

## Classification of Bryde's whale individuals using high-resolution time-frequency transforms and support vector machines

Jean Baptiste Tary<sup>1,a</sup>, Christine Peirce<sup>2</sup>, Richard W. Hobbs<sup>2</sup>

<sup>1</sup> *Geophysics section, School of Cosmic Physics, Dublin Institute for Advanced Studies, Dublin, Ireland*

<sup>2</sup> *Department of Earth Sciences, Durham University, Lower Mountjoy, South Road, Durham, DH1 3LE, UK*

### Abstract

Whales generate vocalizations which may, deliberately or not, encode caller identity cues. In this study, we analyze calls produced by Bryde's whales and recorded by ocean-bottom arrays of hydrophones deployed close to the Costa Rica Rift in the Panama basin. These repetitive calls, consisting of two main frequency components at  $\sim 20$  and  $\sim 36$  Hz, have been shown to follow five coherent spatio-temporal tracks. Here, we use a high-resolution time-frequency transform, the 4<sup>th</sup>-order Fourier synchrosqueezing transform (FSST4), to extract time-frequency characteristics (ridges) from each call to appraise their suitability for identifying individuals from each other. Focusing on high-quality calls recorded less than 5 km from their source, we then cluster these ridges using a Support Vector Machine (SVM) model resulting in an average cross-validation error of  $\sim 11\%$  and balanced accuracy of  $\sim 86 \pm 5\%$ . Comparing these results with those obtained using the standard short-time Fourier transform,  $k$ -means clustering, and lower-quality signals, the FSST4 approach, coupled with SVM, substantially improves classification. Consequently, the Bryde's whale calls potentially contain individual-specific information, implying that individuals can be studied using ocean-bottom data.

---

<sup>a</sup> [tary@cp.dias.ie](mailto:tary@cp.dias.ie)

## 24 I. INTRODUCTION

25 As animals interact with each other, they often, intentionally or not, encode individual-  
26 specific information in their communication (e.g., *Janik, 2009*). There can be different causes of  
27 dissimilarities in acoustic signatures between individuals, such as physical characteristics,  
28 environmental and cultural conditions, and temporal changes in these characteristics and conditions  
29 (*Knight et al., 2024*). This information can then be used by the animals for various purposes (e.g.,  
30 territory definition, conspecifics and offspring identification) and for their study (e.g., population  
31 size, migrations, general behavior). For the study of cetaceans, identifying specific individuals is key  
32 to better understanding species' ecology and evolution over time, their environment, and  
33 anthropogenic impacts (e.g., climate changes, water and acoustic pollution, shipping operations etc.).  
34 Various approaches to individual identification have been developed over the years, ranging from  
35 visual observations, animal tags, and acoustic tags. Alternatively, passive monitoring of animal  
36 underwater acoustic communications provides an opportunity to monitor cetaceans for longer time  
37 periods and over larger areas, and especially so if individuals can be identified. Large datasets of  
38 acoustical signatures including individual animal attribution are, however, challenging to obtain for  
39 different reasons, such as the labor-intensive nature of using animal tags, source attribution for  
40 passive monitoring studies, signal deterioration for long-range applications, and background  
41 acoustical conditions.

42 Identifying individuals using acoustic recordings has been applied to a wide range of animals  
43 such as the South Polar skua (*Charrier et al., 2001*) and gorillas (*Salmi et al., 2014*) in addition to  
44 cetaceans, where the latter studies include Bottlenose dolphins (*Janik & Sayigh, 2013*) and Sperm  
45 whales (*Gero et al., 2016*). Sperm whale vocalizations are characterized by complex signals and  
46 temporal patterns, which define both vocal clans at the scale of an ocean basin and individuals (*Gero*  
47 *et al., 2016; Oliveira et al., 2016; Bermant et al., 2019*). In the case of baleen whales, fewer studies of

48 individual identification exist, and these are mainly of Humpback whales. However, identifying  
49 individuals is challenging for various reasons, ranging from signal source attribution (*Zeh et al., 2024*)  
50 to limited knowledge of their vocal repertoire (e.g., *White & Todd, 2024*), where characteristics are  
51 also species dependent. In the case of Humpback whales, individuals have been identified using  
52 song cepstral content together with a Support Vector Machine (SVM) model (*Mazhar et al., 2007*),  
53 the use of some signal units in songs and their combinations (*Lamoni et al., 2023*), and call temporal  
54 patterns and amplitudes (*Zeh et al., 2024*). In addition, *McDonald et al. (2001)* suggested that some  
55 specific frequency characteristics of A-B calls of Blue whales could be used to identify individuals,  
56 and that frequency features extracted from spectrograms could provide information on individual  
57 North Atlantic right whales (*McCordic et al., 2016*).

58 In this study, we focus on whale calls recorded by Ocean-Bottom Seismographs (OBS) and a  
59 vertical array (VA) of hydrophones deployed in the Panama Basin in January and February, 2015  
60 (*Hobbs & Peirce, 2015; Tary et al., 2024*) ([Figure 1](#)). The calls under consideration are very similar,  
61 short in duration (~3-5 s-long) and consist of two main frequency components: a ~1 s-long  
62 component at ~36 Hz and a ~3 s-long component at ~20 Hz ([Figure 2](#)). These calls are identified as  
63 likely corresponding to Be1 calls attributed to Bryde's whales by *Oleson et al. (2003)*. In this region of  
64 the Eastern Tropical Pacific Ocean, Bryde's whales are common despite their low abundance (*Wade*  
65 *& Gerrodette, 1993; Palacios et al., 2012*). The location of the calls generated by whales within the OBS  
66 network have lateral uncertainties less than a few kilometers (*Tary et al., 2024*). Some of these calls  
67 occur in spatio-temporal sequences and form trajectories across the network ([Figure 1](#)).

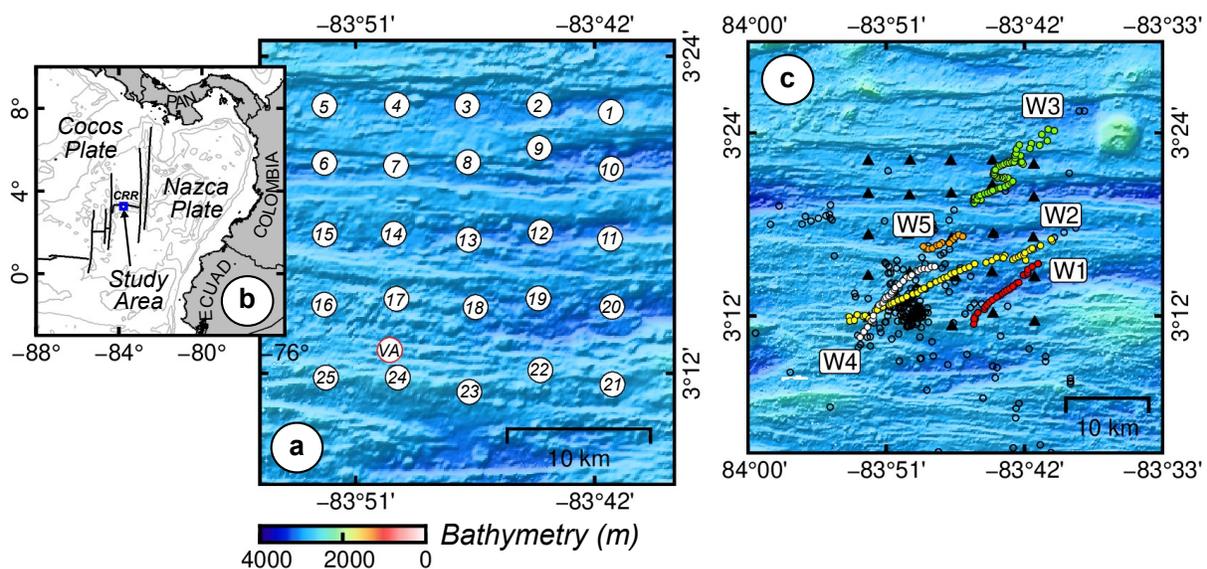
68 Bryde's whales generally travel as individuals or in pairs, and rarely in larger groups. With the  
69 call localization not having the resolution to distinguish between collocated whales (i.e., whales  
70 separated by distances on the order of 10s to 100s of meters), we assume that each of these  
71 trajectories corresponds to a different individual whale and determine if the associated calls can be

72 classified into different groups based on their time-frequency features. If so, this could indicate that  
73 individual information is encoded within these relatively simple calls.

74 The calls under investigation are relatively short with simple characteristics, which contrasts  
75 with other calls used for individual identification such as those of Humpback whales (*White & Todd,*  
76 *2024*). For short and simple calls, general attributes would likely not provide sufficient information  
77 to distinguish between individuals, considering other causes of variability such as animal behaviors,  
78 wave propagation effects, and the impact of background noises (natural or anthropogenic). In order  
79 to capture several attributes at the same time (e.g., signal component durations, mean frequencies,  
80 frequency modulations), here we employ time-frequency representations to classify these calls.

81 Whale calls are generally analyzed using their spectrogram (e.g., *Mellinger & Clark, 2000*). In order to  
82 improve the definition of time-frequency information extracted from each call, instead of using the  
83 full time-frequency representations, we extract time-frequency ridges from representations obtained  
84 using a variant of the short-time Fourier transform (STFT), called the high-order synchrosqueezing  
85 transform (FSSTN - *Pham & Meignen, 2017*).

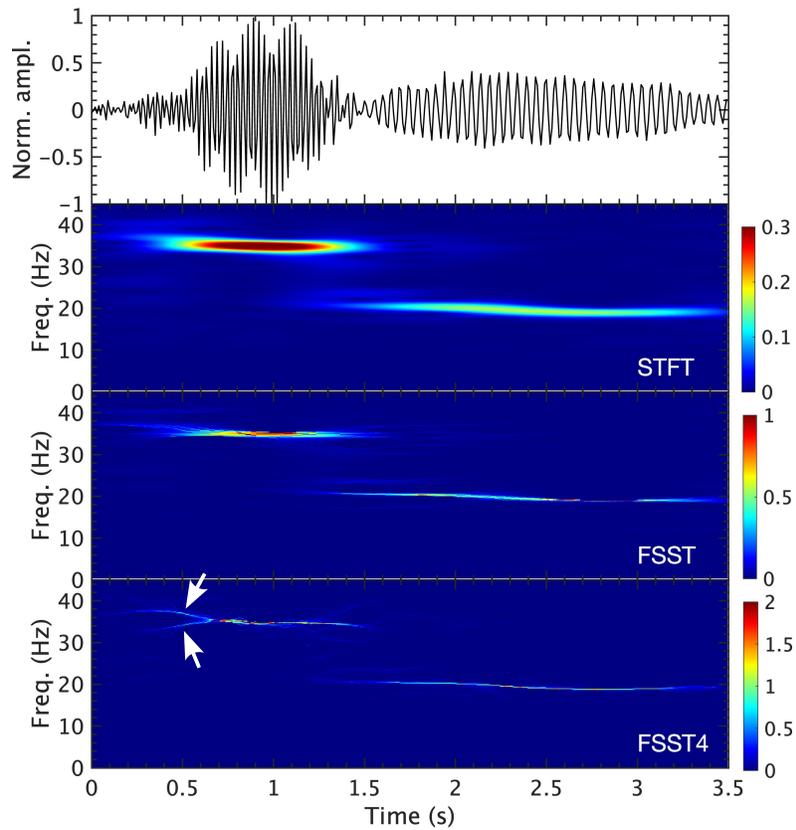
86



87

88 FIG. 1. (Color online). a) Bathymetry map of the survey area showing the ocean-bottom  
89 seismograph (OBS - numbered circles) and vertical hydrophone array (VA - red circle) positions. b)  
90 Location of the survey area (blue rectangle) in the Panama basin, close to the Costa Rica Rift (CRR)  
91 spreading ridge boundary between the Cocos and Nazca tectonic plates. c) Whale call locations from  
92 *Tary et al. (2024)* (open circles), with color-coded circles corresponding to the high-quality whale calls  
93 included in the support vector machine classification (whale tracks 1: red, 2: yellow, 3: green, 4:  
94 white, and 5: orange). OBS locations are indicated by the black triangles.

95



96

97 FIG. 2. (Color online). Whale call recorded by OBS 16 (see [Figure 1](#) for location), contained  
98 within whale track 2 ([Figure 1c](#)), at 15:43:30 on January 29, 2015, and its time-frequency  
99 representations obtained using the short-time Fourier transform (STFT – Gaussian window of

100  $\sim 0.36$  s at half-maximum with 97% overlap), the Fourier synchrosqueezing transform (FSST; 4th-  
101 order – FSST4). The white arrows locate the fork in frequency discussed in Section III.

102

103 Machine learning methods are generally applied to the detection and classification of whale  
104 calls (e.g., *Halkias et al., 2013; Ibrahim et al., 2021; Rasmussen & Širović, 2021; Zhong et al., 2021; Kather et*  
105 *al., 2024*), but seldom to caller identification (*Rendell & Whitehead, 2003; Mahzar et al., 2007; Bermant et*  
106 *al., 2019*) which is usually determined through call statistical analysis. Classifying observations in  
107 different categories using machine learning can be realized using unsupervised and supervised  
108 methods (e.g., *Bergen et al., 2019*). In our case, the whale locations can be transformed into labels to  
109 classify the signals using supervised methods. Of the different existing supervised methods, we  
110 employ support vector machines (SVM) for its demonstrated high performance in various  
111 applications (e.g., *Cervantes et al., 2020*), high generalization ability, and resistance to outliers and  
112 overfitting, even for high-dimensional data (*Kecman, 2005*). The SVM method, originally developed  
113 for binary classification, is a large margin classifier aimed at determining the optimal boundary  
114 between a subset of observations called support vectors (*Cortes & Vapnik, 1995*). To define non-  
115 linear decision boundaries between classes, the original data is often mapped to a higher-dimensional  
116 space, called the feature space, using kernels. In the present case, we demonstrate that different  
117 Bryde’s whales can be distinguished using a SVM classifier using the time-frequency content of their  
118 brief calls recorded in ocean-bottom data.

119

## 120 II. DATA AND METHODS

121 Between January 26, 2015, and February 17, 2015, 25 OBSs and a VA of 12 hydrophones  
122 were deployed close to the Costa Rica Rift in the Panama basin (*Hobbs & Peirce, 2015*) ([Figure 1](#)).

123 This grid of instruments is approximately 20 x 20 km wide, with an instrument spacing of around 5  
124 km. Apart from five OBSs (4, 7, 14, 17, 24) and the VA which recorded during the complete  
125 deployment, the remaining OBSs recorded from January 26, 2015 to their recovery time on February  
126 1 or 2, 2015. The OBSs and VA were equipped with a High-Tech HTI-90-U hydrophone, while  
127 each OBS also had a 3-component short-period geophone package (Sercel L-28 4.5 Hz). Each time  
128 series was sampled at 500 Hz.

129 This dataset was first analyzed to study the structure of the oceanic lithosphere around the  
130 Costa Rica Rift (e.g., *Wilson et al., 2019; Robinson et al., 2020*). It was then re-examined to study the  
131 microseismicity (*Lowell et al., 2020; Tary et al., 2021*) and Bryde's whale calls (*Tary et al., 2024*)  
132 observed in this region. Focusing on the whale calls, two types of calls were observed; a repetitive  
133 call of ~4 s and a less common call consisting of brief signals of 0.5-1 s duration. The main  
134 characteristics of the most common, repetitive call, consist of two main signal components; a first  
135 wave packet of ~1 s duration centered at ~36 Hz, followed by a generally lower-amplitude, longer  
136 signal of ~3 s duration centered at ~20 Hz ([Figure 2](#)). These calls were then detected using two  
137 different methods; an energy method based on the short-term over long-term average ratio  
138 (STA/LTA), and template matching using the subspace detector applied to the calls detected by the  
139 first method. The arrival times of each call at the different instruments were then manually identified  
140 and all calls were located using a measurement-based 1D velocity model of the water column and  
141 the non-linear, probabilistic method implemented in NonLinLoc (*Lomax et al., 2001*). The calls were  
142 finally relocated relative to each other using the double-difference technique (*Waldhauser & Ellsworth,*  
143 *2000*) ([Figure 1c](#)).

144 In order to examine the whale call characteristics and differences between individuals, we  
145 first focus on the whale calls that are well-recorded by the instruments (call-to-instrument distance <  
146 5 km, calls at stations with manually identified start times), and observed in a single time period

147 following a coherent spatial movement. Using these criteria, we identify five whale tracks observed  
148 at various times during the deployment (Figure 1c), whale track 1 with 27 calls recorded 71 times  
149 (January 29, 2015, some calls being recorded by more than one station), whale track 2 with 52 calls  
150 recorded 133 times (January 29, 2015), whale track 3 with 65 calls recorded 199 times (January 29,  
151 2015), whale track 4 with 30 calls recorded 89 times (February 2, 2015), and whale track 5 with 12  
152 calls recorded 38 times (January 26, 2015). All recorded calls for all whale tracks (i.e., 530  
153 waveforms) are then included in the classification. The whale calls are first extracted using a longer  
154 time window of 4.5 s, and aligned using the manually identified start times. They are then re-aligned  
155 using waveform cross-correlation, relative to the event which is the most correlated to all events on  
156 average. Reviewing the call waveforms, the calls were then re-cut to 3 s duration from the call start  
157 times to both include the two main signal components and reject later signal parts with lower signal-  
158 to-noise ratio and/or other wave arrivals. The calls are finally down-sampled to a sampling rate of  
159 100 Hz and band-pass filtered between 10 and 45 Hz.

160

#### 161 **A. Time frequency analysis: high-order synchrosqueezing transforms**

162 The synchrosqueezing transform (SST) is a time-frequency representation which improves  
163 the readability of some time-frequency representations, such as the Continuous Wavelet Transform  
164 (CWT), by reassigning non-zero time-frequency coefficients to previously determined instantaneous  
165 frequencies (IF) (e.g., *Daubechies et al., 2011*). The main purpose of this synchrosqueezing operation is  
166 to significantly reduce frequency smearing (e.g., *Tary et al., 2014*). This has been applied to different  
167 time-frequency transforms such as the STFT (hereafter referred to as the FSST - *Thakur & Wu,*  
168 *2011*) and the S-transform (*Huang et al., 2015*).

169 The SST was originally developed for slowly-varying, well-separated frequency components  
170 (*Daubechies et al., 2011*). However, in the case of the STFT, higher-order SSTs were developed,

171 improving the time-frequency concentration and mode reconstruction of the FSST for strongly  
 172 amplitude-modulated and frequency-modulated multi-component signals (*Oberlin et al., 2015; Pham*  
 173 *& Meignen, 2017*). In this case, we first consider a signal  $s(t)$  that can be decomposed into a series of  
 174  $K$  frequency components as

$$175 \quad s(t) = \sum_{k=1}^K A_k(t) e^{i2\pi\theta_k(t)} + \varepsilon(t), \quad (1)$$

176 where  $A_k(t)$  and  $\theta_k(t)$  are the time-varying instantaneous amplitude and phase of the  $k$ th  
 177 component, respectively, and  $\varepsilon(t)$  is time-varying random noise. Instantaneous frequencies are then  
 178 estimated using

$$179 \quad \hat{\omega}(t, \eta) = \frac{1}{2\pi} \frac{\partial \arg S_F(t, \eta)}{\partial t}, \quad (2)$$

180 where  $\arg S_F(t, \eta)$  is the argument of the complex valued STFT representation  $S_F(t, \eta)$  at time  $t$   
 181 and frequency  $\eta$ . In order to limit the reassignment of noise components, only non-zero frequency  
 182 components above a pre-defined threshold  $\zeta$  are reassigned on  $(t, \hat{\omega}(t, \eta))$  positions following

$$183 \quad T_F^\zeta(t, \omega) = \frac{1}{g^*(0)} \int_{\{\eta, |S_F(t, \eta)| > \zeta\}} S_F(t, \eta) \delta(\omega - \hat{\omega}(t, \eta)) d\eta, \quad (3)$$

184 where  $\delta$  is the Dirac distribution and  $g^*$  the complex conjugate of the window function  $g$ . Here, we  
 185 focus on the main aspects of the method presented by *Pham & Meignen (2017)*. Focusing on signal  
 186 modes of  $s$  having non-negligible phase derivatives  $\theta^{(n)}(t)$  for  $n \geq 3$ , using a high-order Taylor  
 187 expansion of eq. 1 in  $\tau$  close to  $t$  for a mode amplitude and phase gives

$$188 \quad s(\tau) = \exp \left( \sum_{n=0}^N \frac{1}{n!} \left( [\log(A)]^{(n)}(t) + i2\pi\theta^{(n)}(t) \right) (\tau - t)^n \right), \quad (4)$$

189 where  $Z^{(n)}(t)$  is the  $n^{\text{th}}$  derivative of  $Z$  at  $t$ , and  $N$  is the order of the Taylor expansion of phase  
 190  $\theta(\tau)$ . Modifying the STFT representation  $S_F(t, \eta)$  as well as the IF  $\hat{\omega}(t, \eta)$  accordingly, this  
 191 requires the estimation of a frequency modulation operator  $\tilde{q}^{[n, N]}$  and leads to the following  
 192 definition of a  $N^{\text{th}}$  order IF at time  $t$  and frequency  $\eta$

$$193 \quad \tilde{\omega}_\eta^{[N]}(t, \eta) = \begin{cases} \tilde{\omega}(t, \eta) + \sum_{n=2}^N \tilde{q}_\eta^{[n, N]}(t, \eta) (-x_{n,1}(t, \eta)), & \text{if } S_F(t, \eta) \neq 0 \text{ and } \partial_\eta x_{j,j-1}(t, \eta) \neq 0, 2 \leq j \leq N \\ \tilde{\omega}(t, \eta) & \text{otherwise} \end{cases}, \quad (5)$$

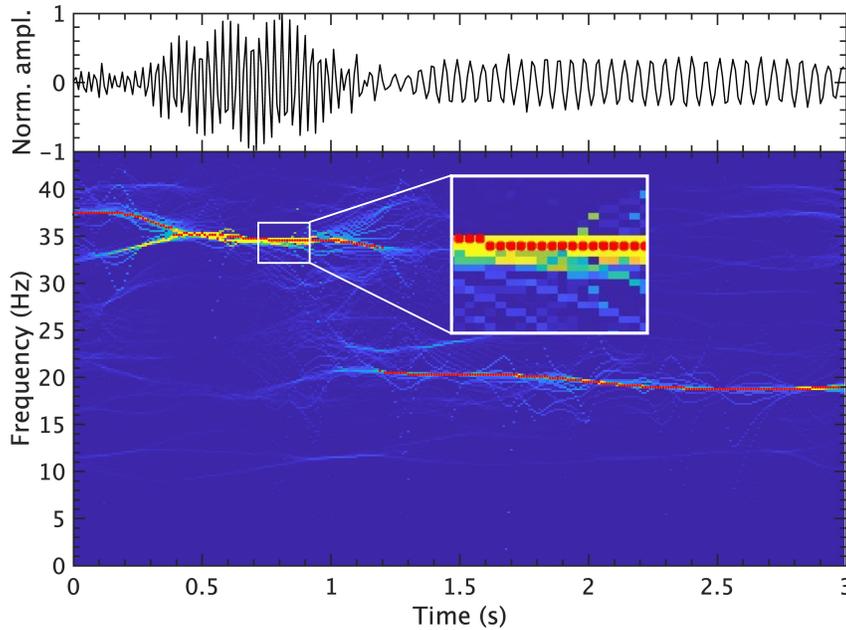
194 with

$$195 \quad \begin{cases} x_{n,1}(t, \eta) = \frac{S_F^{t^{n-1}}(t, \eta)}{S_F(t, \eta)} \text{ for } 1 \leq n \leq N \\ x_{n,j}(t, \eta) = \frac{\partial_\eta x_{n,j-1}(t, \eta)}{\partial_\eta x_{j,j-1}(t, \eta)} \text{ for } 2 \leq j \leq N \text{ and } j \leq n \leq N \end{cases}. \quad (6)$$

196 The real part of  $\tilde{\omega}^{[N]}(t, \eta)$  is incorporated into eq. 3 for  $T_F^\zeta(t, \omega)$  to obtain the  $N^{\text{th}}$  order  
 197 FSST. Frequency ridges are then extracted from the resulting time-frequency representations by  
 198 iteratively searching for energy maxima in the time-frequency plane (*Meignen et al., 2017*). The  
 199 threshold parameter  $\zeta$  on the STFT representation is set to a small value of 0.001 in order to avoid  
 200 removing any signal components. The window function used to calculate the STFT is a Gaussian  
 201 function  $g = \sigma^{-1} e^{-\frac{\pi t^2}{\sigma^2}}$ . The parameter  $\sigma$  controls the width of the Gaussian window and, hence,  
 202 the time and frequency localizations of the STFT (e.g., *Tary et al., 2014*). Its value in our case is 0.11,  
 203 corresponding to the minimum of the Rényi entropy of STFT representations of the different whale  
 204 calls presented in this study (*Stanković, 2001*).

205 To train the SVM model, two ridges are extracted from the  $4^{\text{th}}$  order SST representations  
 206 (FSST4), corresponding to the two primary IF components of the whale calls (*Figure 2*). To  
 207 transform each whale call into a one-dimensional vector, we only keep the IF of maximum  
 208 amplitude for all time samples. As the numbers of training examples per whale is relatively limited,

209 we also reduce the number of training attributes by down-sampling this ridge frequency vector by  
210 three (i.e., one point every 0.03 s) resulting in 100 IF measures per whale call (Figure 3). Finally, the  
211 IF data is normalized to obtain values ranging between -1 and 1.  
212



213  
214 FIG. 3. (Color online). (top) Whale call recorded by OBS 24 at 15:43:30 on January 29, 2015,  
215 and contained within whale track 2. (bottom) The time-frequency representation obtained using the  
216 FSST4, together with the extracted ridge values (red dots), of which one in three values were used  
217 for SVM classification.

218

## 219 B. Classification support vector machine model

220 The SVM method is a supervised learning method that can be used for multi-class  
221 classification and regression (Boser *et al.*, 1992; Vapnik, 1995). In general terms, the SVM method  
222 seeks to separate training examples based on their features, or these features after (non-)linear  
223 mapping, in a number of classes using the largest margin between some of the training examples

224 located close to this margin. These margin-defining training examples are called support vectors and  
 225 they define two optimal hyperplane positions separating the two classes. A multi-class classification  
 226 using SVM is generally obtained using a series of two-class SVMs and combining their classification  
 227 results. An SVM model can then be used to “predict” which class a new training example would  
 228 belong to.

229 Using a training dataset consisting of  $M$  training examples  $\{\mathbf{x}_i, y_i\}, i = 1, \dots, M$ , each  
 230 training example  $\mathbf{x}_i$  is a vector of attributes and has a corresponding class label  $y_i$  ( $y_i \in \{-1, 1\}$ ).

231 For linearly-separable data, a vector  $\mathbf{w}$  and a scalar  $b$  exist so that

$$232 \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, M, \quad (7)$$

233 is defining an optimal hyperplane defined by

$$234 \quad \mathbf{w} \cdot \mathbf{x} + b = 0, \quad (8)$$

235 that separates the data points into two classes of  $y_i$  equal to 1 and -1 with the widest margin  
 236 (*Cortes & Vapnik, 1995*). The geometrical distance of the training examples to the hyperplane is  
 237 given by

$$238 \quad \Delta_i = \frac{y_i(\mathbf{w} \cdot \mathbf{x}_i + b)}{\|\mathbf{w}\|} \geq \frac{1}{\|\mathbf{w}\|}, \quad (9)$$

239 where  $\|\mathbf{w}\|$  is the  $\ell_2$ -norm of  $\mathbf{w}$ . Finding the optimum hyperplane corresponds to  
 240 maximizing  $\Delta_i$  for training examples close to the hyperplane or, equivalently, minimizing  $\|\mathbf{w}\|$ . The  
 241 primal optimization problem can then be expressed as

$$242 \quad \min_{\gamma, \mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (10)$$

$$\text{s. t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, M.$$

243 Using the Lagrangian of this optimization problem we obtain

$$244 \quad \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^M \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1], \quad (11)$$

245 where  $\alpha_i$  are the Lagrange multipliers corresponding to each training example (*Cortes &*  
 246 *Vapnik, 1995*). Minimizing  $\mathcal{L}(\mathbf{w}, b, \alpha)$  implies that  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$  and results in the following dual  
 247 optimization problem

$$\begin{aligned}
 \max_{\alpha} W(\alpha) &= \sum_{i=1}^M \alpha_i - \frac{1}{2} \sum_{i,j=1}^M y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\
 \text{s. t. } \alpha_i &\geq 0, \quad i = 1, \dots, M \\
 \sum_{i=1}^M \alpha_i y_i &= 0.
 \end{aligned}
 \tag{12}$$

249 Instead of directly using the training example attributes in  $\mathbf{x}_i$ , they can be transformed to a  
 250 higher dimensional feature space using function  $\phi(\mathbf{x}_i)$ , modifying  $\mathbf{w} = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)$  and replacing  
 251 the inner product  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$  in eq. 12 by  $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$  which corresponds to the definition of a  
 252 kernel (*Vapnik, 1995*). Most common are the linear kernel, the radial basis function (RBF) or  
 253 Gaussian kernel, and the polynomial kernel. In addition, for variables that are not linearly separable,  
 254 or in the case of data errors, it might not be possible or desirable to obtain a hyperplane that takes  
 255 into account all training examples equally.

256 To overcome these limitations, the optimization problem can be modified to use soft  
 257 margins controlled by a boundary parameter. In this case, using  $\ell_1$  regularization, the dual  
 258 optimization problem of eq. 12 becomes

$$\begin{aligned}
 \max_{\alpha} W(\alpha) &= \sum_{i=1}^M \alpha_i - \frac{1}{2} \sum_{i,j=1}^M y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\
 \text{s. t. } 0 &\leq \alpha_i \leq C, \quad i = 1, \dots, M \\
 \sum_{i=1}^M \alpha_i y_i &= 0,
 \end{aligned}
 \tag{13}$$

260 where the upper bound parameter  $C$  defines the maximum penalty on training examples  
 261 close to the margin boundaries. Besides SVM classifier optimization, we then have different

262 hyperparameters to define to improve the results, namely the choice of kernel function, the upper  
263 bound parameter  $C$ , the kernel scaling parameter  $\gamma$  for the Gaussian kernel, and the polynomial  
264 order for the polynomial kernel function.

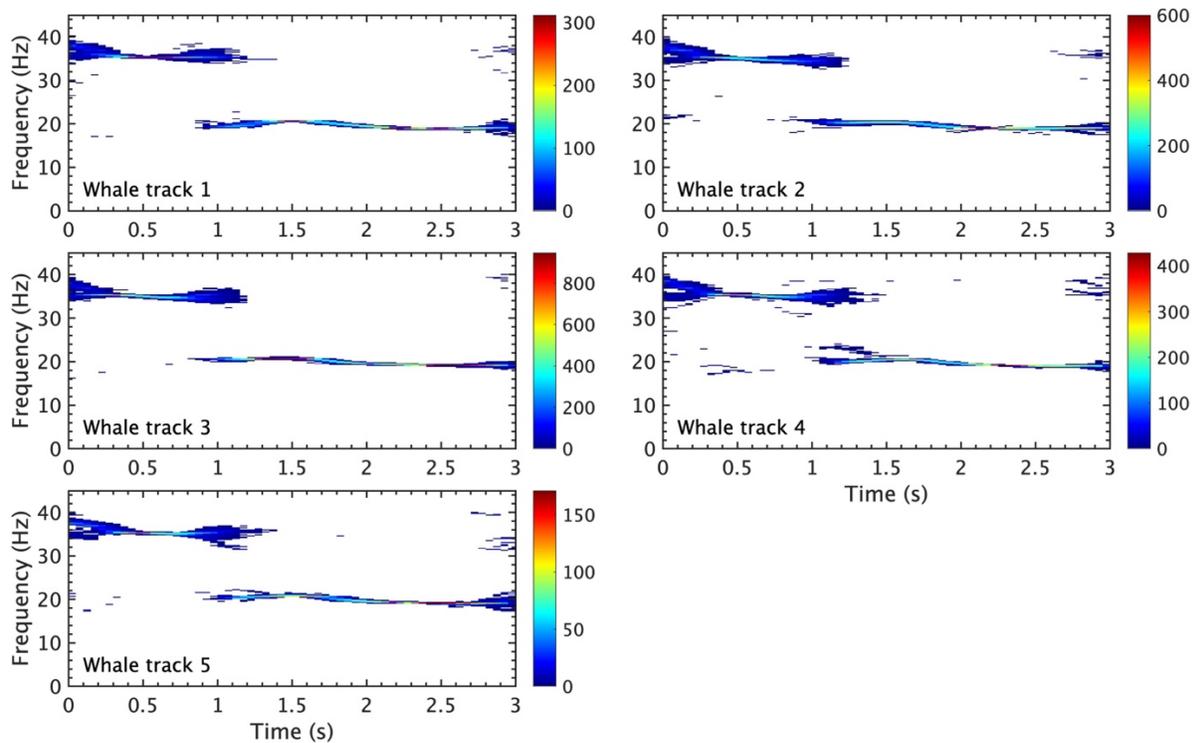
265 In this study, all data examples are first randomly shuffled and then partitioned into a  
266 training and a test set using a 90% - 10% split. To limit the influence of a particular split on the final  
267 results, we run 50 instances of this procedure using different randomizations and splits. The  
268 classification error is then estimated using the average balanced accuracy and its standard deviation,  
269 taking into account class number imbalance, and average cross-validation classification error  
270 (average of 4-fold, 5-fold and 10-fold cross-validation for the 50 runs). The statistical significance of  
271 all model classifications is assessed using permutation tests (*Ojala & Garriga, 2010; Combrisson &*  
272 *Jerbi, 2015*), which compare classifier performances using the original data to its performance using  
273 randomly permuted class labels (i.e., whale track numbers). The model is trained on the training  
274 dataset in the same way as the original model, and its performance measured by its balanced  
275 accuracy on the test set. This procedure is repeated 1000 times to determine the 99.9% percentile  
276 threshold of the balanced accuracy distribution, obtaining a significance level at  $p < 0.001$ .

277

### 278 III. RESULTS

279 The time-frequency content of whale calls is a key element for whale species identification  
280 and, hence, is useful for other purposes such as for their detection. The time-frequency  
281 representation generally in use is the STFT or the spectrogram. [Figure 2](#) shows a Bryde's whale call  
282 example with high signal-to-noise ratio, together with its STFT, FSST and FSST4 representations.  
283 As visible in this example, the well-defined, quasi-harmonic components of these whale calls are  
284 highly suited for analysis using the SST (e.g., *Daubechies et al., 2011*). The frequency reassignment

285 sharpens the two main time-frequency components of the signal and, thus, improves the readability  
286 of the representation. Comparing the FSST4 results with those of the other two transforms, the  
287 frequency resolution is higher and lower-amplitude frequency modulations are better delineated  
288 using this method. For example, in Figure 2 the fork in frequency around time 0.5 s is only clearly  
289 visible in the FSST4 representation. This likely arises from the ability of the FSST4 to better handle  
290 frequency modulated signals relative to the FSST (Pham & Meignen, 2017), which translates into a  
291 better estimation of the time-frequency ridges (Figure 3). The ridges compiled from all calls recorded  
292 by the OBSs located less than 5 km away from the whale location show that their time-frequency  
293 attributes exhibit only slight variations (Figure 4).  
294



295  
296 FIG. 4. (Color online). Density plots showing the ridges extracted from FSST4  
297 representations for all high-quality whale calls contained in the five whale tracks (i.e., calls recorded  
298 within 5 km of the OBSs for periods showing coherent movements). These representations show

299 the two main signal components (i.e., at  $\sim 36$  and  $\sim 20$  Hz), and the similarities and potential  
300 differences of the signals between whale tracks.

301

302 For the first signal component at  $\sim 36$  Hz, more variability is present for the lower amplitude  
303 parts of the signal at the beginning and end of the component. For most of the calls, the higher  
304 amplitude branch at the beginning of the signal shows increasing frequency between  $\sim 36$  Hz and  
305  $\sim 40$  Hz. Less variability is observed for the second signal component. A number of time-frequency  
306 ridges are present at  $\sim 36$  Hz for the second signal component, which likely correspond to calls with  
307 noisy components of lower amplitudes. Another source of signal variability is the timing of the end  
308 and beginning of the two signal parts between  $\sim 1$  and  $\sim 1.5$  s. Each of these variabilities may play a  
309 role in the ability of the classification to distinguish between different whale tracks.

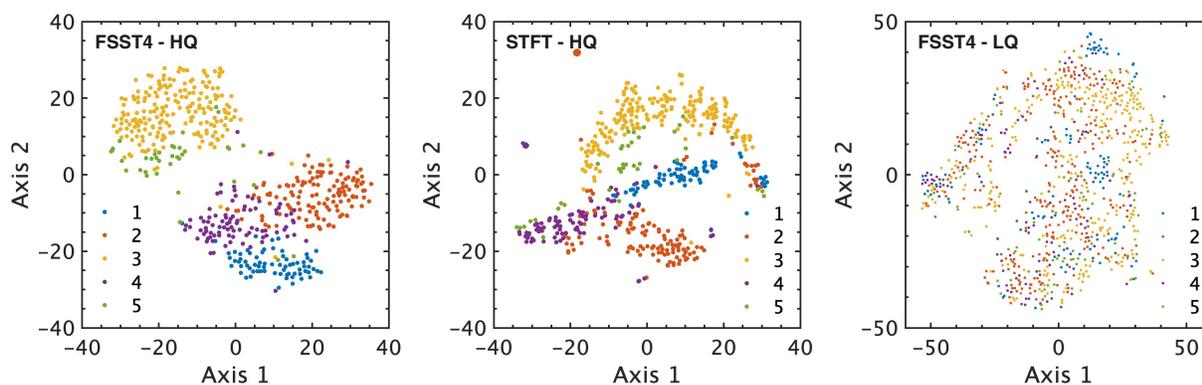
310 To better demonstrate the advantages of the FSST4 and SVM, we consider the following  
311 four cases:

- 312 i) SVM classification using ridges extracted from FSST4 representations,
- 313 ii) SVM classification using ridges extracted from STFT representations,
- 314 iii)  $k$ -means classification using ridges extracted from FSST4 representations, and
- 315 iv) SVM classification using ridges extracted from FSST4 representations for lower-quality calls  
316 recorded at distances between 5 and 10 km from the whale location, enabling evaluation of  
317 the influence of signal quality on the classification performance.

318 The *a priori* clustering of high-dimensional data can be visualized using the t-Distributed  
319 Stochastic Neighbor Embedding (t-SNE) method (*van der Maaten & Hinton, 2008*), which performs a  
320 non-linear mapping of the high-dimensional data to lower dimensions. The t-SNE method is applied  
321 to the aforementioned cases using  $\ell_1$  distance as similarity metric and a perplexity of 30 ([Figure 5](#)).  
322 The t-SNE visualization shows that the best cluster separation is obtained using the FSST4 together

323 with the high-quality signals (calls observed less than 5 km away from its source). No clear clusters  
 324 are visible in the case of the lower-quality signals (calls observed 5 to 10 km away from its source).  
 325 Interestingly, while calls from whale tracks 3 and 5 are well separated from the other calls, whale  
 326 track 1 is slightly separated from whale tracks 2 and 4, and some mixing occurs for whale tracks 2  
 327 and 4 and for whale tracks 3 and 5.

328 The SVM classification models all use Gaussian kernels, the other hyperparameters being  
 329 defined through 500 iterations of Bayesian optimization instead of grid search to reduce training  
 330 time. Allowed values for the upper bound parameter  $C$  and the kernel scaling parameter  $\gamma$  range  
 331 between 0.0001 and 10000. In the case of the  $k$ -means algorithm, the number of clusters is set to 5  
 332 and their initial centroids are set using the  $k$ -means++ algorithm (*Arthur & Vassilvitskii, 2007*). The  
 333 final cluster centroids are then obtained by repeating five times the minimization of the sum of  
 334 absolute distances between data points and centroids ( $\ell_1$  distance) using different initializations,  
 335 keeping the centroids corresponding to the minimum total distance. This procedure reduces the  
 336 probability of obtaining centroids corresponding to a local minimum far from the global minimum.  
 337



338  
 339 FIG. 5. (Color online). t-SNE visualization of the time-frequency ridge data, color-coded by  
 340 whale track number (see [Figure 1c](#)), for ridges extracted using the FSST4 and high-quality signals  
 341 (left), the STFT and high-quality signals (middle), and the FSST4 and lower-quality signals (right).

342 All three visualizations use a perplexity of 30 and  $\ell_1$  distance, and are similar to those obtained using  
343 higher perplexity values.

344

345 The different classification results are presented in Table I. For SVM using the FSST4 and  
346 high-quality calls observed at less than 5 km from the source (Gaussian kernel, upper bound  $C$  of  
347  $\sim 38.1$ , kernel scale  $\gamma$  of  $\sim 25.0$ ), we obtain a training error of  $\sim 0.4\%$  for the model with the highest  
348 balanced accuracy on the test set (2 calls misclassified out of the 477 calls in the training set), a  
349 training average cross-validation error of  $\sim 11\%$ , and a balanced accuracy of  $\sim 86 \pm 5\%$  on the test  
350 set (chance  $\sim 25\%$  at  $p < 0.001$ ). The confusion chart corresponding to this model shows that whale  
351 tracks 1 and 3 are associated with the least misclassification errors (i.e.,  $\sim 8\%$  and  $\sim 5\%$ , respectively),  
352 whereas whale track 5 is associated with the highest misclassification error ( $\sim 35\%$ ) (Figure 6). Using  
353 the STFT instead of the FSST4 slightly increases both the number of incorrectly classified whale  
354 calls during training and the training average cross-validation error to  $14\%$  (Gaussian kernel, upper  
355 bound  $C$  of  $\sim 51.5$ , kernel scale  $\gamma$  of  $\sim 10.2$ ), and decreases the average balanced accuracy to  $\sim 78$   
356  $\pm 8\%$  on the test set (chance  $\sim 25\%$  at  $p < 0.001$ ). The lowest misclassification errors are obtained  
357 for whale tracks 2 and 3, the highest classification error being associated with whale track 5 ( $\sim 46\%$ ).  
358 Using lower quality signals with SVM and the FSST4 (Gaussian kernel, upper bound  $C$  of  $\sim 21.6$ ,  
359 kernel scale  $\gamma$  of  $\sim 9.2$ ), the classification error is  $\sim 1\%$  with a higher average cross-validation error of  
360  $\sim 38\%$  and an average balanced accuracy of  $\sim 53 \pm 5\%$  on the test set (chance  $\sim 27\%$  at  $p < 0.001$ ).  
361 This might indicate that more overfitting is present in this model. In this case, the misclassification  
362 error is greater than  $25\%$  for all whale tracks, with whale track 4 having a misclassification error  
363 reaching  $\sim 78\%$ . Lastly, using  $k$ -means and ridges extracted from high-quality calls using FSST4, we  
364 obtain a training error of  $\sim 39\%$  and a balanced accuracy of  $\sim 64\%$  (chance  $\sim 64\%$  at  $p < 0.001$ ). The

365 misclassification error is consistently over 30% for all whale tracks, with a maximum of ~49% for  
 366 whale track 3 (Figure 6).

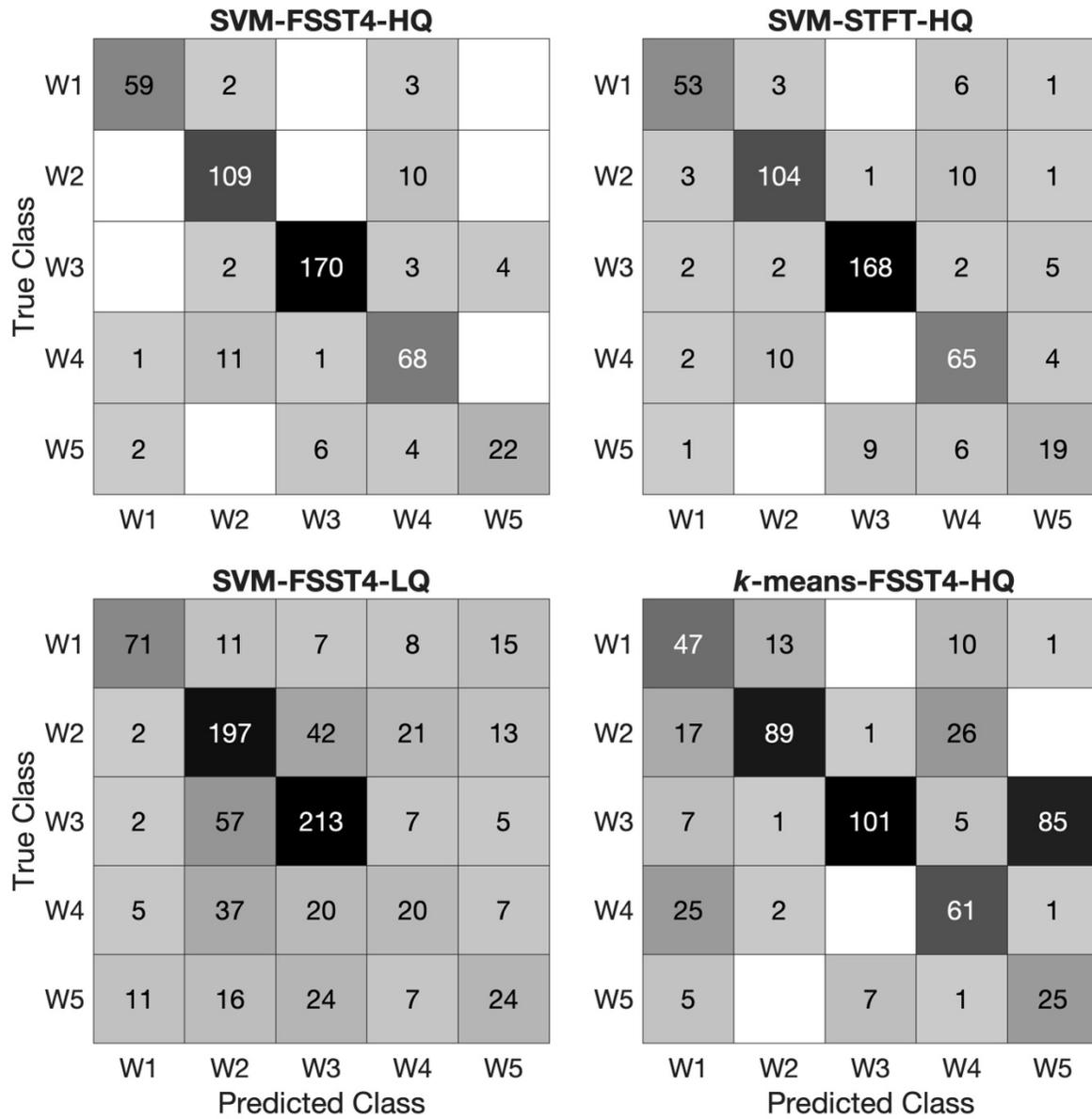
367

368 Table I. Classification results for the different cases using either SVM or  $k$ -means, FSST4 or  
 369 STFT, and high-quality (HQ, calls observed at less than 5 km from its source) or lower-quality (LQ,  
 370 calls observed at distances between 5 and 10 km from its source) signals. For each combination, the  
 371 training classification error for the model with the highest balanced accuracy on the test set is  
 372 shown, together with the average training cross-validation error (C-V error, average of 4-fold, 5-fold  
 373 and 10-fold cross-validation for the 50 runs) for SVM, and the balanced accuracy (Bacc, average and  
 374 standard deviation over the 50 runs for SVM).

	Training error (%)	C-V error (%)	Bacc (%)
SVM + FSST4 + HQ signals	0.4	11	86 $\pm$ 5
SVM + STFT + HQ signals	2	14	78 $\pm$ 8
SVM + FSST4 + LQ signals	1	38	53 $\pm$ 5
$k$ -means + FSST4 + HQ signals	39		64

375

376



377

378

FIG. 6. (Color online). Multiclass confusion matrices obtained from cross-validating the

379

SVM models having the highest balanced accuracy on the test set and the  $k$ -means results for the

380

cases listed in Table I. Rows and columns of each matrix contain the number of calls in their actual

381

class and the number of calls that were classified in each class by the model, respectively. W1, W2,

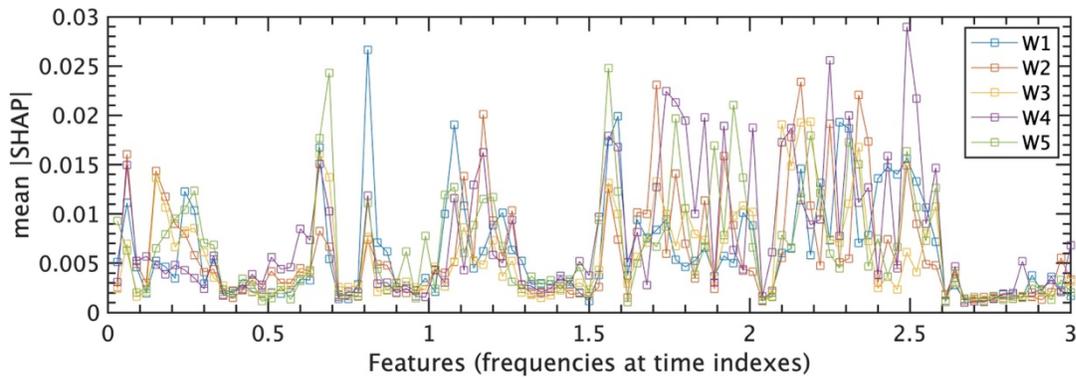
382

W3, W4 and W5 correspond to whale tracks 1, 2, 3, 4 and 5, respectively.

383

384 Focusing on the SVM model using FSST4 and high-quality calls, we calculate SHAP  
 385 (SHapley Additive exPlanations) values (*Lundberg & Lee, 2017*) in order to determine which features  
 386 (i.e., frequency values at different time indexes) have the most effect on the final model classification  
 387 (*Figure 7*). Shapley values quantify the changes in model classification due to a feature, for a given  
 388 training example. Repeating this process for all features considering all training examples, the mean  
 389 absolute Shapley values can be used to calculate SHAP feature importance; larger values  
 390 corresponding to features more significant for the model. SHAP values can, however, be biased due  
 391 to feature correlation, feature interaction and the small size of the training dataset (*Molnar, 2022*).  
 392 Overall, the most significant features for the model classification are similar for the different whale  
 393 tracks (*Figure 7*). Greater SHAP values are obtained at times from the beginning to  $\sim 0.3$  s,  $\sim 0.6-0.7$   
 394 s,  $\sim 0.8-0.9$  s,  $\sim 1-1.3$  s and most of the 2<sup>nd</sup> signal component from  $\sim 1.5$  to  $\sim 2.6$  s. The final part of  
 395 the signal from  $\sim 2.6$  to 3 s shows lower feature effects for all whale tracks.

396



397

398 FIG. 7. (Color online). SHAP feature importance corresponding to mean absolute Shapley  
 399 values for each class (i.e., whale tracks), computed using the training set for the SVM model, FSST4  
 400 transform, and high-quality signals.

401

#### 402 IV. DISCUSSION

403 Time-frequency transforms are important techniques for the study of non-stationary features  
404 of signals emitted by whales. The most commonly used techniques to analyze whale calls are the  
405 STFT and the spectrogram. However, when the signals are constituted by narrow-band time-  
406 frequency components, they are well-suited for analysis by the SST (e.g., *Daubechies et al., 2011*). The  
407 examples presented in [Figures 2](#) and [3](#) show that the FSST and FSST4 provide time-frequency  
408 representations of whale calls with better resolution than the STFT, which then aids the precise  
409 extraction of their characteristics and their interpretation. These SSTs are reversible which means  
410 that signal modes can be extracted and reconstructed. When signals are strongly frequency-  
411 modulated, which is often the case for whale sounds (e.g., for Humpback whales and Blue whales –  
412 *White & Todd, 2024*), high-order SST can be applied to the signal to better delineate the time-  
413 frequency features and avoid some mode mixing (*Pham & Meignen, 2017*). In this study, the whale  
414 calls consist mostly of two frequency components. As such, in principle, more time-frequency  
415 information could be included in the clustering, for example, more ridges or the full time-frequency  
416 representation. This would result in more feature parameters to be included in the clustering and  
417 would also require more training examples. If the full time-frequency representation of any of the  
418 SSTs is used, the thresholding parameter  $\zeta$  would need to be better adjusted to remove noise  
419 components.

420 Using the FSST4 instead of the STFT seems to slightly improve the SVM classification  
421 results. This relatively small improvement might be due to the simple characteristics of the Bryde's  
422 whale calls. On the contrary, using high-quality calls instead of lower-quality calls, and SVM instead  
423 of  $k$ -means, appears to substantially improve the classification results. The large difference in  
424 average cross-validation error and average balanced accuracy, that depends on the signal quality,  
425 suggests that signals recorded close to the whales are needed for their identification. In the present

426 case, this requires the signal to be observed by passive acoustic monitoring a few kilometers away  
427 from the calling whale. These conditions might, however, change depending on the method of  
428 analysis and classification, the recording conditions (i.e., sea bottom *vs.* sea surface, noise levels), and  
429 the type of calls (e.g., using temporal or frequency information, different frequency ranges).

430 For whales, caller identification is actually usually carried out using hydrophones deployed at  
431 the sea surface, close enough to the whales to enable them to be visually identified as well (e.g., *Gero*  
432 *et al.*, 2016; *Lamoni et al.*, 2023), and/or using dedicated instruments such as acoustic tags (*McCordic et*  
433 *al.*, 2016; *Oliveira et al.*, 2016; *Zeb et al.*, 2024). Contrary to other studies performing classification  
434 using frequency measures extracted from spectra or time-frequency representations (e.g., *McDonald et*  
435 *al.*, 2001; *McCordic et al.*, 2016), the time-frequency ridges included in the classification of the present  
436 study implicitly incorporate various spectral measures such as mean frequencies of the different  
437 signal components, component durations, maximum and minimum frequencies, and frequency  
438 modulations over time. Still, other measures could be combined in the classification such as other  
439 types of calls, call temporal patterns (e.g., codas rhythms for Sperm whales), call amplitudes or data  
440 from other instruments (e.g., from geophones in the present case). Feature selection or grouping,  
441 through dimensionality reduction for instance, could also be applied to the time-frequency ridge  
442 values to decrease the classification model complexity and improve its interpretability.

443 The large difference between the results of SVM compared with those of  $k$ -means could  
444 indicate that models using non-linear decision boundaries are more suitable to correctly classify  
445 high-dimensional representations of whale calls. A disadvantage of SVM relative to  $k$ -means is,  
446 however, the difficulty in setting the model hyperparameters (e.g., kernel function, upper bound  
447 parameter). The classification results are mainly limited by the number of training examples  
448 available, especially for whale tracks 1, 4 and 5. The training dataset could be expanded in different  
449 ways including collecting more calls, using data from other instruments, and using data generation

450 and augmentation strategies (e.g., *Zhu et al., 2020*). These strategies could correspond to using the  
451 same call several times with different noise types (real or synthetic), or making a synthetic call using  
452 a statistical description of the call properties (e.g., *Sochelean et al., 2015*). Another common limitation  
453 is the short observation window for each whale, due to both the temporary nature of most  
454 instrument deployments and the migratory behavior of whales, which often question the  
455 representativeness of the calls recorded.

456         While the present methodology and SVM model reach an average cross-validation error of  
457  $\sim 11\%$  and an average balanced accuracy of  $\sim 86 \pm 5\%$ , more observations and further study would  
458 be needed to test its generalization to a larger population of Bryde's whales. Whale calls from whale  
459 tracks 1, 2 and 3 were all recorded during the same period (January 29, 2015). These calls can be  
460 separated in three separate tracks based on their locations and amplitudes. Bryde's whales generally  
461 travel as individuals or in pairs, and seldom in larger groups. In the classification presented in this  
462 study, we assume that each track is generated by a unique vocal whale. The unsupervised t-SNE  
463 clustering seems to show that the observed calls of the different whale tracks are different enough to  
464 define individual clusters. However, the t-SNE clustering also indicates that some calls that are  
465 known to correspond to different whales (i.e., whale tracks 1 and 2 both recorded at the same time  
466 on January, 29, 2015, but localized at different positions) can exhibit similar call features using our  
467 processing. Hence, having more than one vocal whale per whale track is still an open question,  
468 especially for whale track 3 which has the largest number of calls. The close proximity between  
469 whale tracks 1 and 2, shown by their joining tracks ([Figure 1c](#)) and the t-SNE visualization, could  
470 also be indicative of a closer connection between these two whales.

471         Regarding whale tracks 4 and 5 recorded on different days (i.e., February 2 and January 26,  
472 2015, respectively), their changes in signal characteristics resulting in different clusters and  
473 classifications could be interpreted in various ways such as having five different vocal whales, or

474 returns of whales also recorded on other days which would involve temporal changes in signal  
475 characteristics due to spatiotemporal changes in underwater signal propagation. More generally, the  
476 observed call differences resulting in their successful classification could arise from morphological  
477 differences between whales, whales belonging to different populations, and spatiotemporal changes  
478 in environmental conditions impacting signal propagation (e.g., *Knight et al., 2024*). Finally, the SVM  
479 model using subtle differences in call time-frequency characteristics to distinguish between whales,  
480 does not necessarily imply that Bryde’s whales are using this information to identify themselves to  
481 conspecifics (e.g., *Gero et al., 2016*).

482

## 483 V. CONCLUSION

484 The identification of specific whale callers is important information for a range of  
485 applications such as studying whale movements and their change over time, and any external  
486 influence on their behavior. In the present study, we use highly similar low-frequency calls generated  
487 by five Bryde’s whales recorded by ocean-bottom hydrophones and compute time-frequency ridges  
488 using the 4<sup>th</sup>-order SST (FSST4) to extract the main frequency content of each call (i.e., time-  
489 frequency ridges). An SVM model is then trained using these time-frequency ridges to classify the  
490 whale calls. Using calls recorded less than 5 km away from the instruments, the average cross-  
491 validation error associated with the SVM model (Gaussian kernel) is  $\sim 11\%$  with an average balanced  
492 accuracy of  $\sim 86 \pm 5\%$ . Comparing these results with those using either STFT, lower-quality signals,  
493 or  $k$ -means clustering, shows that both the FSST4 and the SVM method improve the final results,  
494 with an increased error when using lower-quality signals and  $k$ -means clustering.

495 These classification results suggest that the short calls produced by Bryde’s whales, and the  
496 time-frequency characteristics embedded in the extracted ridges, contain caller identity cues. They

497 also seem to indicate that caller identity can be determined using ocean-bottom data, albeit using  
498 recordings less than a few kilometers away from the source. A larger number of training examples,  
499 coming from a larger number of well-identified whales and observed over a longer time period,  
500 would be needed to confirm these classification results. However, the methodology presented in this  
501 study does, nevertheless, show promising results and could be applied to other call types, improving  
502 the general understanding of whale vocalizations and ecology.

503

## 504 **ACKNOWLEDGMENTS**

505 This work was supported by the Natural Environment Research Council (NERC) via grant  
506 No. NE/1027010/1. We thank the people who participated in the data acquisition and processing  
507 used in this study (OSCAR research cruise JC114). The data were recorded by the UK's Natural  
508 Environment Research Council Ocean-Bottom Instrumentation Facility (*Minsbull et al., 2005*). The  
509 FSSTn toolbox was used to compute the FSST and FSST4 representations (*Pham & Meignen, 2017*).  
510 The generic mapping tools (GMT) were used for some figures (*Wessel & Smith, 1991*). For the  
511 purposes of open access, the authors have applied a Creative Commons Attribution (CC-BY) license  
512 to any author accepted manuscript arising from this work.

513

## 514 **AUTHOR DECLARATIONS**

### 515 **Conflict of Interest**

516 The authors declare no conflict of interest.

517

## 518 **DATA AVAILABILITY**

519 Continuous data of Ocean Bottom Seismometers from cruise JC114 are archived at the NERC's  
520 British Oceanographic Data Centre and are available following:  
521 [www.bodc.ac.uk/resources/inventories/cruise\\_inventory/report/15036/](http://www.bodc.ac.uk/resources/inventories/cruise_inventory/report/15036/).

522

## 523 REFERENCES

- 524 Arthur, D., and Vassilvitskii, S. (2007). "K-means++: the advantages of careful seeding. In  
525 Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms  
526 (SODA '07)", Society for Industrial and Applied Mathematics, USA, 1027-1035.
- 527 Bermant, P. C., Bronstein, M. M., Wood, R. J., Gero, S., and Gruber, D. F. (2019). "Deep machine  
528 learning techniques for the detection and classification of sperm whale bioacoustics",  
529 Scientific reports, **9**(1), 12588.
- 530 Boser, B. E., Guyon, I. M., and Vapnik, V. N. (1992). "A training algorithm for optimal margin  
531 classifiers", In Proceedings of the fifth annual workshop on Computational learning theory,  
532 144-152.
- 533 Charrier, I., Jouventin, P., Mathevon, N., and Aubin, T. (2001). "Individual identity coding depends  
534 on call type in the South Polar skua *Catharacta maccormicki*", Polar Biology, **24**, 378-382.
- 535 Combrisson, E., and Jerbi, K. (2015). "Exceeding chance level by chance: The caveat of theoretical  
536 chance levels in brain signal classification and statistical assessment of decoding accuracy",  
537 Journal of neuroscience methods, **250**, 126-136.
- 538 Cortes, C., and Vapnik, V. (1995). "Support-vector networks", Machine learning, **20**, 273-297.
- 539 Daubechies, I., Lu, J., and Wu, H. T. (2011). "Synchrosqueezed wavelet transforms: An empirical  
540 mode decomposition-like tool", Applied and computational harmonic analysis, **30**(2), 243-  
541 261.

542 Dunlop, R. A., Noad, M. J., Cato, D. H., and Stokes, D. (2007). “The social vocalization repertoire  
543 of east Australian migrating humpback whales (*Megaptera novaeangliae*)”, *The Journal of*  
544 *the Acoustical Society of America*, **122**(5), 2893-2905.

545 Gero, S., Whitehead, H., and Rendell, L. (2016). “Individual, unit and vocal clan level identity cues in  
546 sperm whale codas”, *Royal Society Open Science*, **3**(1), 150372.

547 Halkias, X. C., Paris, S., and Glotin, H. (2013). “Classification of mysticete sounds using machine  
548 learning techniques”, *The Journal of the Acoustical Society of America*, **134**(5), 3496-3505.

549 Hobbs, R., and Peirce, C. (2015). RRS James Cook JC114 Cruise Report. Online Report.  
550 [https://www.bodc.ac.uk/resources/inventories/cruise\\_inventory/reports/jc114.pdf](https://www.bodc.ac.uk/resources/inventories/cruise_inventory/reports/jc114.pdf).

551 Huang, Z. L., Zhang, J., Zhao, T. H., and Sun, Y. (2015). “Synchrosqueezing S-transform and its  
552 application in seismic spectral decomposition”, *IEEE Transactions on Geoscience and*  
553 *Remote Sensing*, **54**(2), 817-825.

554 Ibrahim, A. K., Zhuang, H., Chérubin, L. M., Erdol, N., O’Corry-Crowe, G., and Ali, A. M. (2021).  
555 “A multimodel deep learning algorithm to detect North Atlantic right whale up-calls”, *The*  
556 *Journal of the Acoustical Society of America*, **150**(2), 1264-1272.

557 Janik, V. M. (2009). “Acoustic communication in delphinids”, *Advances in the Study of Behavior*,  
558 **40**, 123-157.

559 Janik, V. M., and Sayigh, L. S. (2013). “Communication in bottlenose dolphins: 50 years of signature  
560 whistle research”, *Journal of Comparative Physiology A*, **199**, 479-489.

561 Kather, V., Seipel, F., Berges, B., Davis, G., Gibson, C., Harvey, M., ... and Risch, D. (2024).  
562 “Development of a machine learning detector for North Atlantic humpback whale song”,  
563 *The Journal of the Acoustical Society of America*, **155**(3), 2050-2064.

564 Kecman, V. (2005). “Support vector machines—an introduction”, in *Support vector machines: theory and*  
565 *applications* (pp. 1-47). Berlin, Heidelberg: Springer Berlin Heidelberg.

566 Knight, E., Rhinehart, T., de Zwaan, D. R., Weldy, M. J., Cartwright, M., Hawley, S. H., ... and  
567 Kitzes, J. (2024). "Individual identification in acoustic recordings", *Trends in Ecology &*  
568 *Evolution*, **39**(10), 947-960.

569 Lamoni, L., Garland, E. C., Allen, J. A., Coxon, J., Noad, M. J., and Rendell, L. (2023). "Variability  
570 in humpback whale songs reveals how individuals can be distinctive when sharing a complex  
571 vocal display", *The Journal of the Acoustical Society of America*, **153**(4), 2238-2238.

572 Lomax, A., Zollo, A., Capuano, P., and Virieux, J. (2001). "Precise, absolute earthquake location  
573 under Somma–Vesuvius volcano using a new three-dimensional velocity model",  
574 *Geophysical Journal International*, **146**(2), 313-331.

575 Lowell, R. P., Zhang, L., Maqueda, M. A. M., Banyte, D., Tong, V. C. H., Johnston, R. E. R., ... and  
576 Kolandaivelu, K. (2020). "Magma-hydrothermal interactions at the Costa Rica Rift from data  
577 collected in 1994 and 2015", *Earth and Planetary Science Letters*, **531**, 115991.

578 Lundberg, S. M., and Su-In, L. (2017). "A unified approach to interpreting model predictions",  
579 *Advances in neural information processing systems*, **30**, 4765-4774.

580 McCordic, J. A., Root-Gutteridge, H., Cusano, D. A., Denes, S. L., and Parks, S. E. (2016). "Calls of  
581 North Atlantic right whales *Eubalaena glacialis* contain information on individual identity  
582 and age class", *Endangered Species Research*, **30**, 157-169.

583 McDonald, M. A., Calambokidis, J., Teranishi, A. M., and Hildebrand, J. A. (2001). "The acoustic  
584 calls of blue whales off California with gender data", *The Journal of the Acoustical Society*  
585 *of America*, **109**(4), 1728-1735.

586 Mazhar, S., Ura, T., and Bahl, R. (2007). "Vocalization based individual classification of humpback  
587 whales using support vector machine", in *OCEANS 2007* (pp. 1-9). IEEE.

588 Meignen, S., Pham, D. H., and McLaughlin, S. (2017). “On demodulation, ridge detection, and  
589 synchrosqueezing for multicomponent signals”, *IEEE Transactions on Signal Processing*,  
590 **65**(8), 2093-2103.

591 Mellinger, D. K., and Clark, C. W. (2000). “Recognizing transient low-frequency whale sounds by  
592 spectrogram correlation”, *The Journal of the Acoustical Society of America*, **107**(6), 3518-  
593 3529.

594 Minshull, T. A., Sinha, M. C. and Peirce, C. (2005). “Multi-disciplinary subseabed geophysical  
595 imaging - a new pool of 28 seafloor instruments in use by the United Kingdom Ocean  
596 Bottom Instrument Consortium”, *Sea Technology*, **46**, 27-31.

597 Molnar, C. (2022). *Interpretable Machine Learning: A Guide for Making Black Box Models*  
598 *Explainable* (2nd ed.).

599 Oberlin, T., Meignen, S., and Perrier, V. (2015). “Second-order synchrosqueezing transform or  
600 invertible reassignment? Towards ideal time-frequency representations”, *IEEE Transactions*  
601 *on Signal Processing*, **63**(5), 1335-1344.

602 Ojala, M., and Garriga, G. C. (2010). “Permutation tests for studying classifier performance”, *Journal*  
603 *of machine learning research*, **11**(6), 1833-1863.

604 Oleson, E. M., Barlow, J., Gordon, J., Rankin, S., and Hildebrand, J. A. (2003). “Low frequency calls  
605 of Bryde’s whales”, *Mar. Mammal Sci.*, **19**(2), 407-419.

606 Oliveira, C., Wahlberg, M., Silva, M. A., Johnson, M., Antunes, R., Wisniewska, D. M., ... and  
607 Madsen, P. T. (2016). “Sperm whale codas may encode individuality as well as clan identity”,  
608 *The Journal of the Acoustical Society of America*, **139**(5), 2860-2869.

609 Palacios, D.M., Herrera, J.C., Gerrodette, T., Garcia, C., Soler, G.A., Avila, I.C., Bessudo, S.,  
610 Hernandez, E., Trujillo, F., Forez-Gonzalez, L., and Kerr, I. (2012). “Cetacean distribution

611 and relative abundance in Colombia's Pacific EEZ from survey cruises and platforms of  
612 opportunity", *J. Cetacean Res. Manage.* **12**(1), 45-60.

613 Pham, D. H., and Meignen, S. (2017). "High-order synchrosqueezing transform for multicomponent  
614 signals analysis—With an application to gravitational-wave signal", *IEEE Transactions on*  
615 *Signal Processing*, **65**(12), 3168-3178.

616 Rasmussen, J. H., and Širović, A. (2021). "Automatic detection and classification of baleen whale  
617 social calls using convolutional neural networks", *The Journal of the Acoustical Society of*  
618 *America*, **149**(5), 3635-3644.

619 Rendell, L. E., and Whitehead, H. (2003). "Vocal clans in sperm whales (*Physeter macrocephalus*)",  
620 *Proceedings of the Royal Society of London. Series B: Biological Sciences*, **270**(1512), 225-  
621 231.

622 Robinson, A. H., Zhang, L., Hobbs, R. W., Peirce, C., and Tong, V. C. H. (2020). "Magmatic and  
623 tectonic segmentation of the intermediate-spreading Costa Rica Rift—a fine balance  
624 between magma supply rate, faulting and hydrothermal circulation", *Geophysical Journal*  
625 *International*, **222**(1), 132-152.

626 Salmi, R., Hammerschmidt, K., and Doran-Sheehy, D. M. (2014). "Individual distinctiveness in call  
627 types of wild western female gorillas", *Plos One*, **9**(7), e101940.

628 Stanković, L. (2001). "A measure of some time–frequency distributions concentration", *Signal*  
629 *Processing*, **81**(3), 621-631.

630 Socheleau, F. X., Leroy, E., Carvalho Pecci, A., Samaran, F., Bonnel, J., and Royer, J. Y. (2015).  
631 "Automated detection of Antarctic blue whale calls", *The Journal of the Acoustical Society*  
632 *of America*, **138**(5), 3105-3117.

633 Tary, J. B., Herrera, R. H., Han, J., and van der Baan, M. (2014). "Spectral estimation—What is new?  
634 What is next?". *Reviews of Geophysics*, **52**(4), 723-749.

635 Tary, J. B., Hobbs, R. W., Peirce, C., Lesmes Lesmes, C., and Funnell, M. J. (2021). “Local rift and  
636 intraplate seismicity reveal shallow crustal fluid-related activity and sub-crustal faulting”,  
637 Earth and Planetary Science Letters, **562**, 116857.

638 Tary, J. B., Peirce, C., Hobbs, R. W., Bonilla Walker, F, De La Hoz, C., Bird, A., and Vargas, C. A.  
639 (2024). “Application of a seismic network to baleen whale call detection and localization in  
640 the Panama basin – a Bryde’s whale example”, Journal of the Acoustical Society of America,  
641 **155**(3), 2075-2086.

642 Thakur, G., and Wu, H. T. (2011). “Synchrosqueezing-based recovery of instantaneous frequency  
643 from nonuniform samples”, SIAM Journal on Mathematical Analysis, **43**(5), 2078-2095.

644 Van der Maaten, L., and Hinton, G. (2008). “Visualizing data using t-SNE”, Journal of machine  
645 learning research, **9**(11).

646 Vapnik, V. (1995). *The Nature of Statistical Learning Theory*, Springer-Verlag, New York.

647 Wade, P.R., and Gerrodette, T. (1993). “Estimates of cetacean abundance and distribution in the  
648 eastern tropical Pacific”, Report of the International Whaling Commission **43**, 477-493.

649 Waldhauser, F., and Ellsworth, W. L. (2000). “A double-difference earthquake location algorithm:  
650 Method and application to the northern Hayward fault, California”, Bulletin of the  
651 seismological society of America, **90**(6), 1353-1368.

652 Wessel, P., and Smith, W. H. (1991). “Free software helps map and display data”, Eos, Transactions  
653 American Geophysical Union, **72**(41), 441-446.

654 White, P. R., and Todd, V. L. (2024). “Baleen Whale Vocalizations”, in Noisy Oceans: Monitoring  
655 Seismic and Acoustic Signals in the Marine Environment, 183-207.

656 Wilson, D. J., Robinson, A. H., Hobbs, R. W., Peirce, C., and Funnell, M. J. (2019). “Does  
657 intermediate spreading-rate oceanic crust result from episodic transition between magmatic

658 and magma-dominated, faulting-enhanced spreading?—The Costa Rica Rift example”,  
659 Geophysical Journal International, **218**(3), 1617-1641.

660 Zeh, J. M., Perez-Marrufo, V., Adcock, D. L., Jensen, F. H., Knapp, K. J., Robbins, J., ... and Parks,  
661 S. E. (2024). “Caller identification and characterization of individual humpback whale  
662 acoustic behaviour”, Royal Society Open Science, **11**(3), 231608.

663 Zhong, M., Torterotot, M., Branch, T. A., Stafford, K. M., Royer, J. Y., Dodhia, R., and Lavista  
664 Ferres, J. (2021). “Detecting, classifying, and counting blue whale calls with Siamese neural  
665 networks”, The Journal of the Acoustical Society of America, **149**(5), 3086-3094.

666 Zhu, W., Mousavi, S. M., and Beroza, G. C. (2020). “Seismic signal augmentation to improve  
667 generalization of deep neural networks”, in *Advances in geophysics* **61**, 151-177. Elsevier.  
668



**Citation on deposit:** Tary, J.-B., Peirce, C., & Hobbs, R. (in press). Classification of Bryde's whale individuals using high-resolution time-frequency transform and support vector machines. *The Journal of the Acoustical Society of America*

**For final citation and metadata, visit Durham Research Online URL:**

<https://durham-repository.worktribe.com/output/3680785>

**Copyright statement:** This accepted manuscript is licensed under the Creative Commons Attribution 4.0 licence.

<https://creativecommons.org/licenses/by/4.0/>