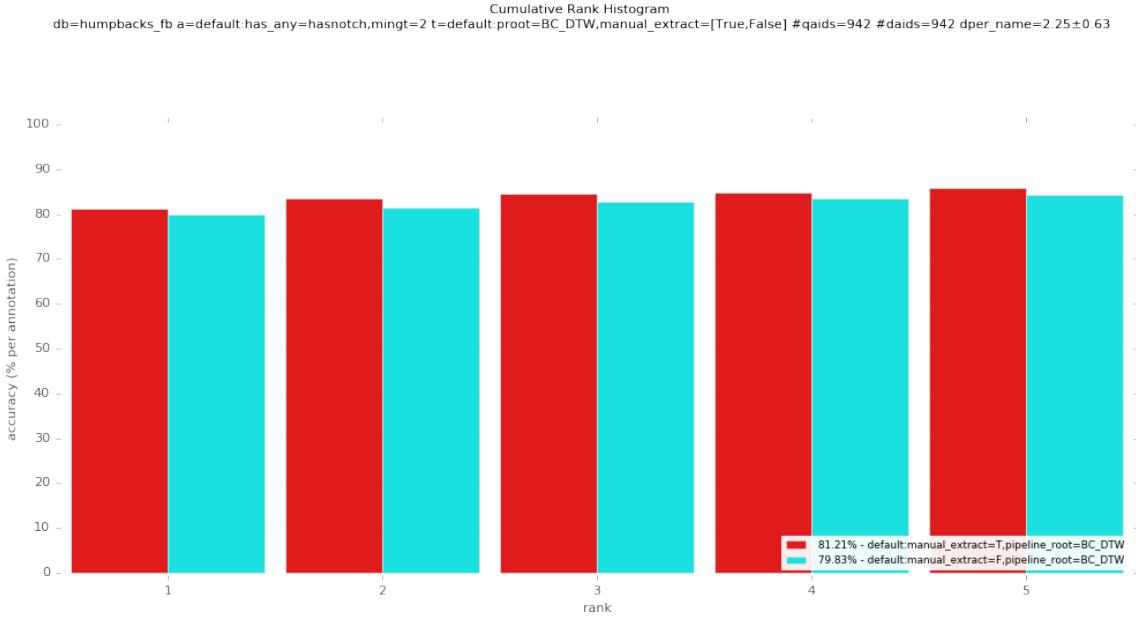


Figure 4.2: Varying Manual Extraction. There is a small difference in matching accuracy between using the manually annotated points (red) provided for this dataset versus the keypoint extractor’s predicted points (cyan). The bottom of the notch keypoint is not used in either of these evaluations.

4.2.2 Effectiveness of Keypoint Extractor

In order to test the effectiveness of the trained fluke keypoint extractor, we give a comparison of matching accuracy to using the manual annotations of these fluke keypoints provided for the Flukebook dataset in Figure 4.2. The manually annotated keypoints provide a slightly better accuracy than the automatically extracted keypoints, and we find that many of these matching differences do not pass the significance (i.e. the difference in query to ground truth distance doesn’t change a lot between configurations). While this does not mean that the keypoints predicted are perfect, it does imply that they are “good enough” to extract a matchable trailing edge, despite being on average 5 and 10 pixels off.

4.2.3 Cropping Width

With dynamic time warping, we theoretically can match sequences of arbitrary lengths, however the distances are distorted by large differences in actual trailing edge length. For this reason we want to fix the length of the trailing edges. Since we are only interested in the width of an image (because of the way the trailing edge

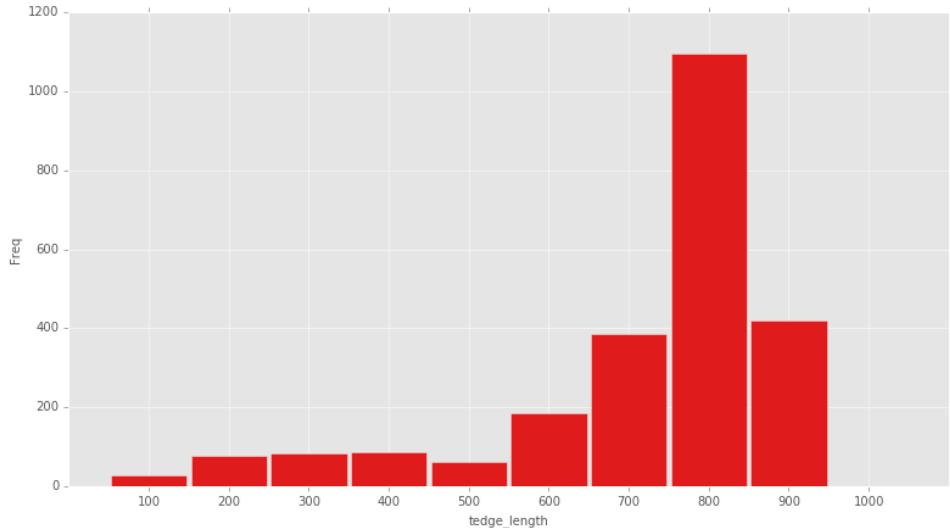


Figure 4.3: Distribution of Unresized Trailing Edge Lengths. This shows a significant distribution of trailing edges centered around a width of 800 pixels.

extraction algorithm works), we can get every trailing edge to have exactly some fixed length w by the following process.

- Crop the image horizontally between the left and right columns found by the keypoint extraction process (or manually determined).
- Resize the cropped image to some fixed width w while preserving the aspect ratio.

In this way, we standardize the trailing edge length so that differences in image size do not affect detection accuracy. One major caveat with this process is of course that using the keypoint extractor's predictions can cause catastrophic failures in this process (e.g. the left and right points are nowhere near a fluke), however in practice we found that this is not an issue.

Ideally, we would choose a w that minimizes interpolation artifacts from making big changes in image size. Figure 4.3 shows the histogram of post-crop widths (i.e trailing edge lengths) from unresized images in the Flukebook dataset, showing

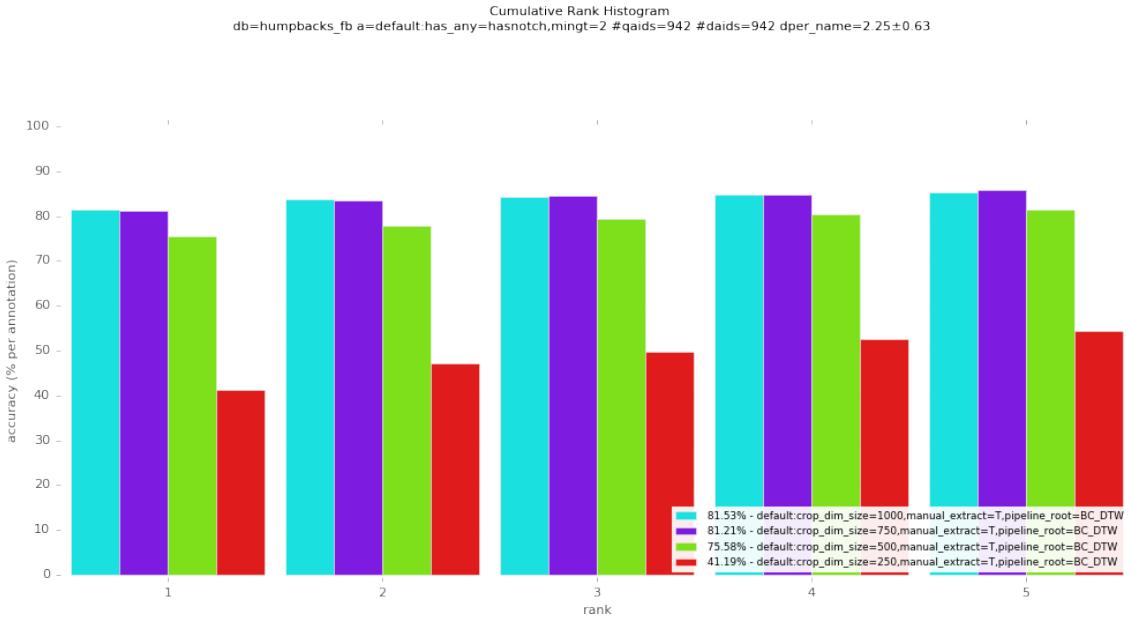


Figure 4.4: Varying w . Note that we use the manually annotated points in this analysis to control for any issues with keypoint extraction.

a large concentration of mass around 800 pixels. Subsequently, Figure 4.4 shows that $w = 750$ performs well, although $w = 1000$ is on par. We select the former for efficiency's sake, as smaller trailing edges vastly improves the speed of the matching algorithm. This provides evidence for the hypothesis that the less the image has to be resized, the better the trailing edge.

4.2.4 Trailing Edge Extraction

One major result that we found was that, when using the averaging method to combine the trailing edge scores with N_y , having a robust trailing edge prediction wasn't as important as having a detailed trailing edge.

4.2.4.1 Trailing Edge Scorer variations

The various trailing edge scorer architectures and their results on the task they were trained for is detailed in the previous chapter. In Figure 4.5 we present the actual matching accuracies that each one produced with the mixing parameter $\beta = 0.5$. We can see that the detailed and higher quality trailing edges produced by the Residual network give a decent performance boost over the other networks,

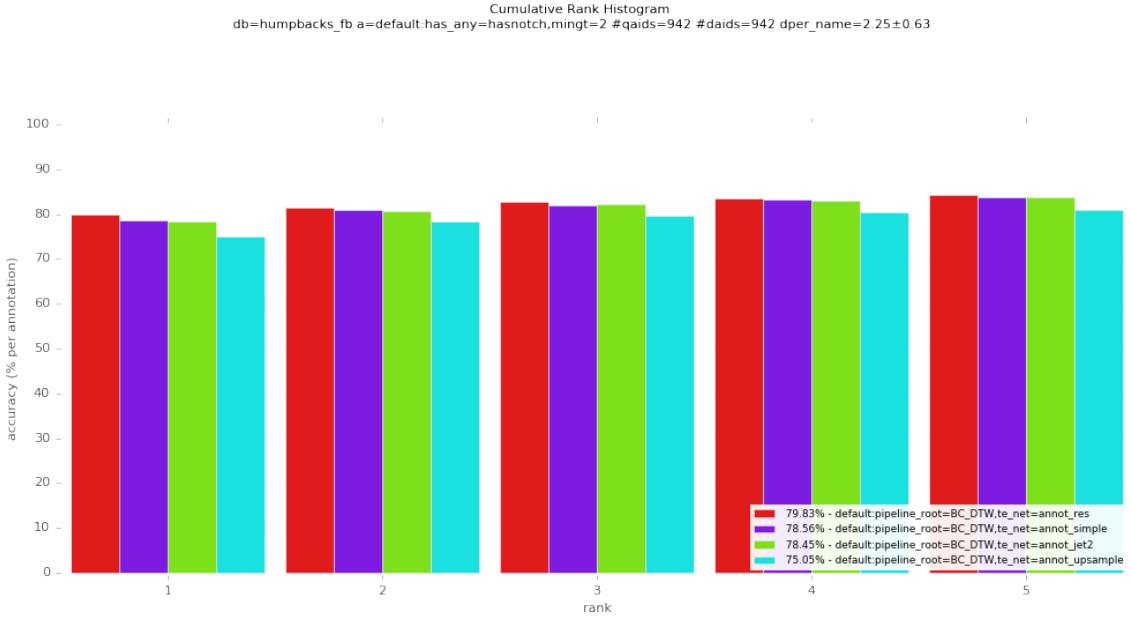


Figure 4.5: Trailing Edge Scorer Architectures. The highest performing trailing edge scorer (Residual) is shown in red, followed by Simple, Jet, and Upsample (in descending order of accuracy).

however this performance boost appears to consist mostly of insignificant changes in query to ground truth distance.

A better evaluation of each trailing edge scoring architecture is given by setting $\beta = 1.0$, which means that only the trailing edge scores are used to extract the trailing edge (see Figure 4.6). Unsurprisingly, since the Upsample network gives blocky trailing edges (see Figure 3.6), it performs very poorly. However in this case, no other network performs poorly, showing that the biggest issue is the low quality trailing edges that are extracted, not so much the lack of breaks in trailing edges.

One caveat with using the Residual network is that, with 64 layers, it consumes a lot of GPU memory when running. This is currently an implementation issue, and the machine that carried out these experiments can handle it. However, given that each network performs comparably, the Simple network is a good choice in general.

4.2.4.2 Using a Trailing Edge Scorer

In Figure 4.7, we can see that simply a pixel-wise average of N_y and $(1 - T_y)$ (i.e. $\beta = 0.5$) produces the best results for the Residual network, though the differences are largely insignificant. However, not using the trailing edge scorer (i.e. $\beta = 0$)

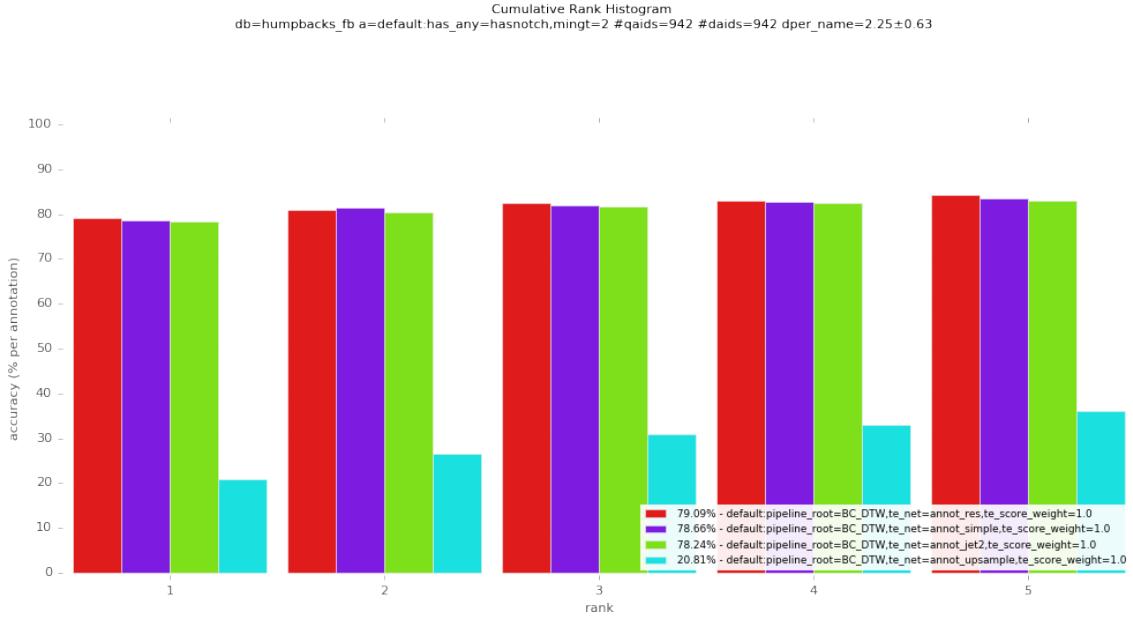


Figure 4.6: Trailing Edge Scorer Architectures at $\beta = 1$. Upsample (cyan) performs significantly worse than the other networks, which all perform comparably.

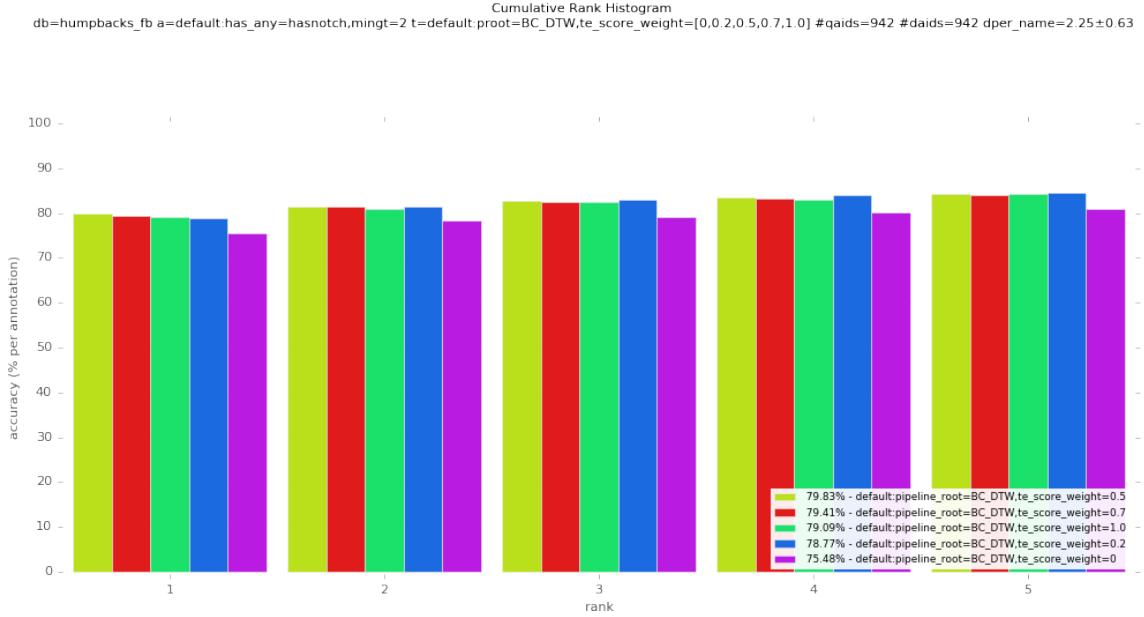


Figure 4.7: Varying β .

performs significantly worse than using one at all (even with a low weight). We show some examples of cases where trailing edge scoring helped in Figure 4.8.

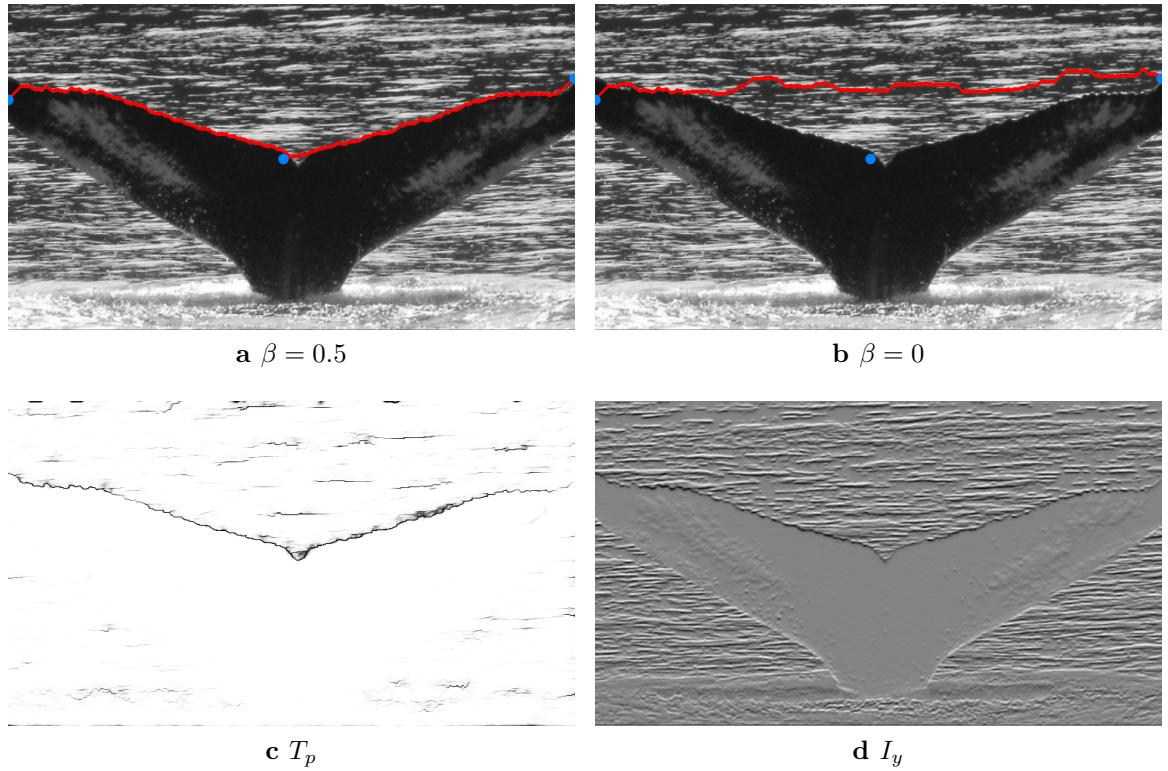


Figure 4.8: Example Use of Trailing Edge Scorer. In (a), we have the trailing edge extracted with the Residual scorer. Compare to (b), which did not use any scorer at all, resulting in a match failure.

4.2.4.3 Number of neighbors in the extraction

The number of neighbors n effectively limits the slope of the trailing edge. We limit it to an odd number for convenience. On the one hand, a lower n can cause the trailing edge to be limited in vertical breadth, but does prevent it from going way off course. Despite this, with trailing edge scoring in place, it might be beneficial to increase n so as to avoid parts of the trailing edge that continually “max out” the number of neighbors.

Ultimately, we can see in Figure 4.9 that limiting the number of neighbors to the immediate neighborhood (i.e. $n = 3$) produces a significant boost over a larger neighborhood. Potentially the trailing edges extracted with $n = 3$ can be somewhat less detailed, and as a result can be more invariant to slight changes in pose of the fluke. While this effect is only present in a small number of cases, the difference in accuracy here is significant, with 78% of the improvements resulting in a score

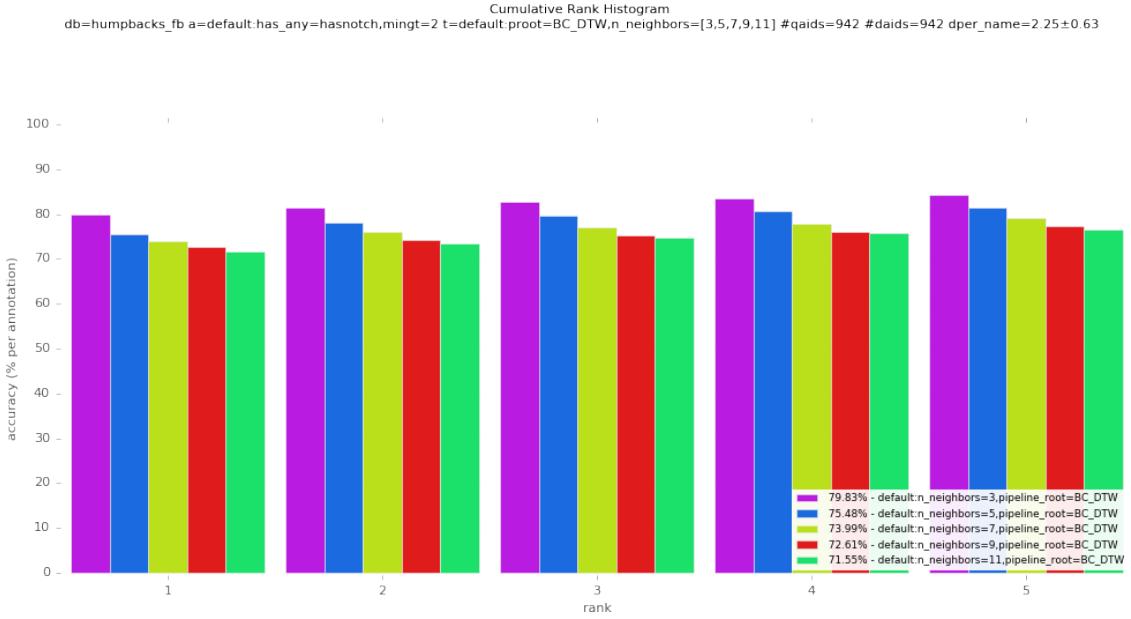


Figure 4.9: Varying n . This shows that the optimal neighborhood constraint is $n = 3$, despite qualitatively producing worse-looking trailing edges. Beyond $n = 5$, the trailing edges can become very noisy.

difference over ϵ .

4.2.4.4 Using the notch

Initially, we believed that using the notch as a control point (i.e. forcing the trailing edge to go through the notch) would give better trailing edges, and therefore an increase in accuracy. As it turns out, this is not the case — using the notch as a control point gives overall worse accuracy, as can be seen in Figure 4.10. Surprisingly, this decrease in matching accuracy holds for both manually annotated and automatically extracted keypoints. However, most of these discrepancies did not result in a score difference above ϵ . This could be because of the small difference in trailing edge between those with a notch and those without, which leads to matching failure but not a large difference in score.

4.2.5 Curvature Scales

Computing the curvature is one of the least parameterized parts of the process, however figuring out what the optimal scales to extract are and how many to compute is non-trivial. Instead of exhaustively exploring these options, we can

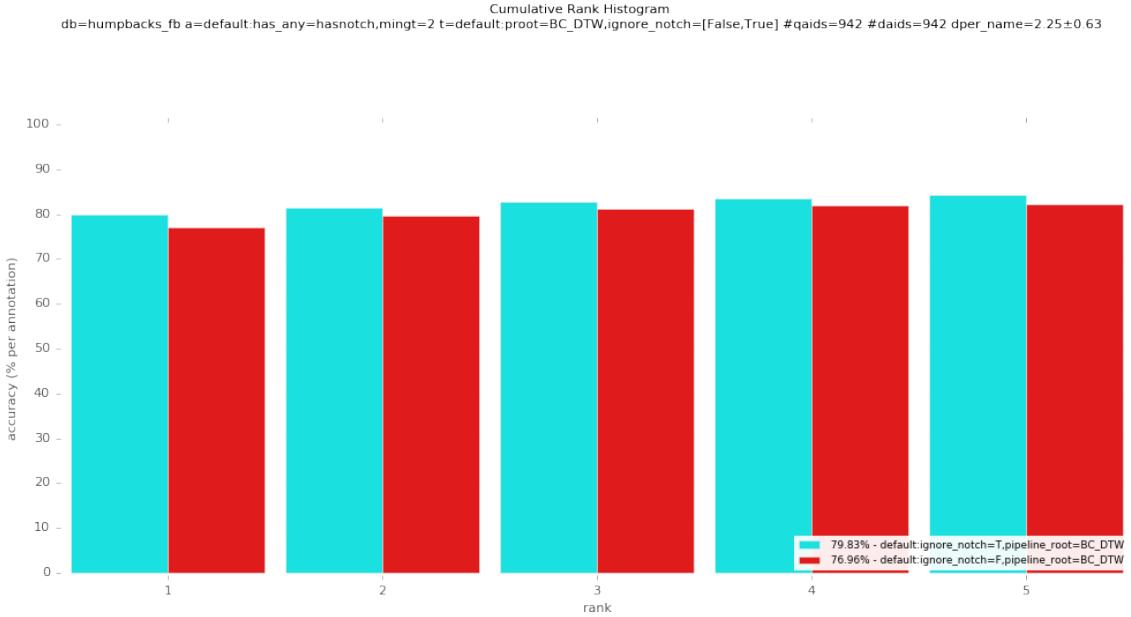


Figure 4.10: Using the notch as a control point.

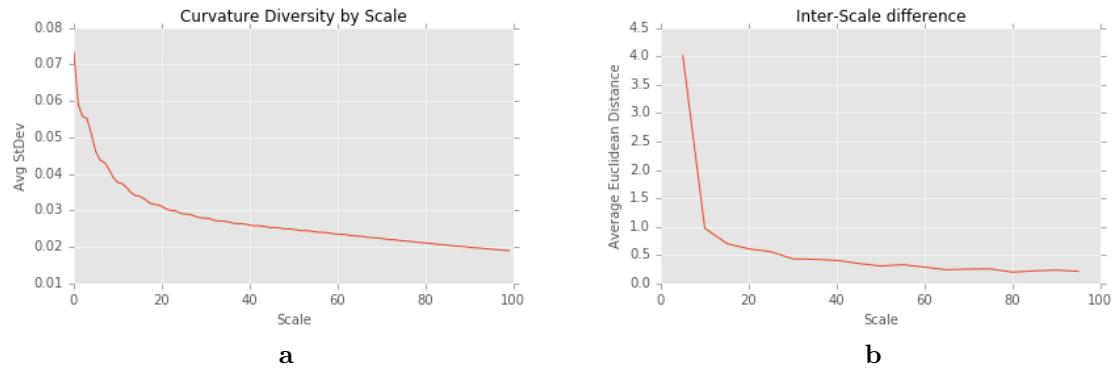


Figure 4.11: Curvature Diversity. Left panel (a) shows the average standard deviation of the (fixed length) curvature at different scales. Right panel (b) shows the average Euclidean distance between successive scales of curvature

use some heuristics to determine which curvature scales to explore. Each scale measures the curvature of some percentage of the trailing edge around each point on the trailing edge. Therefore, in order to determine which scales to measure, we look at how much the actual curvature changes between successive scales, as well as the average variance in curvature at each scale.

We find that (as shown in the Figure 4.11b) as block curvature scale is in-

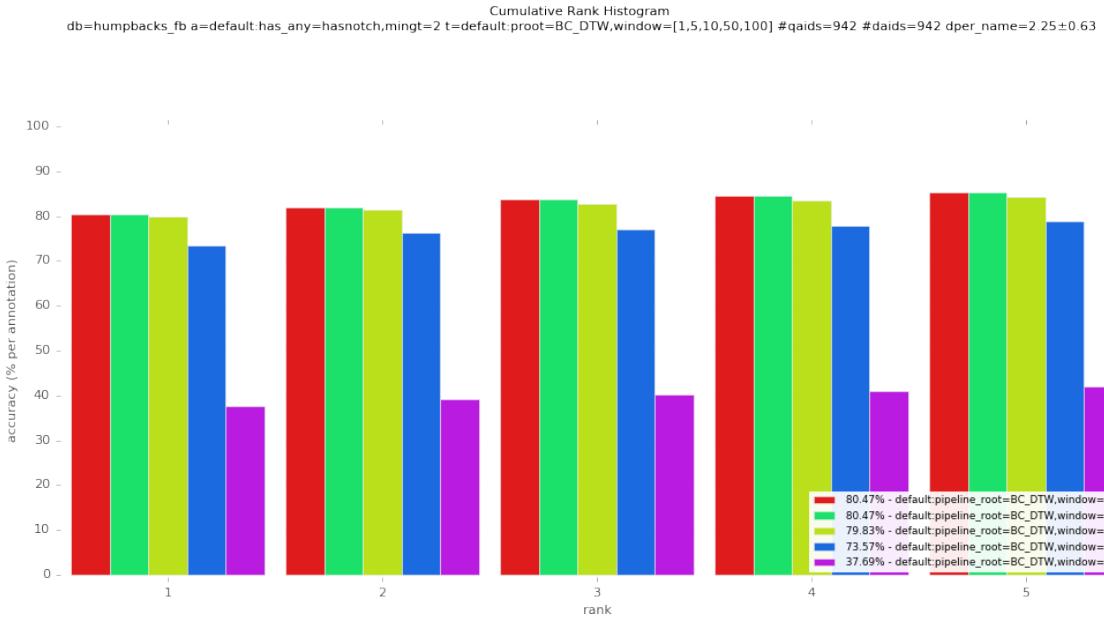


Figure 4.12: Varying Sakoe-Chiba bound.

creased, the diversity at any given point in the curvature goes down drastically (as expected). Therefore, we stick to the lower end of the scale, keeping the curvature scales measured below 10%.

We can also see in Figure 4.11a that successive curvature scales show bigger differences at lower scales than at higher scales, which reinforces using curvatures below 10%, but that it's advantageous to make bigger jumps between scales to maximize diversity while minimizing computation time.

Based on the above, we evaluated scales that run from 1% to 10%, with varying levels of resolution, and found little to no significant effect on matching accuracy compared to using the default set of scales.

4.2.6 Sakoe-Chiba bound

In Figure 4.12, we can see that if we decrease the window (i.e. the Sakoe-Chiba bound) size, at around 10% of the query trailing edge length we maintain only slightly worse accuracy compared to the full window (i.e. 100%), but below this accuracy is severely affected. We find that overall there is a 4 \times slow-down in wall-clock time on our testing machine when going from a window size of 10% to one of 100%. Thus, we use this value for the window size so as to minimize computation

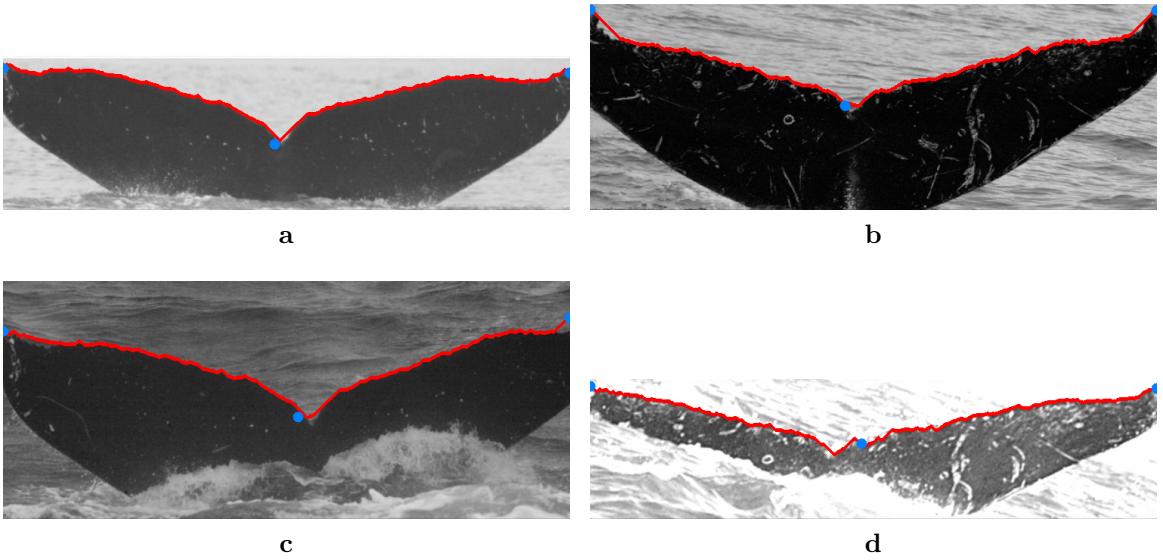


Figure 4.13: Example Disagreements Between Hotspotter and our method. On the left side, (a) was matched correctly to (c) by our method, whereas Hotspotter could not find any matches for (a). On the right hand side however, Hotspotter was able to match two flukes with a large variance in pose and lighting, while our method did not rank (d) in the top 5 matches for (b).

time while maintaining the total accuracy.

Additionally, while it's possible for gross mismatches to occur from there being no window boundary, we can see from Figure 4.12 that this does not pose a problem.

4.3 In Combination with Hotspotter

In a straight comparison, our method gets a slightly better top-1 accuracy than Hotspotter (see Figure 4.14) However, our method targets a specific part of the fluke whereas Hotspotter recognizes general patterns, and thus works on the internal texture of the fluke. Therefore, by combining our method with Hotspotter — if we were able to automatically pick out which algorithm was right for a given ranking — we would expect to see a significant increase in accuracy. We find that in this ideal scenario, we can achieve a 93% top-1 accuracy on the Flukebook dataset. However, automatically deciding which algorithm to use for matching is non-trivial, and we are still exploring it.

See Figure 4.13 for an example where Hotspotter correctly finds a match where

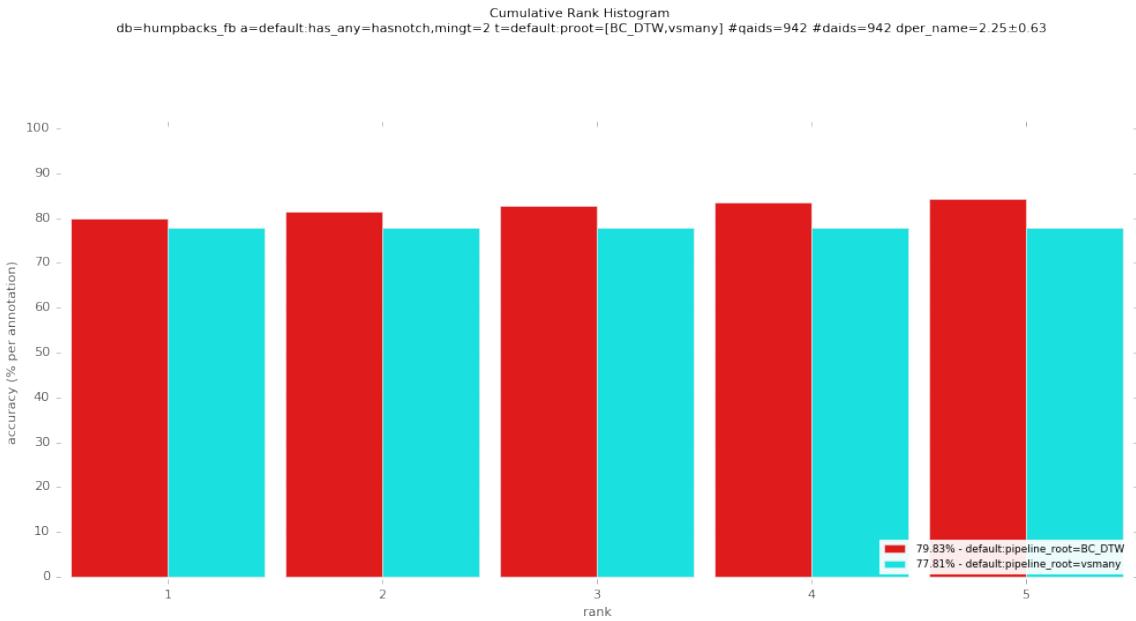


Figure 4.14: Comparison between Hotspotter and our method.

our method fails, and vice-versa.

4.3.1 Characterization of when to use which method

From an intuitive standpoint, it appears that Hotspotter cannot find effective keypoints from trailing edges, which hampers its ability to handle flukes which do not have an apparent pattern.

We hypothesize that this is because the trailing edge — while a distinctive feature — inevitably shares a region with an oceanic background. Since this oceanic background changes from image to image, it cannot verify salient keypoints as matches. As a result, flukes which have little to no internal texture are nearly impossible for Hotspotter to match.

On the other hand, when the trailing edge is unclear or significantly distorted in the image, our method struggles to find an appropriate match. In these cases, Hotspotter can provide a good match, although if the fluke is additionally untextured we have an issue. We were unable to find a single heuristic that correlates with our method failing to find a match. Interestingly, we find that smoothness of the trailing edge, characterized by the standard deviation of the trailing edge slope, does not correlate to match failure, implying that smooth trailing edges are about as easy

to match as rough, distinctive trailing edges. However, since the dataset we use consists primarily of individuals with only two images, there can be a lot of variance in rank that is not necessarily explained by a single property.

CHAPTER 5

Discussion

5.1 Issues with the Proposed Method

The primary issue that we encountered is that the distance between a given pair of trailing edges is very unstable in terms of determining if they match or not (see Figure 4.1)). A side effect of this is that small changes in the trailing edge can lead to matching failures, and in the worst case “push down” the correct match significantly. This “push down” effect severely hampers the effectiveness of our method, as in the use case that we target a human operator verifying a match would ideally only have to look at the top- k flukes for some small value k . However, we find that in many cases the correct match is far down the list, and there is no reasonable value k (i.e. that a human operator would want to find a match in) that all matches are within (see Figure 5.1).

There are several other issues, primarily having to do with the stability and generalization capability of the convolutional networks used for fluke keypoint predictions and trailing edge scores. For the former network, ideally the keypoints could be predicted regardless of our assumption on fluke position and orientation, although it is clear to us that this does not happen. We believe that this is primarily due to the lack of training data that defies these assumptions, although it is also possible that it is beyond the capacity of the models we trained.

Having the fluke horizontally oriented and flat to the camera does not only help the keypoint extraction, but also avoids an obscured trailing edge due to out-of-plane rotations (as seen in Figures 3.2 and 1.3). However imposing this requirement significantly complicates and restricts the actual photography of these flukes.

The network that we trained for scoring trailing edges was ultimately somewhat disappointing. For the most part, these networks were very good at predicting a trailing edge where a strong gradient was present, however this does not significantly improve on the base trailing edge extraction algorithm. Part of the problem for this is similar in nature to the problem for training the fluke keypoint network,

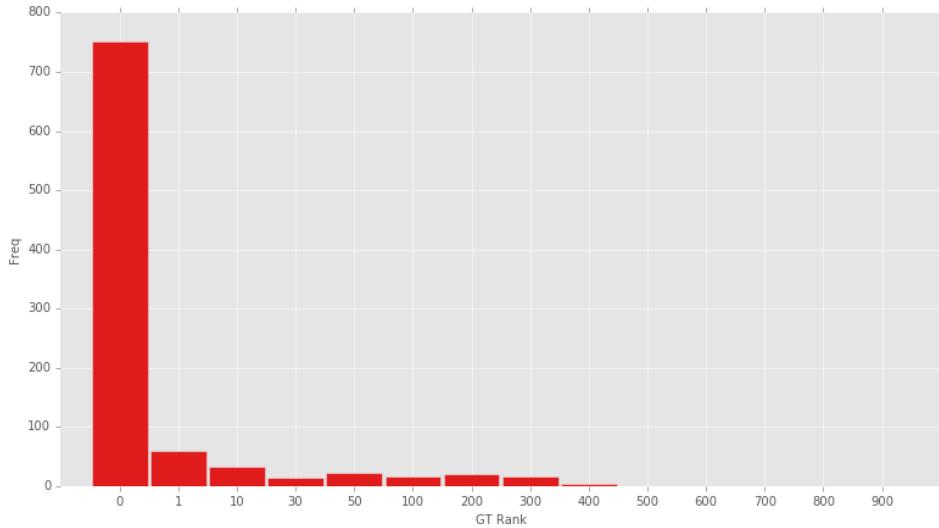


Figure 5.1: Histogram of Ground-Truth Ranks. Note that the histogram ranges are uneven to better show the lower end of the range. In order to have all matches found within the top- k matches we would have to set $k = 414$.

i.e. that a large majority of the dataset represent “easy” cases, making it hard to train the network to handle hard cases. We did use data augmentation to help rectify this bias, but found that it had limited effectiveness. It is possible that more sophisticated data augmentation could help, but primarily we think that spending more time annotating and creating a trailing edge scoring dataset would be ideal.

5.2 Future work

There is a lot of work that can be done based on this method. The immediate focus is to achieve the theoretical 93% accuracy by accurately choosing to use Hotspotter or our method for any given query and result. This would result in a more general method that could be robust to obscured trailing edges and some out-of-plane rotations.

One part of this identification pipeline that we mostly ignored is a detection and orientation step. While orientation is not necessarily important for the curvature

measure (as it is rotation invariant), it is important for the trailing edge extraction. Additionally, being able to detect and crop the fluke automatically from an image would give a much more robust and flexible system. We did not explore this as most of the images in the Flukebook dataset came pre-cropped around the fluke, as well as rotated such that the major axis of the trailing edge was horizontal. These conditions both obviates the need for such a system and make it difficult to train a detection model.

Extracting the trailing edge is currently done with an unsophisticated and restricted algorithm, which can be improved. There are several ways to improve this algorithm, but the main drawbacks are its lack of rotation invariance and the inability to easily draw the second or third best trailing edge. On top of that, the trailing edge scoring networks could be better evaluated with more annotated trailing edges, the creation of which is non-trivial.

There is a lot to be done with the matching algorithm itself, namely in ensuring that it is more tolerant to significant deformations in the trailing edge. One major paradigm that we do not explore in this work is that of extracting and matching multiple features per trailing edge, much in the same way that Hotspotter operates. Hughes et al. [25] take this approach for a more general contour matching system. It would be possible to do something similar but making use of the curvature measures.

5.3 Conclusion

In this thesis we have presented a novel, fully-automated method for photo-identifying humpback flukes that achieves a high top-1 ranking accuracy on a relatively large dataset. This method extracts fine grained trailing edge contours from images of flukes and identifies individuals based on sequence properties of the contour curvature.

REFERENCES

- [1] Y. Akagi, R. Furukawa, R. Sagawa, K. Ogawara, and H. Kawasaki. A facial motion tracking and transfer method based on a key point detection. 2013.
- [2] A. A. Amini, S. Tehrani, and T. E. Weymouth. Using dynamic programming for minimizing the energy of active contours in the presence of hard constraints. In *Computer Vision., Second International Conference on*, pages 95–99. IEEE, 1988.
- [3] S. Avidan and A. Shamir. Seam carving for content-aware image resizing. *ACM Trans. Graph.*, 26(3), July 2007.
- [4] C. Baker, A. Perry, J. Bannister, M. Weinrich, R. B. Abernethy, J. Calambokidis, J. Lien, R. Lambertsen, J. U. Ramirez, and O. Vasquez. Abundant mitochondrial dna variation and world-wide population structure in humpback whales. *Proceedings of the National Academy of Sciences*, 90(17):8239–8243, 1993.
- [5] B. Beekmans, H. Whitehead, R. Huele, L. Steiner, and A. G. Steenbeek. Comparison of two computer-assisted photo-identification methods applied to sperm whales (*physeter macrocephalus*). *Aquatic Mammals*, 31(2):243, 2005.
- [6] S. Berretti, B. B. Amor, M. Daoudi, and A. Del Bimbo. 3d facial expression recognition using sift descriptors of automatically detected keypoints. *The Visual Computer*, 27(11):1021–1036, 2011.
- [7] A. L. Blackmer, S. K. Anderson, and M. T. Weinrich. Temporal variability in features used to photo-identify humpback whales (*megaptera novaeangliae*). *Marine Mammal Science*, 16(2):338–354, 2000.
- [8] T. Branch. Humpback whale abundance south of 60° s from three complete circumpolar sets of surveys. *Journal of Cetacean Research and Management (special issue)*, 3:53–69, 2011.
- [9] J. Calambokidis, E. A. Falcone, T. J. Quinn, A. M. Burdin, P. Clapham, J. Ford, C. Gabriele, R. LeDuc, D. Mattila, L. Rojas-Bracho, et al. *SPLASH: Structure of populations, levels of abundance and status of humpback whales in the North Pacific*. Cascadia Research, 2008.
- [10] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014.

- [11] D. Ciresan, A. Giusti, L. M. Gambardella, and J. Schmidhuber. Deep neural networks segment neuronal membranes in electron microscopy images. In *Advances in neural information processing systems*, pages 2843–2851, 2012.
- [12] J. P. Crall, C. V. Stewart, T. Y. Berger-Wolf, D. I. Rubenstein, and S. R. Sundaresan. HotSpotter - patterned species instance recognition. In *Applications of Computer Vision (WACV), 2013 IEEE Workshop on*, pages 230–237. IEEE.
- [13] F. C. Crow. Summed-area tables for texture mapping. *ACM SIGGRAPH computer graphics*, 18(3):207–212, 1984.
- [14] V. Dumoulin and F. Visin. A guide to convolution arithmetic for deep learning. *arXiv preprint arXiv:1603.07285*, 2016.
- [15] H. Fan, Z. Cao, Y. Jiang, Q. Yin, and C. Doudou. Learning deep face representation. *arXiv preprint arXiv:1403.2802*, 2014.
- [16] P. Fischer and T. Brox. Image descriptors based on curvature histograms. In *Pattern Recognition*, pages 239–249. Springer, 2014.
- [17] K. Fukushima. Neural network model for a mechanism of pattern recognition unaffected by shift in position- neocognitron. *ELECTRON. & COMMUN. JAPAN*, 62(10):11–18, 1979.
- [18] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *Computer vision and pattern recognition, 2006 IEEE computer society conference on*, volume 2, pages 1735–1742. IEEE, 2006.
- [19] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *Computer vision–ECCV 2014*, pages 297–312. Springer, 2014.
- [20] d. J. Hartog and R. Reijns. I3S Contour MANUAL, 2013.
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [22] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- [23] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007.
- [24] R. Huelle, H. U. De Haes, J. Ciano, and J. Gordon. Finding similar trailing edges in large collections of photographs of sperm whales. *Journal of Cetacean Research and Management*, 2(3):173–176, 2000.

- [25] B. Hughes and T. Burghardt. Automated identification of individual great white sharks from unrestricted fin imagery. In M. W. J. Xianghua Xie and G. K. L. Tam, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 92.1–92.14. BMVA Press, September 2015.
- [26] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [27] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pages 2008–2016, 2015.
- [28] M. Kass, A. Witkin, and D. Terzopoulos. Snakes: Active contour models. *International journal of computer vision*, 1(4):321–331, 1988.
- [29] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [30] E. Kniest, D. Burns, and P. Harrison. Fluke matcher: A computer-aided matching system for humpback whale (*megaptera novaeangliae*) flukes. *Marine Mammal Science*, 26(3):744–756, 2010.
- [31] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [32] N. Kumar, P. N. Belhumeur, A. Biswas, D. W. Jacobs, W. J. Kress, I. C. Lopez, and J. V. Soares. Leafsnap: A computer vision system for automatic plant species identification. In *Computer Vision–ECCV 2012*, pages 502–516. Springer, 2012.
- [33] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [34] D. Lemire. Faster retrieval with a two-pass dynamic-time-warping lower bound. *Pattern recognition*, 42(9):2169–2180, 2009.
- [35] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [36] S. A. Mizroch, J. A. Beard, and M. Lynde. Computer assisted photo-identification of humpback whales. *Report of the International Whaling Commission*, 12:63–70, 1990.
- [37] A. Monroy, A. Eigenstetter, and B. Ommer. Beyond straight linesobject detection using curvature. In *Image Processing (ICIP), 2011 18th IEEE International Conference on*, pages 3561–3564. IEEE, 2011.

- [38] M. Müller. *Information retrieval for music and motion*, volume 2. Springer, 2007.
- [39] M. E. Munich and P. Perona. Continuous dynamic time warping for translation-invariant curve alignment with applications to signature verification. In *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*, volume 1, pages 108–115. IEEE, 1999.
- [40] F. Ning, D. Delhomme, Y. LeCun, F. Piano, L. Bottou, and P. E. Barbano. Toward automatic phenotyping of developing embryos from videos. *Image Processing, IEEE Transactions on*, 14(9):1360–1371, 2005.
- [41] D. Nouri. Using convolutional neural nets to detect facial keypoints tutorial, 2014.
- [42] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. *Proceedings of the British Machine Vision*, 1(3):6, 2015.
- [43] H. Pottmann, J. Wallner, Q.-x. Huang, and Y.-l. Yang. Integral invariants for robust geometry processing. 2007.
- [44] Y. Qiao and M. Yasuhara. Affine invariant dynamic time warping and its application to online rotated handwriting recognition. In *Pattern Recognition, 2006. ICPR 2006. 18th International Conference on*, volume 2, pages 905–908. IEEE, 2006.
- [45] H. Sakoe and S. Chiba. Dynamic programming algorithm optimization for spoken word recognition. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 26(1):43–49, 1978.
- [46] S. Salvador and P. Chan. Fastdtw: Toward accurate dynamic time. *Warping in Linear Time and Space*, 2007.
- [47] A. M. Saxe, J. L. McClelland, and S. Ganguli. Exact solutions to the non-linear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [48] 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*, pages 815–823, 2015.
- [49] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
- [50] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

- [51] I. Sobel and G. Feldman. A 3x3 Isotropic Gradient Operator for Image Processing. Never published but presented at a talk at the Stanford Artificial Project, 1968.
- [52] Y. Sun, X. Wang, and X. Tang. Deep convolutional network cascade for facial point detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3476–3483, 2013.
- [53] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [54] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014.
- [55] H. Whitehead. Computer assisted individual identification of sperm whale flukes. *Reports of the International Whaling Commission*, 12:71–77, 1990.

APPENDIX