

1    **Automated tracking of dolphin whistles using Gaussian Mixture Probability**  
2                    **Hypothesis Density (GM-PHD) filters**

3                    Pina Gruden<sup>a)</sup> and Paul R. White

4                    Institute of Sound and Vibration Research

5                    University of Southampton

6                    Highfield, Hants, SO17 1BJ, UK

---

<sup>a)</sup>e-mail: [pg3g12@soton.ac.uk](mailto:pg3g12@soton.ac.uk)

## Abstract

8        This work considers automated Multi Target Tracking (MTT) of odontocete whistle  
9 contours. An adaptation of Gaussian Mixture Probability Hypothesis Density (GM-PHD)  
10 filter is described and applied to the acoustic recordings from six odontocete species. From  
11 the raw data, spectral peaks are first identified and then GM-PHD filter is used to simul-  
12 taneously track the whistles' frequency contours. Overall over 9000 whistles are tracked  
13 with a precision of 85% and recall of 71.8%. The proposed filter is shown to track whis-  
14 tles precisely (with mean deviation of 104 Hz, about one frequency bin, from the annotated  
15 whistle path) and 80% coverage. The filter is computationally efficient, suitable for real-time  
16 implementation, and is widely applicable to different odontocete species.

## I. INTRODUCTION

The detection of marine mammal vocalizations plays an important role in passive acoustic monitoring. The objectives of such studies include species recognition<sup>24,29,9</sup>, species presence and abundance estimation<sup>21</sup>, studying species behaviour<sup>26</sup>, mitigation during industrial activities<sup>35</sup>. Odontocetes (toothed whales) produce a rich variety of high-frequency vocalizations, which can be grouped into three broad categories: whistles, echolocation clicks and burst pulses<sup>2</sup>, all of which have most of their energy above 2 kHz<sup>31</sup>. This work focuses on whistles, which are highly variable, narrowband, frequency modulated, tonal sounds with fundamental frequencies generally between 2 and 30 kHz and are typically used in a social context<sup>14</sup>. Not all odontocete species whistle, but majority of delphinid species do.

Methods used for detection and frequency estimation of odontocete whistles vary from semi-automated methods *e.g.*,<sup>14,24</sup> to fully automated methods *e.g.*,<sup>8,36,11,28,12,9,13</sup>. Most methods are based on spectrogram techniques, although alternative approaches also exist *e.g.*,<sup>12,11</sup>.

Prior to applying a detection algorithm to the signal, some pre-processing of data is typically carried out in order to reduce background noise and interfering signals *e.g.*,<sup>8,36,12,22,28,9</sup>. After the noise removal, spectrogram-based algorithms usually identify strongest spectral peaks *e.g.*,<sup>10,15,28,22,9</sup> or apply image-processing techniques to define the pixels<sup>20</sup> or ridges that represent whistles<sup>13</sup>. The identified peaks are then connected into a continuous whistle

36 contour using different approaches, such as particle filtering<sup>36,28</sup>, Kalman filtering<sup>20</sup>, combi-  
37 nation of polynomial fitting and Kalman filtering<sup>13</sup>; hypothesis tracking with some gating  
38 rules<sup>10,22</sup>; phase tracking<sup>11,12</sup>.

39       The automated methods for whistle contour detection are commonly based on the  
40 algorithms that allow for single target tracking. In this work, an alternative approach is taken  
41 in which the detection and tracking of frequency content of delphinid whistles is considered  
42 as a multi-target tracking (MTT) problem, where whistles are targets that overlap, their  
43 numbers are unknown and vary with time and there are interfering signals present. An MTT  
44 algorithm called Gaussian Mixture Probability Hypothesis Density (GM-PHD) filter<sup>16,32</sup>,  
45 which has been previously used in sonar applications<sup>6</sup>, was adapted here for application of  
46 dolphin whistle contour tracking. The paper is organized as follows. In Section II some  
47 background is given on target tracking and PHD filters. Section III introduces formulation  
48 of the GM-PHD filter for dolphin whistle tracking and derivation of models and parameters  
49 for this particular problem. The performance of the proposed GM-PHD filter is tested on  
50 the acoustic recordings of dolphin whistles, which have been hand-annotated and results  
51 are given in Section IV. Discussion and conclusions may be found in Sections V and VI  
52 respectively. Appendix summarizes the most frequently used symbols and their meanings.

## 53       **II. BACKGROUND**

## A. Target Tracking

Target tracking is a process of estimating a target's state as it evolves in time, from a sequence of noisy measurements. A target is broadly defined as the entity to be tracked and the state vector,  $\mathbf{x}_k$ , contains the information about the properties of the target at time  $k$ . The only available information about the targets is given by the measurement vector,  $\mathbf{z}_k$ , which also typically contains noise. In the case of whistle frequency contour tracking, each whistle represents a target. The target state vector consists of frequency and chirp (rate of change of frequency) information and the measurement vector consists of frequency peaks. Measurements may also be contaminated by the detection of false targets (clutter) and points where there has been a failure to detect a target.

In order to perform target tracking at least two models are required; first a model describing the evolution of the state with time, called the system (or dynamic) model

$$\mathbf{x}_k = \Phi_k(\mathbf{x}_{k-1}, \mathbf{n}_{k-1}) \quad (1)$$

where  $\Phi_k$  is a system function that describes the evolution of the state vector and  $\mathbf{n}_{k-1}$  is a system noise process and is a vector of random variables specifying the random component of the parameter evolution<sup>1,36</sup>. From the system model one can define a state transition density  $f_{k|k-1}(\mathbf{x}_k|\mathbf{x}_{k-1})$ , which characterizes the transition of the state from time  $k-1$  to time  $k$ . The second model required is a model relating the noisy measurement to the state,

71 called the measurement (or observation) model

$$\mathbf{z}_k = \psi_k(\mathbf{x}_k, \boldsymbol{\eta}_k) \quad (2)$$

72 where  $\psi_k$  is a function that defines the measurement process and  $\boldsymbol{\eta}_k$  is the measurement noise  
73 process<sup>1,36</sup>. From the measurement model one can obtain a likelihood function  $g_k(\mathbf{z}_k|\mathbf{x}_k)$ ,  
74 that describes the likelihood that a measurement  $\mathbf{z}_k$  was generated by the target  $\mathbf{x}_k$ . These  
75 models are collectively known as a state-space model.

76 Target tracking is typically achieved with the use of a recursive Bayesian filter where  
77 one attempts to construct the posterior probability density function (pdf) of the state,  
78  $p_k(\mathbf{x}_k|\mathbf{z}_{1:k})$ , based on the set of measurements  $\mathbf{z}_{1:k}$  up to time  $k$ <sup>1</sup>. Such a filter involves  
79 a two stage process; prediction and update, where the system model is used to predict the  
80 state pdf and the measurements are used to refine that prediction<sup>1</sup>. This is implemented  
81 in a recursive manner and at each time step an estimate of the state is obtained from the  
82 posterior pdf.

83 In the case of single target tracking, it is assumed that only one target is present and  
84 that all the observations are generated by that target. If the system and measurement models  
85 are linear and the noise processes are Gaussian, then optimal target tracking is achieved with  
86 the Kalman filter<sup>5</sup>, which in this case represents the optimal solution to Bayesian recursion.  
87 If the models are non-linear and/or the noise is non-Gaussian, particle filters can be used to

88 perform single target tracking<sup>36</sup>.

89 In the majority of real-world applications there are multiple targets present at any given  
 90 time, the number of which will change through time as targets appear (*i.e.* target birth)  
 91 and disappear (*i.e.* target death). At each time  $k$  there are  $n_k$  target states  $\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,n_k}$   
 92 and  $m_k$  measurements  $\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,m_k}$ . The states of the targets and the observations can be  
 93 modelled using the concept of a random finite set. A random finite set is an object in which  
 94 the elements have random values, as in any multivariate random process, but in addition  
 95 to which the number of elements in the set is also random<sup>27</sup>. The random set of states  
 96 (multi-target state),  $\mathbf{X}_k$ , and the random set of measurements (multi-target measurement),  
 97  $\mathbf{Z}_k$ , are represented as follows:

$$\mathbf{X}_k = \{\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,n_k}\} \in \mathcal{F}(\mathcal{X}) \quad (3)$$

$$\mathbf{Z}_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,m_k}\} \in \mathcal{F}(\mathcal{Z}) \quad (4)$$

99 where  $\mathcal{F}(\mathcal{X})$  and  $\mathcal{F}(\mathcal{Z})$  are the finite subsets of the state and observation spaces  $\mathcal{X}$  and  $\mathcal{Z}$ ,  
 100 respectively.

101 In this case the use of multi target tracking (MTT) techniques is required and the  
 102 objective is to jointly estimate the number of targets and their states from the noisy mea-  
 103 surements<sup>32</sup>. Traditional approaches to MTT are based on data association techniques and  
 104 involve explicit associations between measurements and targets that are achieved with the use

105 of single target tracking techniques. Examples of traditional MTT include nearest neighbor  
 106 (NN), joint probabilistic data association (JPDA) and multiple hypothesis tracking (MHT)<sup>4</sup>.

107 However, the uncertainty in the evolution of the multi-target state and the origin of the  
 108 multi-target measurement is naturally modelled by random finite sets<sup>27</sup> and therefore data  
 109 association-free techniques, based on Mahler's finite set statistics (FISST) framework (an  
 110 overview is provided in<sup>19</sup>), have been increasingly used in the last decade for the Bayesian  
 111 multi-target filtering problems. A multi-target Bayesian filter determines at each time step  $k$   
 112 the posterior probability density of multitarget-state  $p_k(\mathbf{X}_k|\mathbf{Z}_{1:k})$ <sup>27</sup>. The high dimensionality  
 113 of the Bayes multi-target filter makes the recursion intractable in practice, a problem which  
 114 is overcome using the Probability Hypothesis Density (PHD) filter<sup>16,17</sup>.

## 115 B. Probability Hypothesis Density (PHD) filter

116 The PHD filter approximates the multi-target Bayes recursion by propagating the first-  
 117 order statistical moment  $v_k(\mathbf{x}|\mathbf{Z}_{1:k})$  of the multi-target posterior  $p_k(\mathbf{X}_k|\mathbf{Z}_{1:k})$ , known as the  
 118 intensity function or the PHD<sup>17,32,25,27</sup>. The PHD is a function whose peaks identify the  
 119 likely positions of the targets. By integrating the PHD on any region of the state space  
 120 one obtains the expected number of targets in that region. It should be noted that PHD  
 121 is a density function but is not a pdf, since its integral over the space of its variable is not  
 122 unity<sup>19</sup>. A target with state  $\mathbf{x}$  is more likely to be present in the region when the PHD  
 123 (intensity function) is large than when it is small, which allows one to obtain state estimates



124 of the targets based on peaks in the PHD.

125 The PHD filter comprises both prediction and update steps. In the prediction step,  
 126 the PHD filter incorporates the motion of individual targets and accounts for disappearance  
 127 of existing targets (by incorporating the probability of target's survival). In addition it  
 128 incorporates the appearance of completely new targets. Hence, the predicted intensity func-  
 129 tion,  $v_{k|k-1}(\cdot)$ , consists of the newborn targets (introduced by the birth intensity function)  
 130 and the existing targets (targets surviving from the previous time step that are represented  
 131 by the posterior intensity function from the previous time step  $v_{k-1}(\cdot)$ ). The abbreviation  
 132  $v_k(\mathbf{x}|\mathbf{Z}_{1:k}) \stackrel{abbr}{=} v_k(\mathbf{x}_k)$  is used and the prediction step can be expressed as<sup>32,25,27</sup>

$$v_{k|k-1}(\mathbf{x}_k) = \gamma_k(\mathbf{x}_k) + \langle p_{S,k}(\mathbf{x}_{k-1})v_{k-1}(\mathbf{x}_{k-1}), f_{k|k-1}(\mathbf{x}_k|\mathbf{x}_{k-1}) \rangle \quad (5)$$

133 where  $\gamma_k(\mathbf{x}_k)$  denotes the PHD of target births between time  $k-1$  and  $k$ ;  $p_{S,k}(\mathbf{x}_{k-1})$  denotes  
 134 the probability of survival, that is probability that a target with state  $\mathbf{x}$  at time  $k-1$  will  
 135 survive until time  $k$ ;  $f_{k|k-1}(\mathbf{x}_k|\mathbf{x}_{k-1})$  denotes single-target state transition density from time  
 136  $k-1$  to  $k$  and  $\langle g, f \rangle = \int f(x)g(x)dx$ . Note that spawning terms, that define how one target  
 137 can become resolved into more than one target, have been omitted from the above equation.  
 138 This is because rarely, if ever, does one observe a dolphin whistle contour which splits into  
 139 two distinct contours.

140 In the update step, the PHD filter incorporates the probability that any given target

141 was not detected (by incorporating the probability of target detection) and updates the pre-  
 142 dicted intensity with a set of measurements by also taking into the account the measurement  
 143 likelihood function and false alarms (clutter). The posterior intensity function  $v_k(\cdot)$  at time  
 144 step  $k$  is given by

$$v_k(\mathbf{x}_k) = [1 - p_{D,k}(\mathbf{x}_k)]v_{k|k-1}(\mathbf{x}_k) + \sum_{\mathbf{z} \in \mathbf{Z}_k} \frac{p_{D,k}(\mathbf{x}_k)g_k(\mathbf{z}|\mathbf{x}_k)v_{k|k-1}(\mathbf{x}_k)}{\kappa_k(\mathbf{z}) + \langle p_{D,k}(\mathbf{x}_k)g_k(\mathbf{z}|\mathbf{x}_k), v_{k|k-1}(\mathbf{x}_k) \rangle} \quad (6)$$

145 where  $p_{D,k}(\mathbf{x}_k)$  denotes the probability of detection, that is the probability that observation  
 146 will be collected at time  $k$  from a target with state  $\mathbf{x}_k$ ,  $\mathbf{Z}_k$  denotes the multi-target measure-  
 147 ment at time  $k$ ,  $\kappa_k(\mathbf{z})$  denotes denotes the PHD of clutter at time  $k$  and  $g_k(\mathbf{z}|\mathbf{x}_k)$  denotes  
 148 the single-target measurement likelihood function at time  $k$ .

149 The computational load of the PHD filter can grow significantly if target births can  
 150 occur uniformly in the state space. One approach to mitigate this is to adapt the birth  
 151 intensity according to the measurements<sup>27</sup>, which results in the prediction and update steps  
 152 being preformed separately for newborn and existing targets. A label  $\beta$  is introduced to  
 153 distinguish between the two types of targets;  $\beta = 0$  refers to existing targets,  $\beta = 1$  refers  
 154 to newborn targets. The prediction stage becomes<sup>27</sup>

$$\begin{aligned}
v_{k|k-1}(\mathbf{x}_k, \beta) &= \gamma_k(\mathbf{x}_k) & \beta &= 1 \\
&= \langle p_{S,k}(\mathbf{x}_{k-1})v_{k-1}(\mathbf{x}_{k-1}), f_{k|k-1}(\mathbf{x}_k|\mathbf{x}_{k-1}) \rangle & \beta &= 0
\end{aligned} \tag{7}$$

155 where  $v_{k-1}(\cdot)$  represents posterior intensity function from the previous time step and consists  
156 of posterior intensity functions of existing and newborn targets from the previous time step  
157 ( $v_{k-1}(\cdot, 0) + v_{k-1}(\cdot, 1)$ ).

158 The update stage of the filter for existing targets ( $\beta = 0$ ) can be expressed as<sup>27</sup>

$$\begin{aligned}
v_k(\mathbf{x}_k, 0) &= [1 - p_{D,k}(\mathbf{x}_k)]v_{k|k-1}(\mathbf{x}_k, 0) \\
&+ \sum_{\mathbf{z} \in Z_k} \frac{p_{D,k}(\mathbf{x}_k)g_k(\mathbf{z}|\mathbf{x}_k)v_{k|k-1}(\mathbf{x}_k, 0)}{\mathcal{L}(\mathbf{z})}
\end{aligned} \tag{8}$$

159 and for newborn targets ( $\beta = 1$ )

$$v_k(\mathbf{x}_k, 1) = \sum_{\mathbf{z} \in Z_k} \frac{g_k(\mathbf{z}|\mathbf{x}_k)\gamma_k(\mathbf{x}_k)}{\mathcal{L}(\mathbf{z})} \tag{9}$$

160 where

$$\mathcal{L}(\mathbf{z}) = \kappa_k(\mathbf{z}) + \langle g_k(\mathbf{z}|\mathbf{x}_k), \gamma_k \rangle + \langle p_{D,k}(\mathbf{x}_k)g_k(\mathbf{z}|\mathbf{x}_k), v_{k|k-1}(\mathbf{x}_k, 0) \rangle \tag{10}$$

161 Note that since newborn targets are created from the measurements, the newborn  
162 targets are always detected, *i.e.*  $p_D(\mathbf{x}, 1) = 1$ <sup>27</sup>.

163 It can be seen from the above equations that in addition to the system and measure-  
 164 ment models (from which the  $f_{k|k-1}(\mathbf{x}_k|\mathbf{x}_{k-1})$  and  $g_k(\mathbf{z}|\mathbf{x}_k)$  are obtained respectively), the  
 165 PHD filter requires definition of additional models and parameters. Specifically, the target's  
 166 survival ( $p_{S,k}(\mathbf{x}_{k-1})$ ) and detection ( $p_{D,k}(\mathbf{x}_k)$ ) probabilities and clutter ( $\kappa_k(\mathbf{z})$ ) and target  
 167 birth ( $\gamma_k(\mathbf{x}_k)$ ) models. The formulation of these is described in the Section III.B.2.

168 The above equations still involve integrals that typically have no closed form solution  
 169 and therefore the PHD filter needs to be approximated<sup>32,25</sup>. Practical implementations of  
 170 PHD filters include Gaussian Mixture PHD (GM-PHD)<sup>32</sup> and Sequential Monte Carlo PHD  
 171 (SMC-PHD)<sup>33</sup> filters. In this work the GM-PHD approach was chosen since it tends to be  
 172 faster and more straightforward than the SMC-PHD approach<sup>18</sup>. The GM-PHD filter and  
 173 its application to a specific problem of dolphin whistle tracking is presented in the next  
 174 section.

### 175 III. METHODOLOGY

#### 176 A. Data, pre-processing steps and obtaining the measurements

177 The data set used in this study was obtained from the 5th Workshop of Detection,  
 178 Classification, Localization and Density Estimation (DCLDE) conference 2011 (available at  
 179 MobySound archive, <http://www.mobysound.org>). This dataset contained raw data and  
 180 analyst-annotated files for six species: long-beaked common dolphin (*Delphinus capensis*),

181 short-beaked common dolphin (*Delphinus delphis*), melon-headed whales (*Peponocephala*  
182 *electra*), spinner dolphin (*Stenella longirostris*), Atlantic spotted dolphin (*Stenella frontalis*)  
183 and bottlenose dolphin (*Tursiops truncatus*). The recordings contained in this dataset were  
184 single-species recordings that were confirmed by trained visual observers. Study areas, data  
185 collection protocols and procedure for hand-annotation of the data are summarized in Roch  
186 *et al.*<sup>28</sup>, Baumann-Pickering *et al.*<sup>3</sup> and Soldevilla *et al.*<sup>30</sup>. The raw data was used for the  
187 GM-PHD filter to track the whistles from and hand-annotations were used to evaluate the  
188 performance of the filter. In addition, a small part of raw data was set aside to be used  
189 as training data for certain parameters of the GM-PHD filter. For this purpose three files  
190 were randomly selected from the annotated dataset and a 1 minute section of each of those  
191 files was taken as the training data. These training files corresponded to three species, *D.*  
192 *capensis*, *D.delphis* and *S.frontalis*, and were obtained using different recording equipment.  
193 This training data was subsequently not used in the performance evaluation.

194 For ease of implementation, where necessary, the data was re-sampled to 192 kHz (before  
195 re-sampling 2.5% of the files were sampled at 300 kHz, 12.5% at 480 kHz and 85% at 192 kHz).  
196 After re-sampling, pre-processing was applied to the data in order to reduce the background  
197 noise and interfering signals. A pre-processing scheme was adapted from Gillespie *et al.*<sup>9</sup>  
198 and was applied with a sliding window that was 2048 points long and had 50% overlap,  
199 resulting in 93.8 Hz spacing between frequency bins. Within each window the following

200 steps were performed as described in Gillespie *et al.*<sup>9</sup>: first echolocation clicks were removed  
201 by applying a weighting function; then spectrogram was computed on a decibel scale, using  
202 2048 point Hanning window, and spectral peaks were enhanced by applying normalization  
203 across frequency based on a 61 point median filter; after that the normalization across time  
204 using exponential moving average (with the weighting constant of 0.02) was performed in  
205 order to remove persistent tones from the spectrogram.

206 In each window, after the noise was removed, spectral peaks were determined by iden-  
207 tifying all frequencies whose normalized magnitude exceeded 8 dB. Only frequency bins  
208 between 2 and 50 kHz were searched for peaks, since most dolphin whistles will lie within  
209 this range and to be consistent with the hand annotations which were also applied to whistle  
210 harmonics. The identified spectral peaks represent the measurement set from which the  
211 whistle contours were tracked using the Gaussian Mixture PHD (GM-PHD) filter.

212 Measurement sets containing spectral peak measurements and a list of all files used  
213 in this study, as well as Matlab implementation of the method for obtaining spectral peak  
214 measurements was released to the MobySound archive.

## 215 **B. Whistle contour tracking with Gaussian Mixture PHD (GM-PHD) filters**

216

217 The GM-PHD filter algorithm<sup>32</sup> was implemented and used to track frequency contours  
218 of whistles from the identified spectral peaks. In this approximation to the PHD filter, the

219 posterior intensity function  $v_k(\mathbf{x}_k)$  is represented by a sum of weighted Gaussian components  
220 whose weights, means and covariances are propagated in time<sup>25</sup>. This strategy is analogous  
221 to Kalman filter<sup>5</sup> for single target tracking, which propagates the first moment (the mean)  
222 of the single-target state<sup>32</sup>. So that each whistle at time  $k$  is represented by a Gaussian  
223 component and is therefore characterized by a mean (consisting of frequency and chirp), a  
224 weight and a covariance. The means and covariances of the existing and newborn whistles  
225 are predicted using the Kalman filter prediction equations and updated with the received  
226 measurements (spectral peaks) also using the Kalman equations. The weights of the whistles  
227 are predicted and updated using the PHD equations and they can be thought of as a measure  
228 of the likelihood of presence of a component. Detailed description of the GM-PHD filter is  
229 given next.

230       The whistle estimates generated by the GM-PHD filter do not inherently contain iden-  
231 tity. In order to assign a particular state to a specific whistle, tracking of Gaussian compo-  
232 nents needs to be carried out. Tracking is achieved by labelling each individual Gaussian  
233 component with a unique tag and the likelihood of each track is then given by the weight of  
234 each component *e.g.*,<sup>7,25,34</sup>.

235       This section is organized as follows. First the GM-PHD algorithm is outlined, then  
236 filter's models and parameters are defined, followed by a description of the performance  
237 evaluation.

238

**1. The GM-PHD algorithm**

239

240

241

242

243

244

The GM-PHD filter approximates the intensity functions (PHDs) with Gaussian mixtures. It should be noted that these do not share the properties of GM approximations to pdfs in terms of weights summing to 1. Here the sum of weights reflects the number of whistles present at each time step. The GM-PHD filter makes the following assumptions. It is assumed that each whistle follows a linear Gaussian dynamical model and that measurements follow a linear model<sup>32</sup>. That is, (1) and (2) can be written as

$$\mathbf{x}_k = F_{k-1}\mathbf{x}_{k-1} + \mathbf{n}_{k-1} \quad (11)$$

$$\mathbf{z}_k = H_k\mathbf{x}_k + \boldsymbol{\eta}_k \quad (12)$$

245

246

247

248

where  $\mathbf{x}_k$  and  $\mathbf{z}_k$  denote the state and measurement vectors respectively,  $F_{k-1}$  and  $H_k$  denote state transition and measurement matrices respectively,  $\mathbf{n}_{k-1}$  denotes system noise with covariance matrix  $Q_{k-1}$  and  $\boldsymbol{\eta}_k$  denotes measurement noise with covariance matrix  $R_k$ . So the state transition density function and measurement likelihood function are Gaussian:

$$f_{k|k-1}(\mathbf{x}_k|\mathbf{x}_{k-1}) = \mathcal{N}(\mathbf{x}; F_{k-1}\mathbf{x}_{k-1}, Q_{k-1}) \quad (13)$$

$$g_k(\mathbf{z}_k|\mathbf{x}_k) = \mathcal{N}(\mathbf{z}; H_k\mathbf{x}_k, R_k) \quad (14)$$

249

where  $\mathcal{N}(\cdot; m, P)$  denotes a Gaussian density with mean  $m$  and covariance  $P$ .



250 It is also assumed that the probability of survival and detection are state independent  
 251 and constant between time steps

$$p_{S,k}(\mathbf{x}) = p_S \quad (15)$$

$$p_{D,k}(\mathbf{x}) = p_D \quad (16)$$

252 The intensity function of target birth is also assumed to be a Gaussian mixture<sup>32</sup>

$$\gamma_k(\mathbf{x}_k) = \sum_{i=1}^{J_{\gamma,k}} w_{\gamma,k}^{(i)} \mathcal{N}(\mathbf{x}; m_{\gamma,k}^{(i)}, P_{\gamma,k}^{(i)}) \quad (17)$$

253 where  $J_{\gamma,k}$ ,  $w_{\gamma,k}^{(i)}$ ,  $m_{\gamma,k}^{(i)}$ ,  $P_{\gamma,k}^{(i)}$ ,  $i = 1, \dots, J_{\gamma,k}$  are given model parameters that determine the  
 254 shape of the birth intensity function, which is derived in Section III.B.2.

255 The algorithm then consists of the following steps:

256 **Step 0: Initialization.** At the initialization (time  $k=0$ ) the intensity function  $v_0$  is a  
 257 mixture of  $J_0$  Gaussian components

$$v_0(\mathbf{x}) = \sum_{i=1}^{J_0} w_0^{(i)} \mathcal{N}(\mathbf{x}; m_0^{(i)}, P_0^{(i)}) \quad (18)$$

258 In this study  $J_0$  is initialized randomly to be between 1 and 10 components, means  $m_0$   
 259 of those components are drawn randomly from a uniform distribution between 2 and 30 kHz  
 260 and the initial covariance  $P_0$  is set to be the same as the system noise covariance,  $Q_{k-1}$ . The

261 initial weights of all components are the same and are set to  $w_0 = 1/J_0$  .

262 Each component is assigned a unique tag (identifier),  $L_0^{(i)}$ , to form a set  $L_0 = \{L_0^{(i)}\}_{i=1}^{J_0}$  7,25.

263 **Step 1: Prediction.** In this step the Kalman filter prediction equations are used to  
 264 predict means ( $m$ ) and covariances ( $P$ ) of the Gaussian components representing existing  
 265 whistles. The weights ( $w$ ) for existing whistles depend on the probability of survival,  $p_S$ .

266 The predicted intensity of existing whistles,  $v_{k|k-1}(\mathbf{x}, 0)$ , at time  $k$  is a Gaussian mixture  
 267 of the form<sup>32,7</sup>:

$$v_{k|k-1}(\mathbf{x}, 0) = p_S \sum_{j=1}^{J_{k-1}} w_{k-1}^{(j)} \mathcal{N}(\mathbf{x}; m_{k|k-1}^{(j)}, P_{k|k-1}^{(j)}) \quad (19)$$

$$m_{k|k-1}^{(j)} = F_{k-1} m_{k-1}^{(j)} \quad (20)$$

$$P_{k|k-1}^{(j)} = F_{k-1} P_{k-1}^{(j)} F_{k-1}^t + Q_{k-1} \quad (21)$$

268 where  $J_{k-1}$  denotes the number of existing whistles derived from the previous time step  
 269 (combination of existing and newborn whistles) and  $w_{k-1}$  denotes the weights from the  
 270 previous time step.

271 In this step  $J_{\gamma,k}$  new Gaussian components, representing newborn whistles, are also  
 272 created according to the birth model (defined in Section III.B.2, Eqs. (38 and 39)).

273 The tags of the Gaussian components in this step are maintained separately; exist-  
 274 ing whistles keep their tags,  $L_{k|k-1}$ , from the previous time step and new tags,  $L_{\gamma,k}^{(i)}$ ,  $i =$

275  $1, \dots, J_{\gamma,k}$ , are assigned to Gaussians introduced by the birth model so that

$$L_{k|k-1} = L_{k-1} \quad (22)$$

$$L_{\gamma,k} = \{L_{\gamma,k}^{(1)}, \dots, L_{\gamma,k}^{(J_{\gamma,k})}\} \quad (23)$$

276 **Step 2: Update.** In this step the predicted means and covariances of existing and  
 277 newborn whistles are updated using the Kalman filter update equations. The predicted  
 278 weights are updated with the PHD equation. The update is performed separately for existing  
 279 and newborn whistles, Eqs. (8) and (9) respectively.

280 For the existing whistles the posterior intensity function at time  $k$  is given by a Gaussian  
 281 mixture<sup>32,7</sup>:

$$v_k(\mathbf{x}, 0) = (1 - p_D)v_{k|k-1}(\mathbf{x}, 0) + \sum_{\mathbf{z} \in Z_k} \sum_{j=1}^{J_{k|k-1}} w_k^{(j)}(\mathbf{z}) \mathcal{N}(\mathbf{x}, m_k^{(j)}(\mathbf{z}), P_k^{(j)}) \quad (24)$$

282 where  $(1 - p_D)$  denotes the probability of missed detection at current time  $k$ ;  $\mathbf{z}$  denotes an  
 283 individual measurement in the measurement set  $Z_k$  at time  $k$  and

$$w_k^{(j)}(\mathbf{z}) = \frac{p_D w_{k|k-1}^{(j)} g_k^{(j)}(\mathbf{z})}{\mathcal{L}(\mathbf{z})} \quad (25)$$

$$g_k^{(j)}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; H_k m_{k|k-1}^{(j)}, R_k + H_k P_k^{(j)} H_k^t) \quad (26)$$

$$m_k^{(j)}(\mathbf{z}) = m_{k|k-1}^{(j)} + K_k^{(j)}(\mathbf{z} - H_k m_{k|k-1}^{(j)}) \quad (27)$$

$$P_k^{(j)} = [I - K_k^{(j)} H_k] P_{k|k-1}^{(j)} \quad (28)$$

$$K_k^{(j)} = P_{k|k-1}^{(j)} H_k^t (H_k P_{k|k-1}^{(j)} H_k^t + R_k)^{-1} \quad (29)$$

284 where  $K_k$  denotes the Kalman gain and  $I$  denotes the identity matrix.

285 For the newborn whistles the posterior intensity function at time  $k$  is also a Gaussian  
 286 mixture

$$v_k(\mathbf{x}, 1) = \sum_{\mathbf{z} \in Z_k} \sum_{j=1}^{J_{\gamma,k}} w_{\gamma,k}^{(j)}(\mathbf{z}) \mathcal{N}(\mathbf{x}, m_{\gamma,k}^{(j)}(\mathbf{z}), P_{\gamma,k}^{(j)}) \quad (30)$$

287 where  $m_{\gamma,k}^{(j)}(\mathbf{z})$  and  $P_{\gamma,k}^{(j)}$  are calculated with Kalman update equations, in the same way as  
 288 in the equations above and the weights are updated according to Eq. (9)

$$w_{\gamma,k}^{(j)}(\mathbf{z}) = \frac{w_{\gamma,k}^{(j)} g_{\gamma,k}^{(j)}(\mathbf{z})}{\mathcal{L}(\mathbf{z})} \quad (31)$$

289 where

$$\mathcal{L}(\mathbf{z}) = \kappa_k(\mathbf{z}) + \sum_{l=1}^{J_{\gamma,k}} w_{\gamma,k}^{(l)} g_{\gamma,k}^{(l)}(\mathbf{z}) + p_D \sum_{l=1}^{J_{k|k-1}} w_{k|k-1}^{(l)} g_k^{(l)}(\mathbf{z}) \quad (32)$$

$$g_{\gamma,k}^{(l)}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; H_k m_{\gamma,k}^{(l)}, R_k + H_k P_{\gamma,k}^{(l)} H_k^t) \quad (33)$$

290 At the end of the update step, there are  $(1 + |Z_k|)J_{k|k-1}$  Gaussian components,  $(1 + |Z_k|)$   
 291 for each predicted Gaussian<sup>32</sup> for existing whistles and  $|Z_k|J_{\gamma,k}$  Gaussian components for  
 292 newborn whistles. The same tag is assigned to each of the associated predicted and updated  
 293 Gaussian components to form the set<sup>7,25</sup>

$$L_k = L_{k|k-1}^{v_{k|k-1}} \cup L_{k|k-1}^{z_1} \cup \dots \cup L_{k|k-1}^{z_{|Z_k|}} \quad (34)$$

294 for existing whistles and for newborn

$$L_{\gamma,k} = L_{\gamma,k}^{z_1} \cup \dots \cup L_{\gamma,k}^{z_{|Z_k|}} \quad (35)$$

295 The intensities and tags of existing and newborn whistles are then joined and predicted  
 296 jointly in the next time step.

297 With every iteration the number of Gaussian terms will increase, increasing the com-  
 298 putational cost of the algorithm. To control this, pruning and merging schemes are applied  
 299 to the mixture at the end of the update step.

300 **Step 3: Pruning and Merging.** Pruning is achieved by truncating all components  
 301 with small weights by applying a pruning threshold,  $T_r$ . In the merging stage, the components  
 302 that are close together are merged into a single Gaussian component based on a merging  
 303 threshold  $U$ . The distance is computed with a Mahalanobis distance measure<sup>32</sup>.

304 Additionally, to further reduce the computational load, if the number of Gaussian  
 305 components exceeds the desired maximum number of components ( $J_{max}$ ), only the  $J_{max}$   
 306 Gaussian components with the largest weights are kept in the recursion.

307 The values for  $T_r$ ,  $U$ ,  $J_{max}$  are discussed in Section III.B.2 and listed in Table I.

308 **Step 4: State estimation and tracking.** At the end of each recursion the pruned  
 309 Gaussian mixture represents the posterior intensity function  $v_k(\cdot)$  and the means of the  
 310 Gaussian components therefore represent local maxima of  $v_k(\cdot)$ . By taking the Gaussians  
 311 that have weights greater than some threshold  $w_{th}$  (derived in Section III.B.2 and listed in  
 312 Table I), the multi-target states are estimated<sup>32,7</sup>. This step does not affect the main GM-  
 313 PHD recursion. The individual whistles are then tracked from the estimated states based  
 314 on their tags. When a track of a whistle exceeds 150 ms then it is labelled as a detection.  
 315 The 150 ms length threshold was selected based on the study by Roch *et al.*<sup>28</sup> and serves to  
 316 reduce the false detections.

## 317 *2. Definition of the models and parameter selection*

### 318 **State space models for dolphin whistles**

319 The whistle state vectors in this study consist of frequency  $f$  and chirp rate  $\alpha$  (rate of  
 320 change of frequency)<sup>36</sup>:

$$\mathbf{x}_k = [f, \alpha]^t \quad (36)$$

321 where  $[\cdot]^t$  denotes the transpose.

322 The system model (11) in current application uses the state transition matrix  $F_{k-1} =$   
 323  $\begin{bmatrix} 1 & \Delta \\ 0 & 1 \end{bmatrix}$ , where  $\Delta$  denotes the time interval between overlapping spectral windows and is  
 324 related to the sampling frequency ( $f_s$ ),  $\Delta = (w_w/2)/f_s$ , where  $w_w$  denotes the length of the  
 325 window. The system noise,  $\mathbf{n}_{k-1}$ , in this model is independent Gaussian white noise with a  
 326 covariance matrix  $Q_{k-1}$ . Initially,  $Q_{k-1}$  was defined as  $Q_{k-1} = \text{diag}[\sigma_f^2, \sigma_\alpha^2]$ , where  $\sigma_f$  and  $\sigma_\alpha$   
 327 denote the standard deviations of the frequency and chirp respectively, here  $\sigma_f = 70.7$  and  
 328  $\sigma_\alpha = 3.2 \times 10^3$ .

329 This noise covariance matrix was then refined by running the GM-PHD filter (described  
 330 in previous Section III.B.1) on the training data and calculating the mean noise covariance,  
 331 resulting in

$$Q_{k-1} = \begin{bmatrix} \sigma_f^2 & \sigma_{f,\alpha} \\ \sigma_{f,\alpha} & \sigma_\alpha^2 \end{bmatrix} \quad (37)$$

332 where the refined standard deviations of frequency and chirp are  $\sigma_f = 70.8$  and  $\sigma_\alpha =$

333  $7.35 \times 10^3$  and the off-diagonal element is  $\sigma_{f,\alpha} = 408.4^2$ .

334 The measurement model (12) uses the measurement matrix  $H_k = [1, 0]$ , indicating that  
 335 only the frequency information is measured. The measurement noise,  $\boldsymbol{\eta}_k$ , is independent  
 336 Gaussian white noise with covariance matrix  $R_k$ .  $R_k$  in this study is defined as a variance of  
 337 a uniform random variable and is therefore  $b_w^2/12$  where  $b_w$  denotes bin width and is equal  
 338 to  $b_w = f_s/w_w$ .

### 339 Other models and parameters

340 In addition to the system (11) and measurement (12) models required by standard  
 341 tracking methods, the PHD filter requires definition of additional models and parameters  
 342 that govern the GM-PHD recursion. All of these are application dependent. Some of the  
 343 parameters can be determined analytically, but some parameters need to be estimated from  
 344 training data.

345 The additional models needed for the GM-PHD filter, model the birth and the clutter  
 346 intensities. The birth model defines where in the state space new whistles are likely to appear.  
 347 If a whistle appears in a region that is not covered by the predefined birth intensity then the  
 348 PHD filter will not detect it<sup>27</sup>. Since dolphin whistles typically occur in a frequency band  
 349 between 2 and 30 kHz<sup>14</sup>, making the birth intensity diffuse over such a large region would  
 350 increase the computational load. Therefore, the birth intensity in this study is based on the  
 351 available measurements<sup>27</sup> and the new whistles are created as follows. In each time step  $k$ ,



352  $J_{\gamma,k}$  newborn whistles are created, where  $J_{\gamma,k}$  corresponds to the number of measurements  
 353 in the measurement set  $\mathbf{Z}_k$  at time  $k$ . Each newborn whistle is a Gaussian component and  
 354 is therefore characterized by a mean ( $m_{\gamma,k}^{(i)}$ ), a weight ( $w_{\gamma,k}^{(i)}$ ) and a covariance ( $P_{\gamma,k}^{(i)}$ ), where  
 355  $i = 1, \dots, J_{\gamma,k}$ . The covariance of the  $i$ -th newborn whistle is set to be  $Q_{k-1}$  (Eq. 37).  
 356 The frequency component of the mean of the  $i$ -th newborn whistle ( $\{m_{\gamma,k}^{(i)}\}_f$ ) is obtained by  
 357 drawing from a Gaussian mixture centred on the measurements and the chirp component of  
 358 the mean ( $\{m_{\gamma,k}^{(i)}\}_\alpha$ ) is set to zero:

$$\begin{aligned} \{m_{\gamma,k}^{(i)}\}_f &\sim \frac{1}{J_{\gamma,k}} \sum_{j=1}^{J_{\gamma,k}} \mathcal{N}(x; z_{f,k}^{(j)}, 0.01 z_{f,k}^{(j)}) \\ \{m_{\gamma,k}^{(i)}\}_\alpha &= 0 \end{aligned} \quad (38)$$

359 where  $z_{f,k}$  denotes frequency measurements at time  $k$ . The weight of the  $i$ -th newborn  
 360 whistle is computed as

$$w_{\gamma,k}^{(i)} = \frac{p_{start}(z_{f,k}^{(i)})}{J_{\gamma,k}} \quad (39)$$

361 where  $p_{start}(z_{f,k}^{(i)})$  is a value of the log-normal pdf of starting frequencies of whistles (that was  
 362 obtained from the training data) at a particular frequency  $z_{f,k}^{(i)}$ .

363 The clutter (false detections) intensity used in the present study was computed as  
 364 follows. It is assumed that clutter is uniformly distributed over the frequency range (2 to 50

365 kHz) and is constant with respect to time. The average number of clutter points ( $r$ ) per time  
 366 step was estimated based on the training data. The training data were pre-processed using  
 367 the technique described in Section III.A. The number of identified spectral peaks per time  
 368 step was compared to the number of annotated whistle peaks from the analyst-annotated  
 369 data. From this the average number of clutter points can be computed. It was determined  
 370 that our pre-processing technique results in  $r = 10$  clutter points per time step, giving the  
 371 clutter intensity of  $\kappa_k = r/A$ , where  $A$  denotes the bandwidth over which clutter can occur,  
 372 which is 48 kHz for this study.

373 In addition to the models for birth and clutter intensities, the GM-PHD filter requires  
 374 the selection of five other parameters;  $p_S$ ,  $p_D$ ,  $U$ ,  $T_r$ ,  $w_{th}$ . Parameters determined analytically  
 375 in this study were probability of survival ( $p_S$ ) and merging threshold ( $U$ ). Probability of  
 376 survival,  $p_S$ , determines how likely the whistle is to survive from one time step to another.  
 377 As such it will depend on the average length of the whistles, specifically one can show that  
 378  $p_S = 1 - (1/\bar{k})$ , where  $\bar{k}$  is the average length of the whistles expressed in time steps.

379 The average length of whistles was calculated from the study by Oswald *et al.*<sup>23</sup>, where  
 380 four species were the same as in the present study. The average length was 0.875 s, which  
 381 equates to 165 time steps (since the time step used in this study is 5.3 ms), giving a  $p_S$  of  
 382 0.994.

383 The merging threshold,  $U$ , determines which components are merged and is based on

384 the Mahalanobis distance between two Gaussians. Mahalanobis distances are characterized  
385 by the Chi-squared distribution with  $d$ -degrees of freedom (where  $d$  equals the number of  
386 variables; in our case, where the state vector consists of frequency and chirp rate,  $d$  is equal  
387 to 2). For a Chi-squared distribution with 2-degrees of freedom, 99% of all the values coming  
388 from this distribution will lie within 9.2. Therefore merging threshold  $U$  was set to 10.

389 Parameters determined experimentally from the data were probability of detection ( $p_D$ ),  
390 pruning threshold ( $T_r$ ) and weight threshold ( $w_{th}$ ). All three parameters were determined  
391 experimentally by running the GM-PHD filter on the training data and by selecting the  
392 values that resulted in the best performance. The parameters used in the GM-PHD for  
393 dolphin whistle tracking are summarized in Table I.

### 394 ***3. GM-PHD performance evaluation***

395 After applying the GM-PHD filter described above to the acoustic recordings of dolphin  
396 whistles, the detected list of time against frequency peaks for each whistle was compared to  
397 the ground truth hand-annotated data in order to evaluate the filter's performance. First  
398 the whistles in the hand-annotated data were evaluated in terms of whistles' duration and  
399 SNR. The ground truth whistle was only expected to be detected if its duration exceeded  
400 150 ms and if its SNR exceeded 10 dB for at least one third of its duration (following Roch  
401 *et al.*<sup>28</sup>). Ground truth whistles meeting these selection criteria were termed valid.

402 Next the output of the GM-PHD filter was compared to the ground truth whistles. The

403 detected whistle was considered a match (true positive) to a ground truth whistle if its timing  
404 overlapped with the ground truth whistle and if the mean difference between the detected  
405 whistle path and ground truth whistle path did not exceed 3 frequency bins (281 Hz). If  
406 the detected whistle exceeded that criteria, it was considered as false positive. It should be  
407 noted that detected whistles were matched to ground truth whistles regardless of whether the  
408 ground truth whistles met the selection criteria (*i.e.* if they were valid). However only the  
409 whistles that matched valid ground truth whistles were considered in the evaluation metrics  
410 that describe the quality and quantity of matches<sup>28</sup>. Also, since the hand-annotations were  
411 only applied to the frequencies between 4.5 kHz and 50 kHz, all the detected whistles that  
412 had over 40% of the contour below the 4.5 kHz were not taken into account in the evaluation.

413       The performance of the GM-PHD filter was measured in terms of recall, precision,  
414 fragmentation, deviation and coverage. For detailed description see Roch *et al.*<sup>28</sup>. Recall  
415 measures the percentage of the expected detections that are retrieved, precision measures the  
416 percentage of the detections that are correct. For the detected whistles that matched valid  
417 ground truth whistles (true positives), three additional performance metrics are computed;  
418 fragmentation, mean deviation and coverage. Fragmentation measures the average number  
419 of detections per ground truth whistle, deviation measures the average frequency deviation  
420 between the path of ground truth whistle and its corresponding detection and coverage  
421 measures the average percentage of a ground truth whistle that is matched.

#### 422 **IV. RESULTS**

423 Across all six species in the selected database, 9192 ground truth whistles met the  
424 selection criteria. The performance of the GM-PHD detector for each species is summarized  
425 in Table II. The GM-PHD detector tracked whistles successfully with overall precision of  
426 85% and overall recall of 71.8%. Across all species, the whistles were tracked precisely  
427 with average deviation from the whistle path of 104 Hz and with coverage of 80.3%. An  
428 example of GM-PHD tracking is shown in Figure 1. The detector tracked the paths of  
429 individual whistles when overlapping whistles were present, although occasional “breaking”  
430 of the whistle contours still occurred (on average there were 1.2 fragments per whistle across  
431 all species). An example is shown in Figure 2, where both successful tracking through a  
432 crossing and some breaking of the whistle track can be observed.

#### 433 **V. DISCUSSION**

434 This study demonstrated the use of a MTT technique for tracking odontocete whistle  
435 contours. The proposed adaptation of the GM-PHD filter successfully simultaneously tracked  
436 whistles in complex environments (overlapping whistles, missed detections, clutter present)  
437 for all species investigated, despite the parameter optimization being performed on only three  
438 of the species in the overall dataset. This suggests that the GM-PHD detector formulation  
439 in this study is widely applicable to whistle tracking problems across a wide range of species.

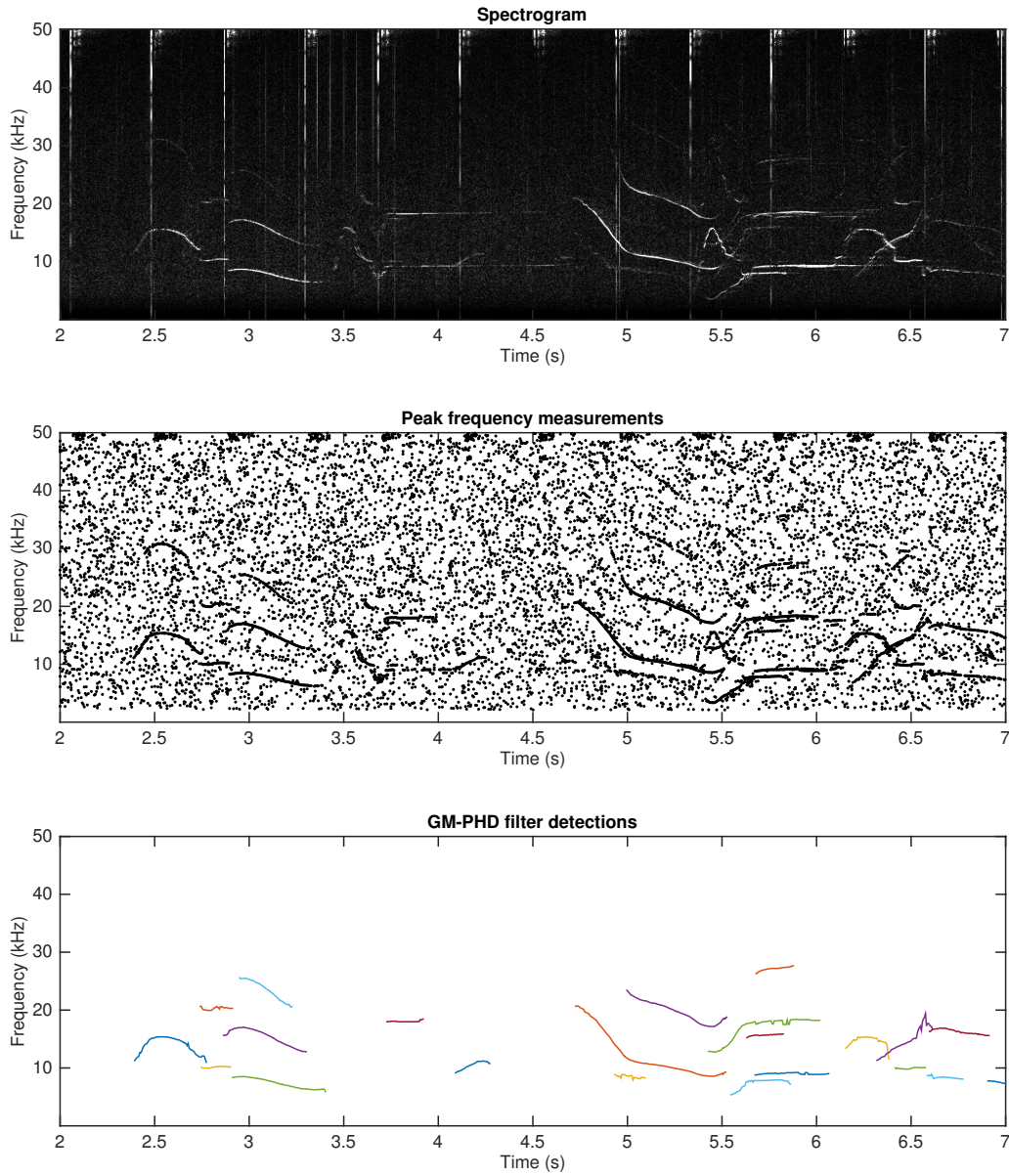


Figure 1: (Color online) Detected whistles with the GM-PHD filter. Spectrogram of raw data is shown (top), peak frequencies measurements (peaks 8 dB above background noise) (middle) and tracked whistles (bottom) where GM-PHD filter detections are shown.

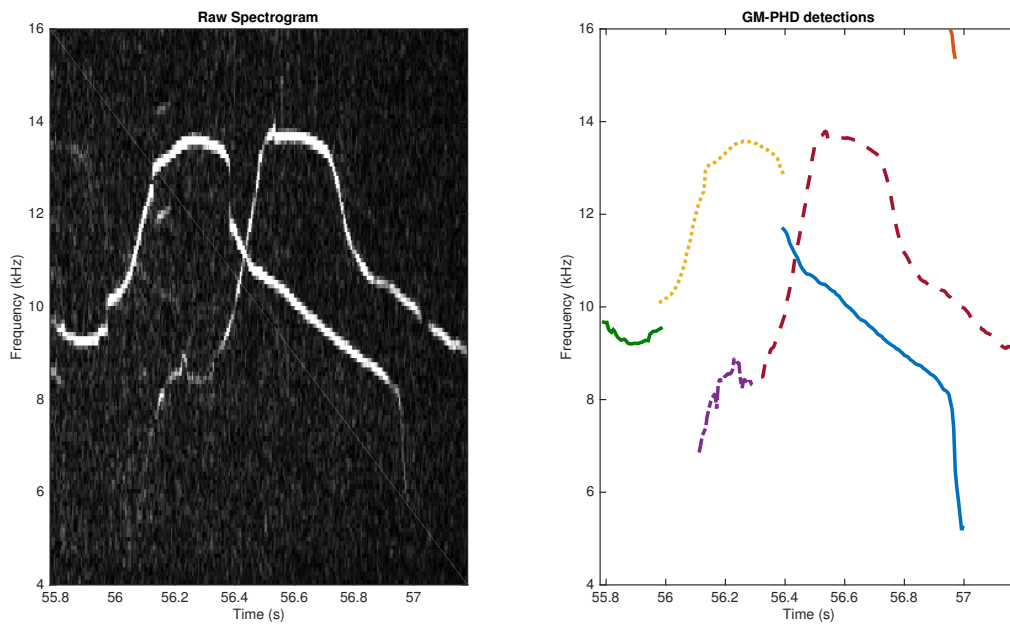


Figure 2: (Color online) Detection of crossing whistles with the GM-PHD filter. Spectrogram of raw data (left) and tracked whistles with GM-PHD filter (right) are shown.

440 The precision for all species was generally higher than the recall. It should be noted  
441 that the precision for *T.truncatus* is slightly lower than the precision for other species (Table  
442 II). When *T.truncatus* files were investigated, it was observed that one file in particular  
443 contained many burst pulses, which the GM-PHD filter detected as whistles, resulting in a  
444 high number of false positives. In general, there is a trade-off between the precision and recall,  
445 and the recall could be increased by allowing shorter fragments to be detected (currently a  
446 150 ms threshold is used). However this would, in turn, lower the precision since it would  
447 likely increase the number of false positive detections.

448 For the detected whistles that matched the ground truth data (true positive detections),  
449 the performance was quite good. The detected whistles followed the path of the annotated  
450 ground truth data closely (within about 1 frequency bin width) and covered the majority  
451 of the individual contours. The whistles were mainly detected as a single contour, but  
452 were occasionally "broken" into more fragments. The breaking of contours mainly occurred  
453 where the amplitude of the whistle dropped below the SNR used to detect spectral peaks and  
454 therefore there were no measurements passed to the GM-PHD detector. While the GM-PHD  
455 filter allows for missed detections, it cannot continue to track a target if the measurements are  
456 absent for several continuous time steps (an example is shown in Figure 2). Also, while the  
457 analyst constructing the ground truth data attempted not to trace whistles where the whistle  
458 path was not obvious, it was observed from manual inspection of some of the annotated files



459 and corresponding spectrograms, that this was not universally applied. This leads to an  
460 increase in the measured fragmentation rate.

461 Comparing the performance of the GM-PHD filter to other filters is difficult, mainly due  
462 to different sound files being used, different pre-processing techniques and different methods  
463 to estimate the SNR. However, the results in Table II demonstrate that performance results  
464 are comparable to those of the graph filter and better than particle filter detailed in Roch  
465 *et al.*<sup>28</sup>. In order to facilitate the comparison between different detectors (that operate on  
466 identified spectral peaks), datasets containing detected spectral peaks used in this study  
467 have been released to MobySound archive.

468 Further improving the GM-PHD filter performance is not a trivial task. In general,  
469 the performance of the filter will greatly depend on the parameter selection and therefore  
470 needs further discussion. In the present study some of the parameters;  $p_D$ ,  $\kappa_k$ ,  $T_r$ ,  $w_{th}$ , were  
471 estimated from the training data set. In particular the probability of detection ( $p_D$ ) was  
472 selected by running the GM-PHD filter on all training data and choosing the value that on  
473 average resulted in the best performance. Since  $p_D$  depends on the SNR, which will change  
474 depending on the environment, the animal's location relative to the sensor and the recording  
475 equipment, significantly different values might have been obtained if different training files  
476 were used. During the GM-PHD recursion the  $p_D$  is assumed to be constant, but between  
477 different recordings the performance could potentially be improved if the  $p_D$  was adjusted

478 to that particular situation. We are currently exploring the methods that would facilitate  
479 this. Another parameter estimated from the training data was the clutter intensity,  $\kappa_k$ .  
480 While the average number of clutter points per time step appeared to be consistent between  
481 species and files in the training set, the value will mainly depend on the threshold used  
482 to generate measurements (detected spectral peaks). The value of  $\kappa_k$  would need to be  
483 adjusted if a different threshold was used or if a different pre-processing or spectral peak  
484 detection strategy was adopted. The selection of pruning ( $T_r$ ) and weight ( $w_{th}$ ) thresholds  
485 mainly affects the computational speed of the algorithm. By selecting higher values for  
486 the two thresholds, fewer Gaussian components remain in the recursion and the speed of the  
487 recursion increases. However, if the selected values are too high, the components representing  
488 whistles start to be excluded from the recursion, which results in a decrease in performance  
489 since fewer whistles are tracked.

490 In addition, the performance of the GM-PHD filter will also crucially depend on the  
491 state-space and birth models used. The birth model in this study was developed from the  
492 proposition by Ristic *et al.*<sup>27</sup>, where the birth model is based on the measurements. The  
493 weights of the newborn whistles were determined based on the probability distribution of  
494 the whistles' start frequencies, which were obtained from the training data. Since training  
495 data encompassed only three species, future work will investigate whether a model based on  
496 more species enhances the performance. The state models, used in this study, describing

497 the evolution of the whistles are based on a simple linear model. Refining this model and  
498 developing a more rigorous method to fit its parameters to the training data should also be  
499 considered.

500 One attraction of the GM-PHD filter is that the formulation of the filter is based on the  
501 mathematical principles and is not *ad-hoc* as some of the other tracking algorithms. Since  
502 the filter is data-association free, it is more computationally efficient than the traditional  
503 MTT methods and can be implemented in real-time. It should be noted that the compu-  
504 tational speed of the algorithm will not only depend on the parameter selection, but also  
505 on the amount of clutter in the measurements. If lower SNR thresholds are used in the  
506 measurement generation (spectral peak detection), more clutter is present in the measure-  
507 ments and the computational load increases, which results in slowing the algorithm. Using  
508 higher thresholds in the spectral peak detection increases the speed of the algorithm, but  
509 some spectral peaks associated with whistles are then missing from the measurements, which  
510 affects the tracking performance. So there is an inherent trade-off between the performance  
511 and computational speed. To illustrate, for the parameters used in this study, the GM-PHD  
512 algorithm implemented in MATLAB (version 8.5 (R2015a)) on a Mac (Os X, processor 2.7  
513 GHz and 8 GB RAM), took 1 min and 48 s to process a file of 1 min duration at 192 kHz  
514 sample rate, that contained 103 hand-annotated whistles.

## 515 VI. CONCLUSIONS

516 The proposed formulation of the GM-PHD filter provides a general and powerful tool  
 517 for simultaneous tracking of odontocete whistle contours. Its performance is comparable  
 518 with the best existing methods, it is computationally efficient and well suited for real-time  
 519 implementation.

## 520 Acknowledgements

521 We would like to thank MobySound archive, DCLDE committee and associated analysts  
 522 for providing the datasets and hand annotations used to test detector's performance in  
 523 this study. We would like to thank Marie Roch and anonymous reviewer for their helpful  
 524 comments on an earlier version of this manuscript. We would also like to thank Slovene  
 525 human resources development and scholarship fund (Ad futura) for funding this research.

## 526 APPENDIX

527 A list of the most frequently used symbols and their meanings.

528	$\eta_k$ and $R_k$	Measurement noise process and its covariance matrix
529	$F_{k-1}$	State transition (system) matrix
530	$f_{k k-1}(\mathbf{x}_k \mathbf{x}_{k-1})$	State transition density
531	$g_k(\mathbf{z} \mathbf{x}_k)$	Likelihood function
532	$\gamma_k(\mathbf{x}_k)$	Intensity function (or PHD) of target births at time $k$

533	$H_k$	Measurement matrix
534	$J_{k-1}$	Number of existing targets deriving from previous time step $k - 1$
535	$J_{\gamma,k}$	Number of newborn targets at time $k$
536	$\kappa_k(\mathbf{z})$	Intensity function (or PHD) of clutter at time $k$
537	$\mathbf{n}_{k-1}$ and $Q_{k-1}$	System noise process and its covariance matrix
538	$p_{D,k}(\mathbf{x}_k) \stackrel{abbr}{=} p_D$	Probability of detection
539	$p_{S,k}(\mathbf{x}_{k-1}) \stackrel{abbr}{=} p_S$	Probability of target's survival from time $k - 1$ to time $k$
540	$p_k(\mathbf{X}_k   \mathbf{Z}_{1:k})$	Posterior pdf of the multi-target state
541	$\beta$	Label $\beta$ denotes newborn targets ( $\beta = 1$ ) or existing targets ( $\beta = 0$ )
542	$v_{k k-1}(\mathbf{x}, \beta)$	Predicted intensity function (or PHD)
543	$v_k(\mathbf{x}, \beta)$	Posterior intensity function (or PHD)
544	$w$ and $w_\gamma$	Weights for existing and newborn Gaussian components (whistles) re-
545		spectively.
546	$\mathbf{x}_k$ and $\mathbf{z}_k$	State and measurement vectors at time $k$
547	$\mathbf{Z}_k$	Multi-target measurement at time $k$

548     **REFERENCES**

- 549 **1.** M. S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, “A tutorial on particle  
550 filters for online nonlinear/non-Gaussian Bayesian tracking”, *IEEE Transactions on*  
551 *Signal Processing* **50**, 174–188 (2002).
- 552 **2.** W. W. Au, “Hearing in whales and dolphins: An overview”, in *Hearing by whales and*  
553 *dolphins*, edited by W. W. L. Au, A. N. Popper, and R. R. Fay, (Springer-Verlag, New  
554 York, 2000), pp. 1–42.
- 555 **3.** S. Baumann-Pickering, S. M. Wiggins, J. A. Hildebrand, M. A. Roch, and H.-U.  
556 Schnitz, “Discriminating features of echolocation clicks of melon-headed whales (*Pe-*  
557 *ponocephala electra*), bottlenose dolphins (*Tursiops truncatus*), and gray’s spinner dol-
- 558 *phins* (*Stenella longirostris longirostris*)”, *Journal of the Acoustical Society of America*  
559 **128**, 2212–2224 (2010).
- 560 **4.** S. Blackman and R. Popoli, *Design and Analysis of Modern Tracking Systems* (Artech  
561 House, Norwood, MA, 1999), p. 1230.
- 562 **5.** S. Bozic, *Digital and Kalman filtering* (Edward Arnold Ltd, London, UK, 1979), p.  
563 157.
- 564 **6.** D. Clark, B.-N. Vo, and J. Bell, “GM-PHD filter multitarget tracking in sonar images”,  
565 in *Proceedings of SPIE 6235 - Signal Processing, Sensor Fusion and Target Recognition*

- 566 XV, edited by I. Kadar, (International Society for Optics and Photonics, Orlando,  
567 Florida, US), pp. 1–8 (2006).
- 568 **7.** D. E. Clark, K. Panta, and B.-N. Vo, “The GM-PHD filter multiple target tracker”,  
569 in *Proceedings of the 9th International Conference on Information Fusion*, (IEEE,  
570 Florence, Italy), pp. 1–8 (2006).
- 571 **8.** S. Datta and C. Sturtivant, “Dolphin whistle classification for determining group iden-  
572 tities”, *Signal processing* **82**, 251–258 (2002).
- 573 **9.** D. Gillespie, M. Caillat, J. Gordon, and P. White, “Automatic detection and classi-  
574 fication of odontocete whistles”, *Journal of the Acoustical Society of America* **134**,  
575 2427–2437 (2013).
- 576 **10.** X. C. Halkias and D. P. Ellis, “Call detection and extraction using bayesian inference”,  
577 *Applied Acoustics* **67**, 1164–1174 (2006).
- 578 **11.** C. Ioana, C. Gervaise, Y. Stphan, and J. I. Mars, “Analysis of underwater mammal  
579 vocalisations using timefrequency-phase tracker”, *Applied Acoustics* **71**, 1070–1080  
580 (2010).
- 581 **12.** A. T. Johansson and P. R. White, “An adaptive filter-based method for robust, auto-

- 582 matic detection and frequency estimation of whistles”, The Journal of the Acoustical  
583 Society of America **130**, 893–903 (2011).
- 584 **13.** A. Kershenbaum and M. A. Roch, “An image processing based paradigm for the  
585 extraction of tonal sounds in cetacean communications”, The Journal of the Acoustical  
586 Society of America **134**, 4435–4445 (2013).
- 587 **14.** M. O. Lammers, W. W. L. Au, and D. L. Herzing, “The broadband social acoustic  
588 signaling behavior of spinner and spotted dolphins”, Journal of the Acoustical Society  
589 of America **114**, 1629–1639 (2003).
- 590 **15.** S. Madhusudhana, E. Oleson, M. Soldevilla, M. Roch, and J. Hildebrand, “Fre-  
591 quency based algorithm for robust contour extraction of blue whale B and D calls”,  
592 in *OCEANS 2008 - MTS/IEEE Kobe Techno-Ocean*, (IEEE, Kobe, Japan), pp. 1–8  
593 (2008).
- 594 **16.** R. P. Mahler, “A theoretical foundation for the Stein-Winter ”Probability Hypothesis  
595 Density (PHD)” multitarget tracking approach”, in *Proceedings of the 2000 MSS Na-*  
596 *tional Symposium on Sensor and Data Fusion* , San Antonio, Texas, US, pp. 99–117  
597 (2000).
- 598 **17.** R. P. Mahler, “Multitarget bayes filtering via first-order multitarget moments”, IEEE  
599 Transactions on Aerospace and Electronic Systems **39**, 1152–1178 (2003).



- 600 **18.** R. P. Mahler, “A survey of PHD filter and CPHD filter implementations”, in *Proceed-*  
601 *ings of SPIE 6567 - Signal Processing, Sensor Fusion and Target Recognition XVI*,  
602 edited by I. Kadar, (International Society for Optics and Photonics, Orlando, Florida,  
603 US), pp. 1–12 (2007).
- 604 **19.** R. P. Mahler, *Statistical multisource-multitarget information fusion*, (Artech House,  
605 Norwood, MA, 2007), p. 856.
- 606 **20.** A. Mallawaarachchi, S. Ong, M. Chitre, and E. Taylor, “Spectrogram denoising and  
607 automated extraction of the fundamental frequency variation of dolphin whistles”, *The*  
608 *Journal of the Acoustical Society of America* **124**, 1159–1170 (2008).
- 609 **21.** T. A. Marques, L. Thomas, J. Ward, N. DiMarzio, and P. L. Tyack, “Estimating  
610 cetacean population density using fixed passive acoustic sensors: An example with  
611 blainvilles beaked whales”, *The Journal of the Acoustical Society of America* **125**,  
612 1982–1994 (2009).
- 613 **22.** D. K. Mellinger, S. W. Martin, R. P. Morrissey, L. Thomas, and J. J. Yosco, “A method  
614 for detecting whistles, moans, and other frequency contour sounds”, *The Journal of*  
615 *the Acoustical Society of America* **129**, 4055–4061 (2011).
- 616 **23.** J. N. Oswald, J. Barlow, and T. F. Norris, “Acoustic identification of nine delphinid

- 617 species in the eastern tropical pacific ocean”, Marine Mammal Science **19**, 20–37  
618 (2003).
- 619 **24.** J. N. Oswald, S. Rankin, J. Barlow, and M. O. Lammers, “A tool for real-time acous-  
620 tic species identification of delphinid whistles”, Journal of the Acoustical Society of  
621 America **122**, 587–595 (2007).
- 622 **25.** K. Panta, D. E. Clark, and B.-N. Vo, “Data association and track management for  
623 the Gaussian mixture probability hypothesis density filter”, IEEE Transactions on  
624 Aerospace and Electronic Systems **45**, 1003–1016 (2009).
- 625 **26.** N. J. Quick and V. M. Janik, “Whistle rates of wild bottlenose dolphins (*Tursiops*  
626 *truncatus*): influences of group size and behavior”, Journal of Comparative Psychology  
627 **122**, 305–311 (2008).
- 628 **27.** B. Ristic, D. Clark, B.-N. Vo, and B.-T. Vo, “Adaptive target birth intensity for  
629 PHD and CPHD filters”, IEEE Transactions on Aerospace and Electronic Systems  
630 **48**, 1656–1668 (2012).
- 631 **28.** M. A. Roch, T. S. Brandes, B. Patel, Y. Barkley, S. Baumann-Pickering, and M. S.  
632 Soldevilla, “Automated extraction of odontocete whistle contours”, Journal of the  
633 Acoustical Society of America **130**, 2212–2223 (2011).

- 634 **29.** M. A. Roch, M. S. Soldevilla, J. C. Burtenshaw, E. E. Henderson, and J. A. Hildebrand,  
635 “Gaussian mixture model classification of odontocetes in the Southern California Bight  
636 and the Gulf of California”, *Journal of the Acoustical Society of America* **121**, 1737–  
637 1748 (2007).
- 638 **30.** M. S. Soldevilla, E. E. Henderson, G. S. Campbell, S. M. Wiggins, J. A. Hildebrand,  
639 and M. A. Roch, “Classification of risso’s and pacific white-sided dolphins using spec-  
640 tral properties of echolocation clicks”, *Journal of the Acoustical Society of America*  
641 **124**, 609–624 (2008).
- 642 **31.** P. L. Tyack and C. W. Clark, “Communication and acoustic behavior of dolphins and  
643 whales”, in *Hearing by whales and dolphins*, edited by W. W. L. Au, A. N. Popper,  
644 and R. R. Fay, (Springer-Verlag, New York, 2000), pp. 156–224.
- 645 **32.** B.-N. Vo and W.-K. Ma, “The Gaussian mixture probability hypothesis density filter”,  
646 *IEEE Transactions on Signal Processing* **54**, 4091–4104 (2006).
- 647 **33.** B.-N. Vo, S. Singh, and A. Doucet, “Sequential Monte Carlo methods for multitarget  
648 filtering with random finite sets”, *IEEE Transactions on Aerospace and Electronic*  
649 *Systems* **41**, 1224–1245 (2005).
- 650 **34.** Y. Wang, H. Meng, H. Zhang, and X. Wang, “Improved GM-PHD tracker with de-

- 651        layed decision”, in *10th International Conference on Signal Processing*, (IEEE, Beijing,  
652        China), pp. 255–258 (2010).
- 653 **35.** C. R. Weir and S. J. Dolman, “Comparative review of the regional marine mammal  
654        mitigation guidelines implemented during industrial seismic surveys, and guidance  
655        towards a worldwide standard”, *Journal of International Wildlife Law and Policy* **10**,  
656        1–27 (2007).
- 657 **36.** P. White and M. Hadley, “Introduction to particle filters for tracking applications in  
658        the passive acoustic monitoring of cetaceans”, *Canadian Acoustics* **36**, 146–152 (2008).

Table I: Summary of parameters used in GM-PHD filter for odontocete whistle tracking.  $p_S$  and  $p_D$  denote probabilities of survival and detection respectively;  $U$ ,  $T_r$  and  $w_{th}$  denote merging, pruning and weight thresholds respectively and  $J_{max}$  denotes maximum allowed number of Gaussian components in one iteration.

$p_S$	$p_D$	$U$	$T_r$	$w_{th}$	$J_{max}$
0.994	0.85	10	0.001	0.009	100

Table II: Performance of the GM-PHD filter for detection of odontocete whistle contours.  $N$  files denotes number of audio files used, *Valid whistles* denotes the number of ground truth whistles that met the selection criteria,  $\mu$ *Deviation* denotes average deviation,  $SD$  denotes standard deviation. The summary performance is computed across all ground truth whistles that met the criteria and is not the average of file performances.

Species	N files	Valid whistles	Recall	Precision	Coverage $\pm$ SD (%)	Fragments $\pm$ SD	$\mu$ Deviation $\pm$ SD (Hz)
<i>D.capensis</i>	7	1859	72.1	91.1	80.6 $\pm$ 22.3	1.2 $\pm$ 0.4	94 $\pm$ 51
<i>D.delphis</i>	10	1931	71.6	85.7	79.2 $\pm$ 23.2	1.2 $\pm$ 0.4	96 $\pm$ 53
<i>P.electra</i>	3	756	66.8	91.3	79.8 $\pm$ 21.6	1.1 $\pm$ 0.3	92 $\pm$ 54
<i>S.longirostris</i>	3	869	76.4	93.5	77.2 $\pm$ 22.2	1.2 $\pm$ 0.5	100 $\pm$ 51
<i>S.frontalis</i>	2	242	70.7	88.6	86.1 $\pm$ 19.4	1.1 $\pm$ 0.3	117 $\pm$ 63
<i>T.truncatus</i>	15	3535	71.7	78.3	81.2 $\pm$ 21.2	1.2 $\pm$ 0.5	117 $\pm$ 53
<b>OVERALL</b>	<b>40</b>	<b>9192</b>	<b>71.8</b>	<b>85.0</b>	<b>80.3<math>\pm</math>22.0</b>	<b>1.2<math>\pm</math>0.4</b>	<b>104<math>\pm</math>54</b>

659

## Figure Captions

660 Figure 1. Detected whistles with the GM-PHD filter. Spectrogram of raw data is shown  
661 (top), peak frequencies measurements (peaks 8 dB above background noise) (middle) and  
662 tracked whistles (bottom) where GM-PHD filter detections are shown.

663 Figure 2. Detection of crossing whistles with the GM-PHD filter. Spectrogram of raw data  
664 (left) and tracked whistles with GM-PHD filter (right) are shown.