


The Journal of the Acoustical Society of America
Automated extraction of dolphin whistles- a Sequential Monte Carlo Probability Hypothesis Density (SMC-PHD) approach
--Manuscript Draft--

Manuscript Number:	JASA-05519R1
Full Title:	Automated extraction of dolphin whistles- a Sequential Monte Carlo Probability Hypothesis Density (SMC-PHD) approach
Article Type:	Regular Article
Corresponding Author:	Pina Gruden Research Corporation of the University of Hawaii (RCUH) Honolulu, HI UNITED STATES
First Author:	Pina Gruden
Order of Authors:	Pina Gruden Paul White, Prof.
Section/Category:	SPECIAL ISSUE ON MACHINE LEARNING IN ACOUSTICS
Keywords:	automated whistle tracking; multi-target Bayesian; probability hypothesis density; sequential Monte Carlo
Abstract:	<p>The need for automated methods to detect and extract marine mammal vocalizations from acoustic data has increased in the last few decades due to the increased availability of long-term recording systems. Automated dolphin whistle extraction represents a challenging problem due to the time-varying number of overlapping whistles present in, potentially, noisy recordings. Typical methods utilize image processing techniques or single target tracking, but often result in fragmentation of whistle contours and/or partial whistle detection. This study casts the problem into a more general statistical multi-target tracking framework, and uses the probability hypothesis density (PHD) filter as a practical approximation to the optimal Bayesian multi-target filter. In particular, a particle version, referred to as a Sequential Monte Carlo PHD (SMC-PHD) filter, is adapted for frequency tracking and specific models are developed for this application. Based on these models, two versions of the SMC-PHD filter are proposed and their performance is investigated on an extensive real-world dataset of dolphin acoustic recordings. The proposed filters are shown to be efficient tools for automated extraction of whistles, suitable for real-time implementation.</p>

CONFIDENTIAL



Click here to access/download

**Rebuttal Letter / Helpful/Supporting Material for
Reviewer**

Response_letter_JASA-05519.pdf



Click here to access/download

Reviewer PDF with line numbers, inline figures and captions

Manuscript_JASA-05519_Revision.pdf



Automated extraction of dolphin whistles - a Sequential Monte Carlo Probability

Hypothesis Density (SMC-PHD) approach

Pina Gruden^{1, a)} and Paul R. White¹

*Institute of Sound and Vibration Research, University of Southampton, Highfield,
Hants, SO17 1BJ, UK*

1 The need for automated methods to detect and extract marine mammal vocalizations
2 from acoustic data has increased in the last few decades due to the increased availabil-
3 ity of long-term recording systems. Automated dolphin whistle extraction represents
4 a challenging problem due to the time-varying number of overlapping whistles present
5 in, potentially, noisy recordings. Typical methods utilize image processing techniques
6 or single target tracking, but often result in fragmentation of whistle contours and/or
7 partial whistle detection. This study casts the problem into a more general statistical
8 multi-target tracking framework, and uses the probability hypothesis density (PHD)
9 filter as a practical approximation to the optimal Bayesian multi-target filter. In
10 particular, a particle version, referred to as a Sequential Monte Carlo PHD (SMC-
11 PHD) filter, is adapted for frequency tracking and specific models are developed for
12 this application. Based on these models, two versions of the SMC-PHD filter are
13 proposed and their performance is investigated on an extensive real-world dataset of
14 dolphin acoustic recordings. The proposed filters are shown to be efficient tools for
15 automated extraction of whistles, suitable for real-time implementation.

^{a)} pgruden@hawaii.edu; Currently at: Research Corporation of the University of Hawaii (RCUH), Honolulu, HI, 96822, US.

16 I. INTRODUCTION

17 The detection and extraction of marine mammal calls is a crucial first step in many
18 applications, such as abundance estimation¹, species identification²⁻⁴, behavioural studies⁵,
19 and is used in mitigation during industrial activities⁶. These applications can involve data
20 collection over extended periods of time and result in large quantities of data accumulating,
21 in which case automated analysis tools become a necessity. This work proposes and validates
22 a multi-target tracking approach for automated whistle extraction using Sequential Monte
23 Carlo Probability Hypothesis Density (SMC-PHD) filters, including specific models tailored
24 for tracking multiple whistles in a real-world dataset.

25 When extracting tonal sounds, such as narrowband frequency modulated delphinid whis-
26 tles, the aim is to describe the contour of each call - *i.e.*, the frequency evolution of a
27 tonal signal through time. This process can be referred to as extraction^{7,8}, detection^{3,9}
28 or tracking^{10,11}. The typical signal processing work-flow involves a pre-processing stage,
29 where the effect of the background noise and interfering signals is reduced, followed by the
30 extraction of whistles^{3,7,11}. Most methods for automated whistle extraction are based on
31 spectrogram techniques and aim to identify the strongest spectral peaks, which are then
32 connected to form continuous whistle contours^{3,7-9,11}. For the purpose of this work, a "mea-
33 surement" is defined to be the frequency associated with a single spectral peak identified
34 within a given spectral window. The number of measurements within a spectral window
35 varies between windows and is unknown *a priori*.

36 Whistle extraction represents a challenging problem for several reasons. One is that the
37 amplitude of a call changes rapidly throughout the whistle duration, which may cause the
38 energy in the whistle to rise above, and then fall below, the detection threshold³, resulting
39 in sections of the whistle being missed. Further, there are usually many overlapping whistles
40 and other interfering sounds present³, which can mask the signal being tracked. The end
41 result is a partial extraction and fragmentation of the contours.

42 This can hamper certain applications, such as classifiers, if they require the extraction
43 of the full whistle contours². While it is still possible to use whistle fragments to identify
44 species^{3,12}, it is expected that as the length of the detected whistle contour increases, the
45 species specific information contained in that detection improves and therefore enhances the
46 classification. For instance, Ref.³ found that as the fragment length of the whistle contour
47 increased, the classification performance increased as well. This enhancement can prove
48 significant in situations where a mix of rare and abundant species are present¹³.

49 The goal of this study is thus to improve on the whistle extraction process, by casting
50 it into a multi-target tracking (MTT) framework, which allows for simultaneous tracking of
51 multiple objects of interest from the noisy measurements in the presence of missed detec-
52 tions, and false alarms (*i.e.*, clutter, additional measurements not generated by a whistle).
53 Additionally, in contrast to the majority of automated methods for whistle contour extrac-
54 tion, the MTT accounts for the time-varying number of whistles by modelling their birth
55 (when a whistle starts) and their death (when a whistle ends).

56 **A. Background**

57 A tracker is based on defining a system, whose configuration is defined at time k , by the
 58 parameter values in the state vector \mathbf{x}_k . The dynamics of the system describe how the states
 59 evolve with time, and are encapsulated in the state equation (1). The vector, \mathbf{z}_k , contains
 60 the value of the quantities that are measured at time k . The measurements are related to
 61 the system states through the measurement equation (2):

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

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

62 where F_s is the function which combines the previous state vector and the system noise
 63 process \mathbf{n}_k to generate the current state and G_m is the function computing the measure-
 64 ment, \mathbf{z}_k , combining the current state vector with a measurement noise process, $\boldsymbol{\eta}_k$. For
 65 the whistle tracking problem, we employ a state vector consisting of two parameters - the
 66 instantaneous frequency and its derivative (the chirp rate) - whilst the only measurements
 67 available are measurements of the instantaneous frequency. The specific form of the state
 68 and measurement equations used in this paper are discussed in Section II C.

69 A Bayesian recursive filtering approach is frequently adopted in single-target tracking
 70 problems to estimate the state of a system from a sequence of noisy measurements¹⁴. The
 71 Bayes filter aims to compute the posterior probability density function (pdf) of the state
 72 estimate at each time step, and is based on a two stage recursion¹⁴. The first stage, the
 73 prediction step, uses the state dynamics (1) to compute an *a priori* estimate of the state's

74 density function. Whilst the second stage, the update step, updates that density based
 75 on the newly available measurement, leading to an estimate of the posterior pdf for the
 76 state vector¹⁴. An analytic solution to the Bayes filter under the assumptions of: a single
 77 target, Gaussianity of the noise processes and linearity of the underlying models, can be
 78 obtained. The resulting method is the Kalman filter¹⁵. A more general approach for a single
 79 target, avoiding the need for the assumptions of linearity and Gaussianity, is offered by the
 80 Sequential Monte Carlo (SMC) filter (or particle filter)¹⁴, which is the basis of much of what
 81 follows here.

82 Finite Set Statistics (FISST) provides a suitable framework within which an MTT
 83 Bayesian filter can be constructed^{16,17}. FISST models the states of the targets and the
 84 measurements using the concept of random finite sets (RFS), and transforms a multi-sensor,
 85 multi-target problem into a mathematically equivalent single-sensor, single-target problem.
 86 A RFS is an object in which the unordered elements have random values, as in any multivari-
 87 ate random process, but in addition to which the number of elements (the set cardinality)
 88 is also random¹⁷.

89 A multi-target Bayesian filter¹⁷ determines, at each iteration, the full posterior pdf of
 90 the multi-target state, which makes it computationally intractable in practice, especially
 91 when there are a large number of targets. To overcome this problem, one solution is to
 92 use a filter based on the Probability Hypothesis Density (PHD), $v_k(\mathbf{x}|\mathbf{Z}_{1:k})$ (where $\mathbf{Z}_{1:k}$ is
 93 the set of measurements \mathbf{z} at times 1 to k), which is the first-order moment of the multi-
 94 target posterior^{18,19}. Note that for compactness and clarity, the dependence of v_k on $\mathbf{Z}_{1:k}$ is
 95 suppressed in subsequent equations. The majority of the practical applications of the PHD

96 filter involve spatial tracking of moving objects or targets²⁰⁻²³. Herein, the PHD filter is
 97 applied to track dolphin whistles, and in the following the PHD recursion is outlined in that
 98 context.

99 A PHD is a function whose peaks identify the likely positions of the whistle contours. A
 100 whistle with a state \mathbf{x} is more likely to be present in the region where the PHD is large. It
 101 should be noted that the PHD is a density function but not a pdf: a point made apparent
 102 by noting that its integral over the space of its variables is not unity, but is the expected
 103 number of whistles.

The PHD filter is implemented in a recursive manner using prediction and update steps. The goal is to determine the number and the states of the whistle contours at each time k . In the prediction step, the predicted PHD, $v_{k|k-1}$, consists of the information regarding newborn whistles and persistent whistles (whistles surviving from the previous time step, represented by the posterior PHD, v_{k-1}). In the update step the predictions of whistles are refined by incorporating the most recent measurements to obtain the posterior PHD, v_k . The prediction and update steps can be written as^{17,19}:

$$v_{k|k-1}(\mathbf{x}_k) = \gamma_k(\mathbf{x}_k) + p_S \int v_{k-1}(\mathbf{x}_{k-1}) f_{k|k-1}(\mathbf{x}_k | \mathbf{x}_{k-1}) d\mathbf{x}_{k-1}, \quad (3)$$

$$v_k(\mathbf{x}_k) = [1 - p_D] v_{k|k-1}(\mathbf{x}_k) + \sum_{\mathbf{z} \in Z_k} \frac{p_D g_k(\mathbf{z} | \mathbf{x}_k) v_{k|k-1}(\mathbf{x}_k)}{\kappa_k(\mathbf{z}) + p_D \int g_k(\mathbf{z} | \mathbf{x}_k) v_{k|k-1}(\mathbf{x}_k) d\mathbf{x}_k}, \quad (4)$$

104 where $\gamma_k(\mathbf{x}_k)$ denotes the PHD of whistle births between time $k-1$ and k (*i.e.*, the integral
 105 of $\gamma_k(\mathbf{x}_k)$ over a given region gives the expected number of new whistles appearing in that
 106 region at a given time); p_S denotes the probability of survival, that is the *a priori* probability
 107 that a whistle at time $k-1$ will survive until time k ; and $f_{k|k-1}(\mathbf{x}_k | \mathbf{x}_{k-1})$ denotes single-target

108 state transition density (*i.e.*, probability density of a transition to the state \mathbf{x}_k given the
 109 state \mathbf{x}_{k-1}). The probability of detection is denoted by p_D and it represents the probability
 110 that a measurement will be detected from a whistle, $\kappa_k(\mathbf{z})$ denotes the PHD of clutter, and
 111 $g_k(\mathbf{z}|\mathbf{x}_k)$ denotes the single-target measurement likelihood function (*i.e.*, a likelihood that a
 112 measurement \mathbf{z} was generated by a whistle with a state \mathbf{x}_k). Eqs. (3) and (4) have been
 113 adapted to exclude the spawning terms¹⁷, since contour splitting is not typically observed
 114 in dolphin whistles.

115 A closed form solution to (3) - (4) can be obtained assuming the PHD is a mixture of
 116 weighted Gaussian components leading to the so-called Gaussian Mixture PHD (GM-PHD)
 117 filter²⁴. Advantages of this method are that it is straightforward to implement, however it
 118 requires one to assume a linear model for the system and a Gaussian assumption for the noise
 119 processes. Despite this limitation it has been successfully used to track dolphin whistles¹¹.
 120 A more general approximation to (3) - (4) can be achieved with a particle filter, in what is
 121 known as the Sequential Monte Carlo PHD (SMC-PHD) filter^{25,26}, where weighted particles
 122 (random samples) are used to approximate the PHD function. The SMC-PHD filter is a
 123 direct generalization of the approach employed for a single-target particle filter, and particles
 124 are propagated over time using importance sampling and re-sampling strategies^{25,26}. This
 125 implementation of the PHD filter imposes no constraints on the underlying models, and is
 126 the focus of the current work.

127 **B. Contributions**

128 In this paper we propose a complete multi-target frequency tracking scheme to track
129 frequency modulated narrowband signals from audio recordings using a SMC-PHD filter.
130 To achieve this we develop new models for this application, and a particle labeling scheme
131 to allow associations of the tracks between frames. Further, this paper reports the outcome
132 of performance tests on a real-world dataset, and benchmarks the performance against the
133 previously mentioned GM-PHD filter.

134 This paper is organized in the following manner. Section II describes the dataset, the
135 SMC-PHD algorithm for dolphin whistle tracking, along with the developed models and op-
136 timized parameters, as well as the method’s evaluation procedure. Section III contains eval-
137 uation of the proposed methods on the real-world dataset comprising of dolphin recordings.
138 The discussion and the conclusions can be found in Section IV and Section V respectively.

139 **II. METHODS**

140 **A. Data, pre-processing steps and obtaining the measurements**

141 The dataset for evaluation of the filter’s performance was from the 5th Workshop of
142 Detection, Classification, Localization and Density Estimation (DCLDE) conference in 2011,
143 obtained from the MobySound archive (<http://www.mobysound.org>), which has been used
144 in Refs.^{3,7,11,27}. For this work, a subset containing raw recordings and hand-annotated
145 files of whistle contours for six delphinid species (*Delphinus capensis*, *Delphinus delphis*,
146 *Peponocephala electra*, *Stenella longirostris*, *Stenella frontalis*, and *Tursiops truncatus*) was

147 used, and is the same dataset used in Ref.¹¹. The majority of recordings were sampled at
148 192 kHz, but a small portion (15%) of the files had higher sampling rates, and were re-
149 sampled to 192 kHz for consistency. The data collection protocols and study areas are given
150 in Refs.^{7,28}. In this study, the hand-annotated files supplied with the dataset were used
151 as a ground truth data for the filter’s performance evaluation (described in Section IID).
152 Whistles which are 150 ms long and have a Signal to Noise Ratio (SNR) exceeding 10 dB
153 for at least one third of their duration are termed valid⁷ and primarily only valid whistles
154 are used in the following analysis. The raw recordings were used for the filter to track the
155 whistles and, in addition, a part of raw data was set aside as a training set for parameter
156 selection in the SMC-PHD filters. This training set consisted of three one minute duration
157 files chosen randomly: they were recordings of *Delphinus capensis*, *Delphinus delphis*, and
158 *Stenella frontalis* that contained 67, 55 and 63 valid whistles, respectively. This training
159 data was not subsequently used in the performance evaluation.

160 To implement the whistle contour tracking, a set of measurements for each time instance
161 is needed, which is a standard procedure for any multi-target tracking²⁹. The measurement
162 sets, that in our application comprise of the spectral peaks, are obtained using established
163 methods^{3,11}. This pre-processing reduces the background noise and the impact of interfering
164 signals, and is based on a spectrogram using a sliding window of 2048 points (frequency bin
165 width 93.8 Hz) with 50% overlap. The spectral peaks were identified using an 8 dB threshold
166 applied to the normalized spectrogram, converted to a dB scale. Only spectral peaks in the
167 frequency range (2 - 50 kHz) were selected, since that encompassed most dolphin whistles and
168 their harmonics. The precision of the location of the spectral peaks was improved by fitting

169 a quadratic through points surrounding the peak and using the location of the maximum of
170 that fitted quadratic as the refined peak location. These spectral peaks represent the mea-
171 surement set from which the whistle contours were tracked. The pre-processing code used
172 to generate the measurement set from the raw files and the measurement set itself are avail-
173 able at <https://doi.org/10.5258/SOTON/D0316> to facilitate the comparisons with other
174 algorithms that operate on spectral peaks. Moreover, the SMC-PHD filter implementation
175 is available at https://github.com/PinaGruden/SMCPHD_whistle_contour_tracking.

176 **B. Sequential Monte Carlo PHD (SMC-PHD) filter for the whistle contour de-** 177 **tection**

178 The SMC-PHD filter^{25,26} consists of the basic prediction and update steps seen in the
179 Bayes filter. Following the principles of sequential Monte-Carlo methods (or particle fil-
180 ters) the underlying integrals are solved recursively using point-wise approximations. These
181 methods rely upon a set of particles and their associated weights which are propagated from
182 one time step to the next.

183 The standard formulation of the SMC-PHD filter^{25,26} suffers from two limitations that
184 relate to initiating newborn particles (*i.e.*, target birth) and to estimating the state³⁰. The
185 location where targets are born is typically known *a priori* and a large number of parti-
186 cles are required in that region. Further, state estimates are typically constructed using
187 *ad-hoc* clustering of particles^{25,30}. To make the filter more computationally efficient and
188 increase accuracy, a data driven variation of the SMC-PHD filter has been proposed³⁰⁻³²,
189 which generates new particles based on the measurements. This reduces the number of par-

190 ticles required, and eliminates the need for clustering techniques during the state estimation
 191 process by exploiting properties of the PHD update equation. Further, the state estimates
 192 do not contain identities (*i.e.*, it is not known which state estimate belongs to which target
 193 being tracked), and either particle labeling^{33,34} or an external algorithm^{22,35} can be used to
 194 achieve the temporal association of the estimates and so obtain target tracks.

195 The algorithm description and the pseudo-code of the SMC-PHD filter used for whistle
 196 contour tracking are given in Section II B 1, Alg. 1, the temporal association procedure is
 197 detailed in Section II B 2, and the specific models and parameters are presented in Section
 198 II C.

199 1. The SMC-PHD algorithm

200 The SMC-PHD filter propagates through time the weighted particle system $\mathcal{P}_k \equiv$
 201 $\{w_k^{(i)}, \mathbf{x}_k^{(i)}\}_{1 \leq i \leq N_k}$, which approximates the PHD function, $v_k(\mathbf{x}_k)$ ³². At time k , each whis-
 202 tle contour is represented by a cluster of particles representing the state vectors $\mathbf{x}_k^{(i)}$ and
 203 the corresponding weights $w_k^{(i)}$. At each time step, the filter produces an estimate of the
 204 multi-whistle state, $\hat{\mathbf{X}}_k$, which contains state estimates, $\hat{\mathbf{x}}_k$, of whistles. Further, the sum
 205 of particle weights represents an estimate of the number of whistles¹⁷.

206 The prediction step in the SMC-PHD filter starts with the persistent particles from the
 207 previous time step, along with newborn particles being drawn from a proposal density to
 208 form the predicted particles $\mathbf{x}_{k|k-1}^{(i)}$. The predicted particle weights, $w_{k|k-1}^{(i)}$, are computed
 209 by scaling them using p_S . These estimates are then refined in the update step using the set
 210 of measurements, \mathbf{Z}_k .

211 The update process consists of multiple elements. First, the weights and particles are
 212 partitioned on the basis of the measurement set using the probability that the i -th particle
 213 is associated with the j -th measurement, denoted $P_{i,j}$. To allow for the possibility of a
 214 missed detection an additional category is added to the measurements, and the probability
 215 that the i -th particle is associated with the missed detection is denoted $P_{i,0}$. Partitioning
 216 into the clusters is performed by randomly drawing an index for each particle according
 217 to $P_{i,j}, j = 0, \dots, |\mathbf{Z}_k|$ (where $|\cdot|$ denotes the cardinality). The computation of $P_{i,j}$ and
 218 partitioning are based on Eq. (4), and are detailed in Eq. (50) in Ref.³². This partitioning
 219 creates a cluster of particles $C_{k|k-1}(z)$, one for each measurement, plus one cluster for the
 220 missed detection class $C_{k|k-1}(\emptyset)$. Note that there is the possibility that any cluster could
 221 be empty. The method treats the particles associated with the measurements and missed
 222 detections in different fashions.

223 For non-empty clusters associated with a measurement, all the particle weights in the
 224 cluster are updated according to the second term in Eq. (4), and particles are resampled
 225 through a stratified resampling process³⁶. Then a probability of the cluster existing, p_e , is
 226 computed, by summing all the weights of particles in that cluster. If that probability is
 227 greater than a predefined threshold, η , then the resampled particles within a cluster are
 228 averaged to give a state estimate $\hat{\mathbf{x}}_k$.

229 In the cluster corresponding to missed detections, $C_{k|k-1}(\emptyset)$, the particles and associated
 230 weights have no measurements on which to base the update. Their weights are scaled
 231 by $(1 - p_D)$, *i.e.* probability of missed detection. Only particles whose weights exceed a
 232 threshold, ξ , are kept to reduce the computational burden.

233 Finally, the algorithm takes into account the possibility that a whistle starts at time k ,
 234 *i.e.* a whistle birth. At the end of each iteration a set of N_b particles are drawn, focussed on
 235 regions in the state-space where measurements, not associated with state estimates, denoted
 236 $\mathbf{Z}_{b,k}$, were made. This process is detailed in Section II C 3.

237 2. *Temporal association of the whistle estimates*

238 For each time step, the SMC-PHD filter in Alg. 1 outputs a set of the estimated states
 239 $(\hat{\mathbf{X}}_k)$ that represent whistle contour peaks for that time step. However, these do not have
 240 identities associated with them, *i.e.*, one does not know which estimate in one time step
 241 links to which estimate in the next time step, something that is required to be able to form
 242 continuous whistle tracks (contours).

243 Two broad approaches for temporal association of the estimated states are used in the
 244 literature; one is to use a separate algorithm for the association^{22,35}, the other is to label
 245 the individual particles^{33,34}. The particle labeling approach tends to be computationally
 246 more efficient since it only requires an additional set of labels to be propagated alongside
 247 the weighted particle set, and does not need a separate algorithm. The particle labeling
 248 approach was adopted in this work, and each particle i was assigned a label $T^{(i)}$, which was
 249 propagated through time. Unlike other labeling approaches, which typically require cluster-
 250 ing methods^{33,34}, the procedure employed in this study grouped particles by exploiting the
 251 properties of the PHD update equation³². The procedure is outlined below, with reference
 252 to the relevant steps in Alg. 1.

253 On initialization, each particle is assigned a null label $T^{(i)}$. In subsequent prediction steps
 254 (line 2 Alg. 1), the set of the predicted particle labels remains the same, $\mathbf{T}_{k|k-1} = \mathbf{T}_{k-1}$. In
 255 the update step, after resampling (line 12 Alg. 1), the resampled particles within a given
 256 cluster retain the labels of the particles from which they were derived. After that, the cluster
 257 identity is determined, based on the maximum sum of weights of the particles with the same
 258 label. If a cluster originates from a newborn whistle, then the cluster is assigned a new
 259 identity, with all previously unlabelled particles in this cluster being assigned the new label.

260 The identity of the state estimate, is determined from the identity of the cluster from
 261 which the state estimate was derived (line 16 Alg. 1).

262 Each newborn particle is assigned the label $T_{b,k}^{(i)} = 0$ (line 27 Alg. 1). The labels are added
 263 to the labels of the persistent particles and are predicted and updated together in the next
 264 time step.

265 The individual whistles are then tracked from the estimated states based on their iden-
 266 tities. So that all the states with the same identity are linked together into a continuous
 267 whistle contour (track). Finally, the duration of the track is examined and if it falls below a
 268 threshold then it is rejected^{3,7}. This condition is called the track length criterion and various
 269 values for the threshold were investigated ranging from 53 to 150 ms (10 to 28 time steps in
 270 the spectrogram).

271 If there is more than one state estimate with the same identity at a given time, this
 272 conflict is resolved in the following way: in the case that this is a new whistle (no previous
 273 state estimates had this identity), the mean of the states is taken and it becomes the state
 274 estimate for that identity. If this is not a new whistle, *i.e.*, there were previously some

275 state estimates with this identity, the last state estimate with the same identity is projected
 276 forward to the current time step k , using the system function, Eq. (1). The Mahalanobis
 277 distance between the predicted states and those with the conflicting identities is computed.
 278 It is the state closest to the prediction which is retained and the other states with conflicting
 279 identities are discarded.

280 C. Models and parameters for the whistle SMC-PHD

281 As stated in Section IA, the state vectors chosen for this study consist of frequency f
 282 [Hz] and chirp rate α (rate of change of frequency, \dot{f} [Hz/s]), so that $\mathbf{x}_k = [f, \alpha]^t$, where $[\cdot]^t$
 283 denotes the transpose^{10,11}. The following subsections describe the models and parameters
 284 for the SMC-PHD filter employed in this study.

285 1. Measurement model

286 It is assumed that the only measurement available is the frequency and that model for
 287 the noise is additive. Accordingly, the measurement model, (2), can be simplified to¹¹:

$$z_k = \mathbf{H} \mathbf{x}_k + \eta_k = [1, 0] \mathbf{x}_k + \eta_k, \quad (5)$$

288 where z_k [Hz] denotes the available measurements. The measurement noise, η_k [Hz], is
 289 assumed to be Gaussian white noise with a variance R . The value of R is set to the variance
 290 of a uniform random variable covering a single frequency bin¹¹, and is thus dependent on
 291 the width of the frequency bin. Specifically, $R = 732 \text{ Hz}^2$ for the parameters used in this
 292 study.

293 From Eq. (5) the measurement likelihood function, $g_k(z|\mathbf{x}_k)$, is determined to be
 294 $g_k(z|\mathbf{x}_k) = \mathcal{N}(z; \mathbf{H}\mathbf{x}_k, R)$, where $\mathcal{N}(\cdot; m, \Sigma)$ denotes a Gaussian density with a mean m
 295 and a covariance Σ .

296 2. System model

297 The system (motion) model describes how the state develops with respect to time. In the
 298 case of whistle contour tracking, the motion model should contain information on how the
 299 frequency and chirp component of a given whistle evolve with time. The choice of the motion
 300 model can be crucial for the performance of the tracking algorithms, and many applications,
 301 such as surveillance tracking, may have a good understanding of the underlying dynamics³⁷.
 302 However, this is not the case for the frequency evolution of whistle contours.

303 In this work two different motion models are explored. The first model used is a lin-
 304 ear model described in Ref.¹¹, similar to the “nearly-constant-velocity models” in Ref.³⁷,
 305 specifically:

$$\mathbf{x}_k = \mathbf{F}\mathbf{x}_{k-1} + \mathbf{n}_k = \begin{bmatrix} 1 & \Delta \\ 0 & 1 \end{bmatrix} \mathbf{x}_{k-1} + \mathbf{n}_k, \quad (6)$$

306 where Δ denotes the time interval between spectral windows (here 0.0053 s, based on $f_s =$
 307 192 kHz, window length = 2048, 50% overlap). The system noise, \mathbf{n}_k , in this model is
 308 Gaussian white noise with a covariance matrix \mathbf{Q} .

309 While the SMC-PHD allows for non-linear models the linear model in Eq. (6) provides a
 310 baseline for comparisons, since it was successfully applied to track dolphin whistles¹¹ albeit
 311 with a different PHD filter.

312 The covariance matrix, \mathbf{Q} , is assumed to be diagonal, which is equivalent to assum-
 313 ing that the noise processes acting on the frequency and chirp rate are uncorrelated.
 314 In which case there are two degrees of freedom when selecting \mathbf{Q} , namely the two vari-
 315 ances, σ_f^2 and σ_α^2 , representing the noise processes driving the frequency and chirp rate,
 316 respectively. These were determined experimentally by running the SMC-PHD filter on
 317 the training data and selecting the value that gave the best performance (see Section
 318 IID for performance metrics description). A range of values were tested, specifically,
 319 $\{\sigma_f^2, \sigma_\alpha^2\} = (10, 10^2), (10, 10^3), (10^2, 10^3), (10^2, 10^4), (10^3, 10^4), (10^3, 10^6)$. The best perfor-
 320 mance was achieved for the pair $(10^2, 10^4)$.

321 A second motion model was developed based on training a neural network to learn the
 322 temporal relationships defining whistle evolution. The training is based on the set of hand-
 323 annotated data and the idea is similar to the one used in video tracking³⁸, where a set of
 324 hand-annotated traffic trajectories are used to learn how the objects in the scene typically
 325 move, and thus construct a prior to help predict the vehicle motion.

326 The neural network structure adopted here is that of a Radial Basis Function (RBF)
 327 network³⁹. This form of network has the advantage of a comparatively simple structure, but
 328 retains a good ability to generalize. For our application, the RBF can be expressed as:

$$\mathbf{x}_k = \sum_{j=0}^M w_j \phi_j(\mathbf{x}_{k-1}; \mathbf{c}_j, \mathbf{Q}_j) + \mathbf{n}_k, \quad (7)$$

329 where ϕ_j is the set of $M + 1$ basis functions. Herein we use Gaussian functions which are
 330 parametrised by \mathbf{c}_j and \mathbf{Q}_j , these control the location and width of the basis functions respec-
 331 tively (these are closely related to the mean and covariance matrix of a multi-dimensional
 332 Gaussian probability density function). A diagonal form for the matrices \mathbf{Q}_j allows one to
 333 reduce the dimensionality of the model, so limits the amount of training data necessary, but
 334 does restrict the ability of a network of a given size to generalise. The basis function for
 335 $j = 0$, is included as a special case and represents a bias term, realised by fixing $\phi_0 = 1$.

336 In this study, the training data comprised 13,688 data points from 185 whistles. The
 337 hand annotations only measure the frequencies in whistles. To train the network to learn
 338 the full state model it is necessary to know the the chirp rates as well as the frequencies. The
 339 chirp rates were estimated from the hand annotations using a one point backward difference
 340 formula to approximate the frequency derivative (chirp rate).

341 For a given network size, M , the network is trained by first using a k-means clustering
 342 algorithm⁴⁰ to determine a suitable set of centres, \mathbf{c}_j . The \mathbf{Q} matrices are then determined
 343 on the basis of the Euclidean distances between those centres. The final step is to compute
 344 the weights w_j , which only requires the solution of a linear system of equations³⁹. To
 345 select a suitable value for M a cross-validation process is used. This cross-validation process
 346 randomly sub-divided the training data into thirds. Two thirds were used for training during
 347 validation and one third for cross-validation. The network was trained with various choices
 348 of M and the process repeated 10 times for different sub-divisions of the training data with
 349 the results of the 10 repeats averaged. The value of M yielding the smallest mean squared
 350 error was chosen, in this case $M = 60$.

351 As in the linear case the noise \mathbf{n}_k was assumed to be Gaussian, white and uncorrelated
 352 between the two state variables. Based on the statistics of the residuals the noise variances
 353 $\{\sigma_f^2, \sigma_\alpha^2\}$ were (39 Hz², 7326 Hz²/s²).

354 3. *Birth model*

355 The birth model defines where in the state space new whistles are likely to appear, and
 356 how many appear at each time step, as characterized by the birth PHD, $\gamma_k(\cdot)$. If a whistle
 357 appears in a region that is not covered by the birth PHD then the filter may fail to track
 358 it³¹. There are two main challenges associated with determining a suitable birth model, one
 359 is to determine the birth region (*i.e.*, where do new whistles appear) and the other is to
 360 determine the birth magnitude (*i.e.*, how many new whistles appear at each step). Since
 361 the birth PHD in the SMC-PHD filter is represented by a cloud of weighted particles, these
 362 challenges translate into determining the regions in state space from which the new particles
 363 are initiated, and determining their weights.

364 *a. The birth region.* In many other tracking applications the birth region is assumed
 365 to be a single point or uniform across a region of state space^{33,41}. An alternative is to use a
 366 data-driven approach, which is effective when the birth region is not known in advance³¹. In
 367 this approach every measurement initializes newborn targets, with gating techniques used to
 368 divide the measurements into those originating from persistent targets and those originating
 369 from newborn targets²¹.

370 The efficiency of the data driven approach arises because it only introduces particles close
 371 to measurements where the likelihood of a new target appearing is high. This approach is

372 extended and developed further in the current work. The algorithm used here partitions the
 373 particles based on measurements and only the clusters with sufficient weighting are used to
 374 estimate states of persistent whistles (lines 15-18 Alg. 1). This defines the set of persis-
 375 tent whistles and the measurements associated with them. Any remaining measurements,
 376 denoted $\mathbf{Z}_{b,k}$, are then considered when generating newborn particles, denoted $\mathbf{x}_{b,k}$.

377 The process by which newborn particles are generated is as follows. For every $z \in \mathbf{Z}_{b,k}$
 378 a fixed number of particles, N_b , is drawn. The frequency and chirp rate components of the
 379 state are drawn independently from different distributions. For the frequency element:

$$\{x_{b,k}^{(i)}\}_f \sim \mathcal{N}(x; z^{(j)}, R), \quad (8)$$

380 Whereas for the chirp rate there is no direct measurement on which to base the initial
 381 state estimate. To overcome this, the distribution of chirp rates at the start of a whistle
 382 was approximated using a Gaussian Mixture Model (GMM)⁴². The GMM was fitted to the
 383 distribution of chirp rates measured in the hand annotated training data set. The starting
 384 chirp rates for the annotated dataset were computed based on the difference between the
 385 first two frequency samples on a whistle. When fitting the GMM, model order was selected
 386 on the basis of the Bayesian Information Criterion⁴², leading to a choice of a mixture of
 387 three Gaussians. Formally:

$$\{x_{b,k}^{(i)}\}_\alpha \sim \sum_{n=1}^3 a_n \mathcal{N}(x; \mu_n, \sigma_{\alpha,n}^2), \quad (9)$$

388 where a_n , μ_n and $\sigma_{\alpha,n}$ are the weights, means and variances of the GMM respectively. For
 389 our dataset these parameters were $\mathbf{a} = [0.28, 0.02, 0.71]$, $\boldsymbol{\mu} = [1190, -113887, 12999]$ Hz/s,
 390 and $\boldsymbol{\sigma}_\alpha^2 = [9.74, 32.6, 1180] \times 10^6$ Hz²/s².

391 *b. The birth magnitude.* The birth magnitude, ν_b , is the expected number of object
 392 births at a given time³², and is commonly chosen in an *ad-hoc* manner or based on *a priori*
 393 knowledge on the expected number of newborn objects^{24,31}.

394 In this study an alternative approach of computing ν_b adaptively was investigated, based
 395 on the idea that not all measurements are equally likely to generate newborn whistles.
 396 For this purpose, a distribution of the start frequencies of the whistles (the first frequency
 397 in each whistle contour) from the training data (see Section II A) was first computed. The
 398 start frequencies of the whistles in the training data had a skewed, non-Gaussian distribution
 399 (Jarque-Bera test⁴³, $p = 0.001$ at 5% significance level), with the majority of whistles starting
 400 between 8 and 13 kHz. The start frequencies were fitted to a log-normal distribution,
 401 $p_{\text{start}}(z)$, with the log of the start frequencies having a mean $\mu = 9.4$ and standard deviation
 402 $\sigma = 0.4$. The weight of the i^{th} newborn particle, $w_{b,k}^{(i)}$, is then:

$$w_{b,k}^{(i)} \propto \frac{p_{\text{start}}(z)}{N_b}, \quad (10)$$

403 where N_b denotes the number of particles per newborn whistle. Thus the weights of the
 404 newborn particles reflects the *a priori* likelihood of a whistle starting at that frequency.

405 **4. Other parameters**

406 Beside the models discussed in the preceding sections, there are additional parameters
 407 required to implement the SMC-PHD filter. Some parameters are determined based on
 408 properties of the dataset so the choices here match those used in Ref.¹¹. The constant
 409 defining the PHD of clutter (κ_k) was chosen to be uniform across the frequency band of
 410 interest. The observed mean number of false spectral peaks (clutter) per time step, r , was
 411 set to 10, leading to $\kappa_k = r/48000 = 0.0002$. Further, the probability of a whistle surviving
 412 from one time step to the next (p_S) was set to $p_S = 0.994$ based on the mean length of
 413 whistles observed in the training data.

414 The particle elimination threshold (ξ) prevents the number of particles increasing with-
 415 out bounds. It needs to be chosen in a way that the particles on the undetected persistent
 416 whistles, that are collected in a cluster $C_{k|k-1}(\emptyset)$, are not eliminated. Following recommen-
 417 dations elsewhere⁴⁴, this study used $\xi = 100(1 - p_D)/M$ where M denotes the number of
 418 particles in $C_{k|k-1}(\emptyset)$.

419 The remaining parameters were optimized by running the SMC-PHD filter on the training
 420 data and choosing the value that resulted in the best performance (defined in Section IID).
 421 The parameters evaluated in this way are: the probability of detection (p_D), number of
 422 particles per persistent (M_p) and newborn (N_b) whistle, and the threshold used in the state
 423 estimation (η). The values used are summarized in Table I.

424 D. Performance evaluation

425 To evaluate the SMC-PHD filter’s performance, the outputs of the algorithms (which con-
 426 sist of time against frequency peaks for each whistle) were compared to the hand-annotated
 427 ground truth data. Only valid whistles (see Section II A) were expected to be detected.
 428 A detected whistle was considered a match (true positive) to a ground truth whistle if its
 429 timing overlapped with the ground truth whistle and if the mean difference between the
 430 detected whistle path and ground truth whistle path did not exceed 3 frequency bins (281
 431 Hz). If the detected whistle exceeded that criteria, it was considered a false positive. It
 432 should be noted that detected whistles were matched to ground truth whistles regardless
 433 of whether the ground truth whistles were considered valid. However, only the whistles
 434 that matched valid ground truth whistles were considered in the evaluation metrics that de-
 435 scribe the quality and quantity of matches⁷. Additionally, whilst the algorithm searched for
 436 whistles between 2 and 50 kHz, the hand-annotations were only applied to the frequencies
 437 between 4.5 and 50 kHz, therefore any detected whistle that had over 40% of its contour
 438 below 4.5 kHz was not taken into account in the evaluation process¹¹.

439 The performance was measured in terms of recall, precision, fragmentation, mean de-
 440 viation and coverage^{3,7,11}. Recall measures the percentage of the valid whistles that are
 441 retrieved, whilst precision measures the percentage of the detections that are correct⁷. A
 442 high precision therefore indicates a low false alarm rate and a high recall indicates high de-
 443 tection efficiency³. For the detected whistles that matched valid ground truth whistles (true
 444 positives), three additional performance metrics were computed that describe the quality

445 of the detections: fragmentation, mean deviation, and coverage. Fragmentation measures
 446 the average number of detections per ground truth whistle, mean deviation measures the
 447 average frequency deviation between the path of ground truth whistle and its corresponding
 448 detection and coverage measures the average percentage of a ground truth whistle that is
 449 matched⁷.

450 To further evaluate the SMC-PHD filter, its performance was benchmarked against the
 451 GM-PHD filter¹¹. The parameters used in the GM-PHD filter were the same as in Ref.¹¹,
 452 namely: $p_S = 0.994$, $p_D = 0.85$, $U = 10$, $T_r = 0.001$, $w_{th} = 0.009$, and $J_{max} = 100$.

453 The sensitivity of the SMC-PHD filter to the input parameter values in Table I was also
 454 investigated. The sensitivity analysis was carried out by drawing 30,000 random samples
 455 for each parameter from their respective parameter ranges. Each of these randomly selected
 456 parameter sets were then used in the SMC-PHD filter and applied to a single representative
 457 5 s long segment of data containing multiple overlapping whistles, echolocation clicks and
 458 echosounder pulses. Note that, since the value of the parameter ξ changes during the
 459 recursion automatically, it was not included in the sensitivity analysis. To evaluate the
 460 performance of each parameter set, the $F1$ score was computed, which is a harmonic mean
 461 of the precision and recall, and reaches its best value at 100:

$$F1 = \frac{2 \times Recall \times Precision}{Recall + Precision}. \quad (11)$$

462 Afterwards, a pseudo-marginal distribution of each parameter was obtained by creating
 463 equally spaced bins across a given parameter range and computing an average performance
 464 in each bin. A pronounced peak in the pseudo-marginal distribution indicates the filter is

465 sensitive to parameter values around the peak location, while a flat distribution indicates
466 no sensitivity for the specified range.

467 III. RESULTS

468 In total, 9,192 whistles from six different dolphin species were tracked with the SMC-
469 PHD algorithm. The SMC-PHD was considered using two system models one linear and
470 one based on an RBF (see Section II C 2). The performance was investigated across a range
471 of track length criteria (10 - 28 time steps in the spectrogram; 53 - 150 ms).

472 The overall performance results, across all species, are summarized in Fig. 1, and it can be
473 seen that the SMC-PHD filter that utilized the RBF motion model appeared to have better
474 precision (with similar recall), for shorter track lengths, compared to the filter using the
475 linear motion model. For longer track lengths there was a trade-off between precision and
476 recall, with the filter using the RBF motion model having higher precision but lower recall
477 compared to the filter using the linear motion model (Fig. 1). There was also a difference
478 in the coverage, fragmentation, and mean deviation between the two filter types. While
479 the filter that used the RBF motion model had slightly lower coverage and slightly higher
480 fragmentation, it had lower mean deviation from the annotated whistle paths (Fig. 1). In
481 both filters the track length criteria appeared to mainly influence the precision, recall, and
482 fragmentation metrics (Fig. 1). A shorter track length criterion resulted in a higher recall,
483 lower precision, and a higher fragmentation compared to when a longer criterion was applied
484 (Fig. 1).

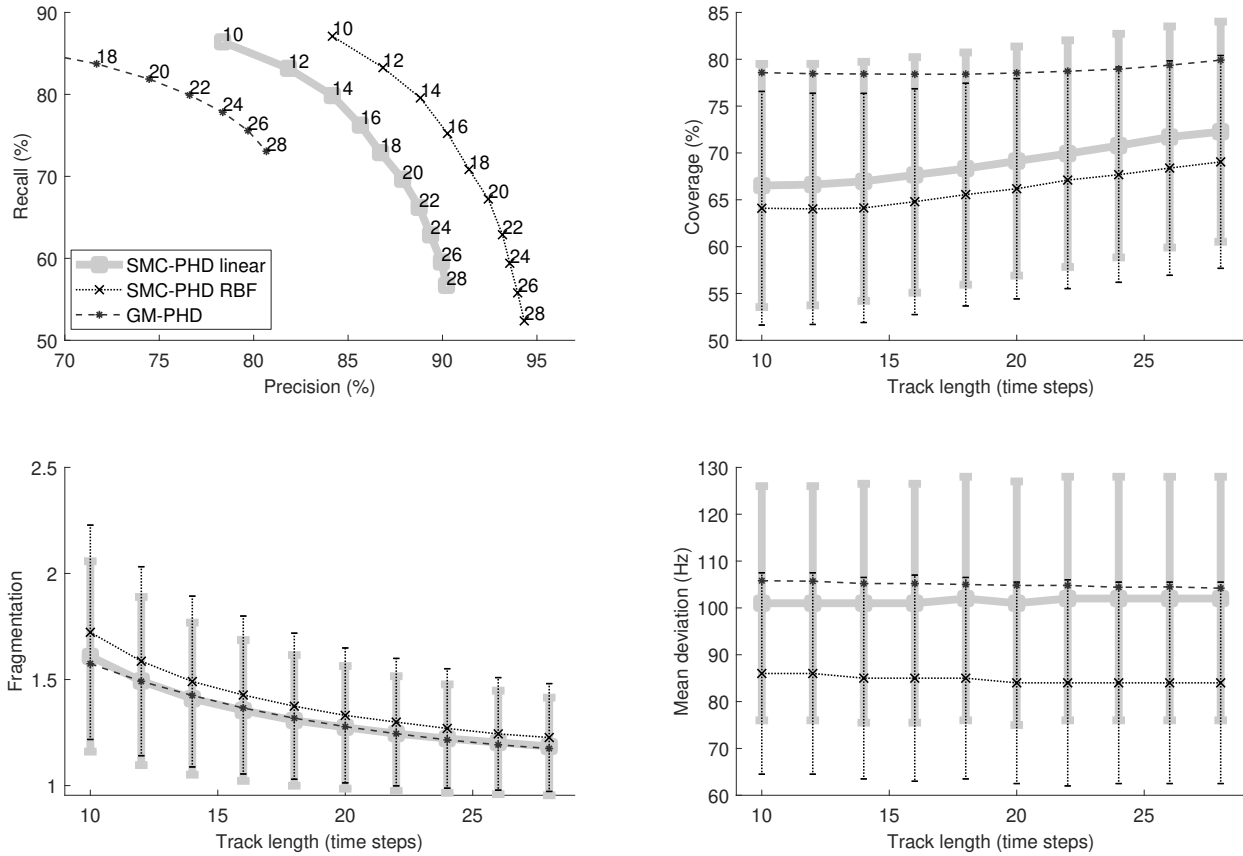


FIG. 1. (Color online) The performance of the SMC-PHD using a linear and RBF motion models across a range of track length criteria (from 10- 28 time steps; 53 - 150 ms). Error bars indicate 1 SD of a given metric. For comparison, the performance of the GM-PHD filter is plotted, but without the corresponding errorbars to preserve the figure's clarity. The performance was computed across all ground truth whistles that met the criteria.

485 Compared to the GM-PHD filter, both SMC-PHD versions had a better precision, but at
 486 the cost of a lower recall when longer track lengths were considered (Fig. 1). The GM-PHD
 487 recall and precision values for shorter track lengths are not displayed in Fig.1, to preserve
 488 the figure's clarity, but they change smoothly, reaching recall of 90% and precision of 40%
 489 for track length 10. This is a comparable recall to both SMC-PHD filters, but at much

490 lower precision. Both SMC-PHD filter versions had a smaller coverage of the individual
491 whistles compared to the GM-PHD filter, but the whistles were tracked more accurately,
492 with a smaller mean deviation between the detection and the ground truth whistle (Fig. 1).
493 All filters had a similar fragmentation rate (Fig. 1).

494 In terms of computational speed, both versions of the SMC-PHD algorithm were capable
495 of tracking the whistles in real time. For example, a two minute file sampled at 192 kHz,
496 containing 795 hand-annotated whistles, took 92.5 s and 117.5 s to be processed with the
497 SMC-PHD filter with linear motion model and SMC-PHD filter with RBF motion model,
498 respectively (implemented in MATLAB, Release R2016b, on a Mac, Os X, processor 2.7
499 GHz and 8 GB RAM).

500 An example of tracking by both versions of the SMC-PHD filter using a short and a long
501 track lengths is shown in Fig. 2. It can be seen that the measurements, from which the filters
502 tracked the whistles, contained a large amount of clutter, *i.e.*, measurements not associated
503 with the whistles (Fig. 2, B). In agreement with the performance results in Fig. 1, it was seen
504 that the SMC-PHD that used a linear motion model produced more false positive detections,
505 for example it detected some of the echosounder pulses when using a track length of 10 time
506 steps (Fig. 2, C) compared to the SMC-PHD that used RBF motion model (Fig. 2, E). It
507 also had higher deviation from the annotated whistle path (for both track lengths) compared
508 to the SMC-PHD filter that used RBF motion model (Fig. 2). However, in some whistles
509 better coverage was achieved (a higher proportion of a given whistle was detected) compared
510 to the SMC-PHD filter that used RBF motion model (Fig. 2). Moreover, in both filters the

511 longer track length criteria resulted in fewer false positives compared to when shorter track
 512 length criteria were used (Fig. 2).

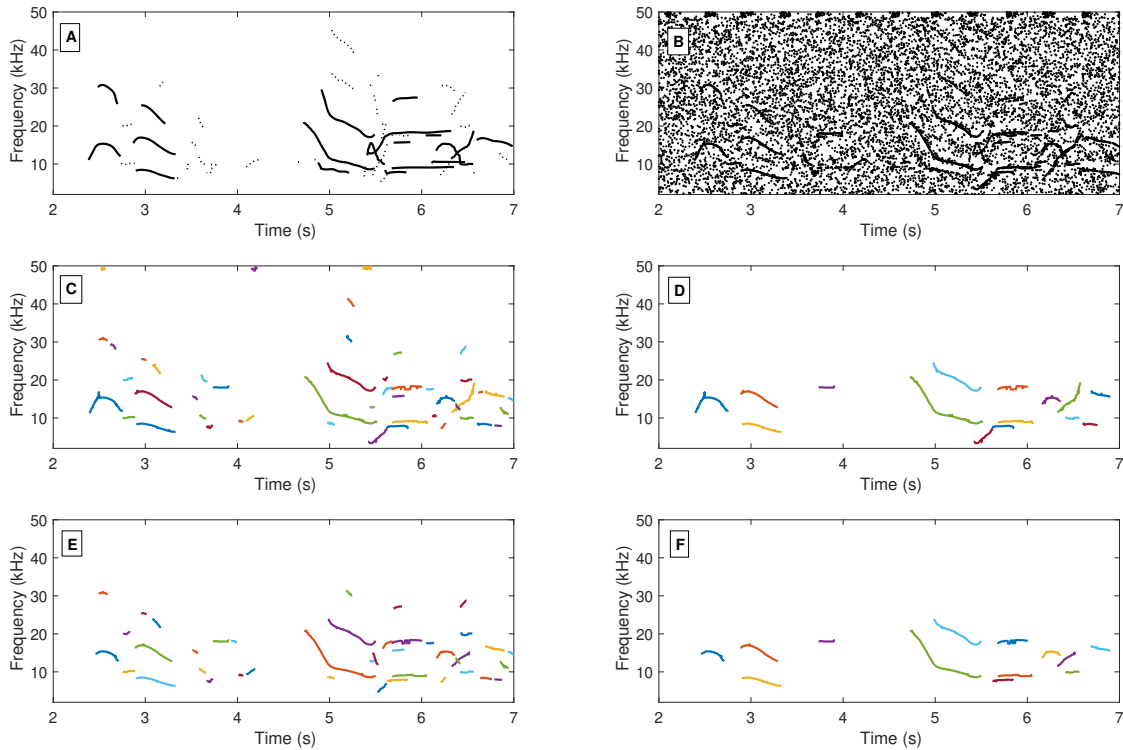


FIG. 2. (Color online) An example of the whistle tracking scenario. (A) Hand-annotated data (solid lines denote valid whistles, dashed lines not-valid whistles; for definition see Section IID). (B) Measurements (spectral peaks) - inputs for the SMC-PHD filters. Extracted whistles with the SMC-PHD filter that utilised linear motion model for track length 10 (C) and 28 (D) time steps. Extracted whistles with the SMC-PHD filter that utilised RBF motion model for track length 10 (E) and 28 (F) time steps.

513 In order to investigate the sensitivity of the SMC-PHD filter to the values of the input
 514 parameters, the best-performing version (using the RBF motion model and track length
 515 criteria of 10 time steps) was evaluated on the example shown in Fig. 2. The filter appeared

516 to be insensitive to small deviations from the values in Table I, as seen by the similar per-
517 formance for bins adjacent to these values in Fig. 3 (indicated with an “x”). However, large
518 changes in some parameters (p_S , p_D , r , η) lead to significant drops in average performance
519 and increase in performance variance. Changes in M_p and N_b did not appear to influence the
520 average F1, but the performance dropped for $M_p < 15$. These results are representative of
521 the behaviour of both SMC-PHD versions, but are not shown here due to space constraints.

522 It should be noted that the performance distribution in each bin of a given parameter
523 in Fig. 3 represents multiple random draws of all the other parameters and therefore is
524 not optimized. As such, the average F1 is lower than the performance obtained with the
525 optimized parameter values in Table I, which result in F1= 87.6 and F1= 85 for the SMC-
526 PHD filter that uses RBF and linear motion models, respectively.

527 The false positive detections were also examined in more detail. These were mainly due
528 to interference from an echosounder and burst pulses, which can display a tonal quality in
529 a spectrogram with the resolution chosen here, see Fig. 4.

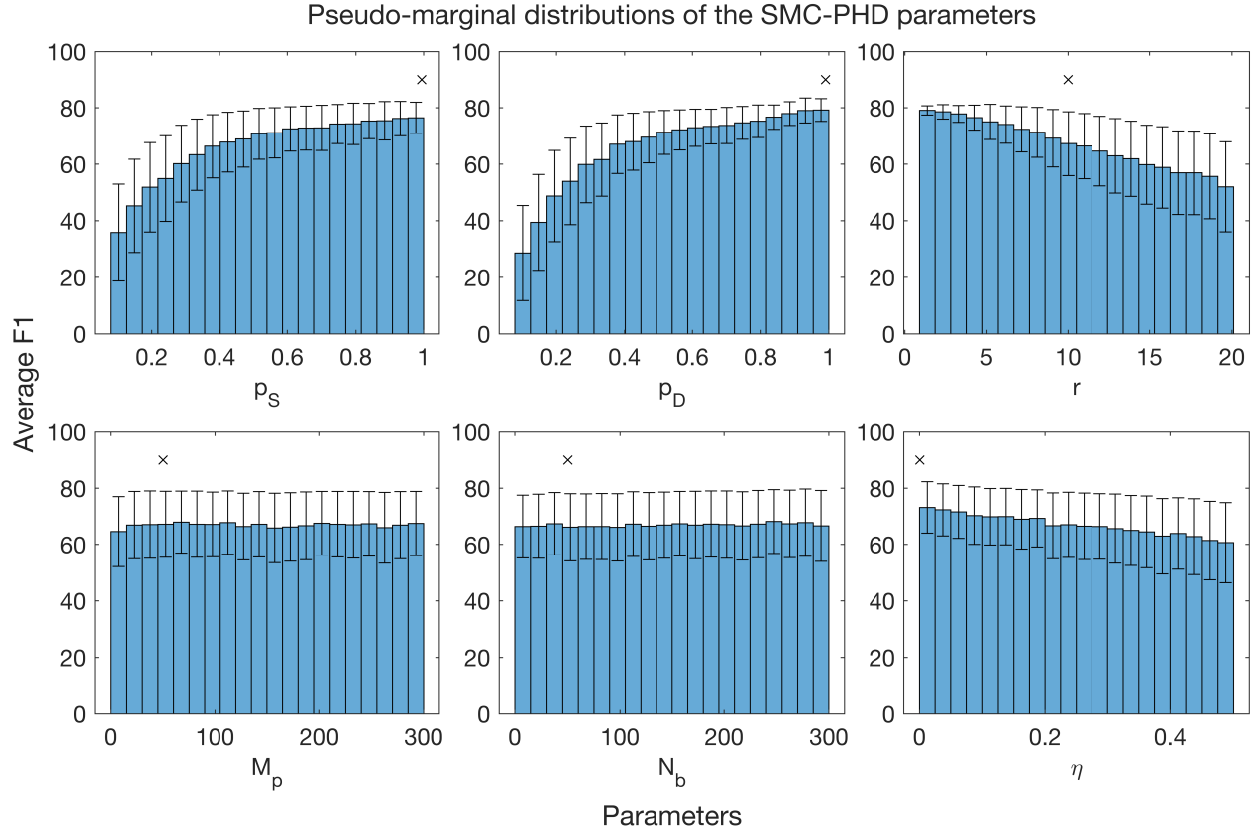


FIG. 3. (Color online) Pseudo-marginal distributions of the SMC-PHD parameters listed in Table I. Average F1 score per bin is shown, with the error bars indicating 1 SD. The values of the parameters that were used in Table I are denoted by “x”.

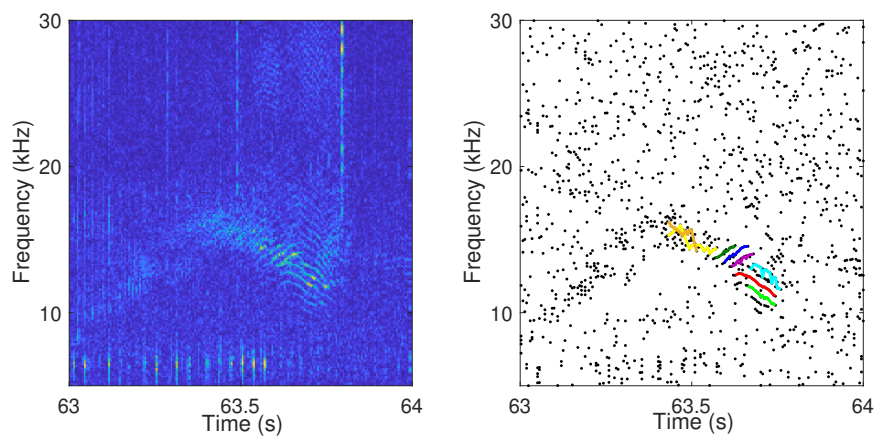


FIG. 4. (Color online) An example of the false positive detection of a burst pulse. A spectrogram of raw data is shown (left) and the measurements (dots) with detected false positives with the SMC-PHD filter (lines) are shown (right).

530 **IV. DISCUSSION**

531 The use of the RFS-based filters is a new approach in the field of the bioacoustics. This
532 paper presents the first attempt to adapt these techniques, specifically the SMC-PHD filter,
533 for the purpose of the frequency tracking of narrowband, frequency modulated signals from
534 the underwater recordings. These methods provide a flexible framework for multi-target
535 tracking, since they impose no restrictions on the form of the character of the underlying
536 motion and measurement models nor do they assume a form for the various noise pro-
537 cesses. The underlying models and parameters for this application are developed and the
538 two proposed schemes are tested on a real world dataset comprising of dolphin whistles. The
539 results showed the proposed filters are able to simultaneously extract multiple whistles from
540 complex acoustic environments; they are able to track the whistles from highly cluttered
541 measurements, through crossings with other whistles and points of missed detections; and
542 they are suited for real time implementation. While the proposed filter implementation in
543 Matlab was able to run in real-time on a typical desktop/laptop computer, it is assumed
544 that significant gains in processing speed can be obtained by carefully implementing the
545 filter in a more optimized programming language.

546 The proposed filters appear to generalize well. While the training data for the filters'
547 parameters and models in Section II C consisted of only three delphinid species, the eval-
548 uated performance returned good results for all six delphinid species in the dataset (Fig.
549 1). Moreover, the recordings used for evaluation contained different whistle types, different
550 noise conditions and amount of interfering signals, and were obtained with different record-

551 ing equipment. This gives additional evidence that the method generalizes for different
552 delphinid species, noise and recording conditions. However, the optimal values for param-
553 eters R , κ_k , and p_D are inherently linked to the choice of pre-processing parameters used to
554 obtain the measurements. For example, the measurement noise variance R depends on the
555 frequency bin resolution used in the spectrogram computation, and the clutter PHD κ_k and
556 the probability of detection p_D depend on the spectrogram amplitude threshold.

557 The use of the SMC-PHD filter requires the development of specific models and param-
558 eters that govern the recursion. The sensitivity of the filter to the input parameter values
559 was evaluated on a representative example that contained multiple overlapping whistles and
560 noise sources, and that was not part of the training data. The filter appeared to be robust
561 to small changes from the trained parameter values. However, large changes in parameters
562 that influence the particle weights (p_S, p_D, r, η) led to a significant drop in performance and
563 increased the variance in the performance results. The value of the parameters that control
564 the number of particles per persistent and newborn whistles, M_p and N_b , did not appear
565 to have a significant influence on the F1 performance. It should be noted, however, that
566 increased number of particles will affect the computational speed of the algorithm, with
567 larger number of particles slowing down the recursion. While the parameter values used
568 in this study appear to give a good performance, there always remains some potential of
569 performance improvement through the selection of a better parameter set. One alternative
570 approach is to modify the filter so that it can adaptively adjust the parameter values during
571 processing⁴⁵.

572 The proposed versions of the SMC-PHD filter were benchmarked against each other and
573 against a different approximation to the PHD filter, the GM-PHD filter¹¹. Both versions
574 of the SMC-PHD filter appeared to have better precision and similar recall compared to
575 the GM-PHD for short track length criterion, but at the cost of lower recall when longer
576 track lengths are considered. Both versions of the SMC-PHD filter tracked whistles more
577 accurately, with smaller mean deviation from the annotated whistle path, but at the cost of
578 having smaller coverage of individual whistles compared to the GM-PHD filter.

579 The performance of the filters depends on the underlying models. A linear motion model
580 describing the evolution of whistle contours was used in the GM-PHD filter and in one
581 of the SMC-PHD filter versions. However, since the true motion model is unknown, it is
582 advantageous to consider learning it from data rather than arbitrarily adopting a linear
583 model. Learning the model from data, as was done for the SMC-PHD filter with RBF
584 motion model, results in a non-linear model and thus requires the use of the SMC-PHD
585 filter. It was seen that the precision of the filter with the non-linear RBF model was better,
586 and this filter tracked individual whistles more accurately (with less deviation) compared
587 to the two filters using linear models. The trade-off was a smaller coverage of individual
588 whistles and slightly higher fragmentation compared to filters using linear models. It should
589 be noted that the non-linear model employed here was trained on a relatively small subset
590 of data, and future studies should consider models trained on larger datasets and consider
591 employing non-Gaussian statistics for the noise processes.

592 The performance was measured based on the hand annotated ground truth data, that was
593 subjective, as with all hand annotations, but at the same time reflected on the performance

594 of the filters in the practical scenarios. As such the values of the performance of the filters
595 should be taken as a guide, not an absolute measure of performance. For both system
596 models, there was a general trade-off between the precision and the recall depending on
597 the track length criteria (which specifies the minimum whistle contour length before it is
598 classed as a detection). Shorter track lengths led to better recall but lower precision, since
599 the number of false positive detections are increased. A shorter track length criterion also
600 increased fragmentation in both instances. Depending on the requirements of the study, the
601 track length criteria can be chosen appropriately.

602 To further improve the performance the following could be considered. This study utilized
603 measurements that consisted only of the frequency peaks from a spectrogram, which makes
604 this problem similar to that of bearing-only tracking in other applications. Adding additional
605 information to the measurements, such as the amplitude or the chirp rate, and expanding the
606 measurement model could potentially improve the performance⁴⁶ and should be investigated
607 further.

608 Furthermore, in the present work the particle labeling approach for temporal association
609 was chosen, since it does not add significantly to the computational load of the recursion.
610 With the proposed labeling scheme, the identity conflicts (when multiple estimates were as-
611 signed the same identity at a given time step) were resolved outside the main PHD recursion
612 and the particles from conflicting clusters propagated freely with the same labels. Although
613 not reported here, a different approach was also tested, where the particles associated with
614 the estimate that did not get assigned to the track were renamed (assigned a new identity).
615 However, this did not produce better results. Another approach could be that instead of

616 discarding the remainder of the conflicting estimates (estimates that are not assigned to a
617 given track), these estimates would be compared against other tracks (with different labels)
618 and assigned to different tracks as appropriate.

619 While the proposed filters successfully tracked dolphin whistles, it should be noted that
620 any frequency modulated signals in the measurements would be extracted. On one hand,
621 this can result in some false alarms that lower the precision of the filter. This was seen
622 with the echosounder and burst pulses, which displayed a tonal quality due to the temporal
623 resolution adopted in this study. It may be possible to remove these false alarms in post-
624 processing steps. On the other hand, having the ability to detect burst pulses could be
625 beneficial in certain applications. Moreover, these filters can be adapted for the extraction
626 of baleen whale sounds or other frequency modulated sounds of interest.

627 **V. CONCLUSIONS**

628 This study considered the frequency tracking of dolphin whistle contours in the context
629 of multi-target tracking. This was achieved with the use of the SMC-PHD filter, a practical
630 approximation to the multi-target Bayesian filter. The filter was adapted and extended
631 for the purpose of frequency tracking and specific models were introduced, resulting in two
632 versions of the filter. The proposed SMC-PHD filters successfully tracked a time-varying
633 number of overlapping whistles from highly cluttered measurements in the presence of false
634 alarms and missed detections. The high degree of flexibility provided by these methods,
635 allied to acceptable computational requirements, means that they are well-suited to real-
636 time tracking of narrowband frequency modulated signals.

637 In addition, to facilitate comparisons of different methods, the measurement sets, a
638 list of all raw audio files used in this study, as well as MATLAB implementation of the
639 method for obtaining spectral peak measurements are openly available from the Univer-
640 sity of Southampton repository at <https://doi.org/10.5258/SOTON/D0316>. Moreover, the
641 SMC-PHD filter implementation is available at https://github.com/PinaGruden/SMCPHD_
642 [whistle_contour_tracking](#).

643 ACKNOWLEDGMENTS

644 We would like to thank MobySound archive, DCLDE committee and associated analysts
645 for providing the datasets and hand annotations used to test detectors performances in this
646 study. We would also like to thank Slovene human resources development and scholarship
647 fund (Ad futura) for funding this research.

648 REFERENCES

- 649 ¹T. A. Marques, L. Thomas, J. Ward, N. DiMarzio, and P. L. Tyack, “Estimating cetacean
650 population density using fixed passive acoustic sensors: An example with Blainvilles beaked
651 whales,” *J. Acoust. Soc. Am.* **125**(4), 1982–1994 (2009).
- 652 ²J. N. Oswald, S. Rankin, J. Barlow, and M. O. Lammers, “A tool for real-time acoustic
653 species identification of delphinid whistles,” *J. Acoust. Soc. Am.* **122**(1), 587–595 (2007).
- 654 ³D. Gillespie, M. Caillat, J. Gordon, and P. R. White, “Automatic detection and classifi-
655 cation of odontocete whistles,” *J. Acoust. Soc. Am.* **134**(3), 2427–2437 (2013).

- 656 ⁴P. Gruden, P. R. White, J. N. Oswald, Y. Barkley, S. Cerchio, M. Lammers, and
657 S. Baumann-Pickering, “Differences in oscillatory whistles produced by spinner (*Stenella*
658 *longirostris*) and pantropical spotted (*Stenella attenuata*) dolphins,” *Marine Mammal Sci.*
659 **32**(2), 520–534 (2016).
- 660 ⁵N. J. Quick and V. M. Janik, “Whistle rates of wild bottlenose dolphins (*Tursiops trun-*
661 *catus*): influences of group size and behavior,” *J. Comp. Psychol.* **122**(3), 305–311 (2008).
- 662 ⁶C. R. Weir and S. J. Dolman, “Comparative review of the regional marine mammal miti-
663 gation guidelines implemented during industrial seismic surveys, and guidance towards a
664 worldwide standard,” *J. Int. Wildlife Law Policy* **10**(1), 1–27 (2007).
- 665 ⁷M. A. Roch, T. S. Brandes, B. Patel, Y. Barkley, S. Baumann-Pickering, and M. S.
666 Soldevilla, “Automated extraction of odontocete whistle contours,” *J. Acoust. Soc. Am.*
667 **130**(4), 2212–2223 (2011).
- 668 ⁸A. Mallawaarachchi, S. H. Ong, M. Chitre, and E. Taylor, “Spectrogram denoising and au-
669 tomated extraction of the fundamental frequency variation of dolphin whistles,” *J. Acoust.*
670 *Soc. Am.* **124**(2), 1159–1170 (2008).
- 671 ⁹D. K. Mellinger, S. W. Martin, R. P. Morrissey, L. Thomas, and J. J. Yosco, “A method
672 for detecting whistles, moans, and other frequency contour sounds,” *J. Acoust. Soc. Am.*
673 **129**(6), 4055–4061 (2011).
- 674 ¹⁰P. R. White and M. L. Hadley, “Introduction to particle filters for tracking applications
675 in the passive acoustic monitoring of cetaceans,” *Can. Acoust.* **36**(1), 146–152 (2008).

- 676 ¹¹P. Gruden and P. R. White, “Automated tracking of dolphin whistles using Gaussian
677 mixture probability hypothesis density filters,” *J. Acoust. Soc. Am.* **140**(3), 1981–1991
678 (2016).
- 679 ¹²S. Rankin, F. Archer, J. L. Keating, J. N. Oswald, M. Oswald, A. Curtis, and J. Barlow,
680 “Acoustic classification of dolphins in the california current using whistles, echolocation
681 clicks, and burst pulses,” *Marine Mammal Sci.* **33**(2), 520–540 (2017).
- 682 ¹³M. Caillat, L. Thomas, and D. Gillespie, “The effects of acoustic misclassification on
683 cetacean species abundance estimation,” *J. Acoust. Soc. Am.* **134**(3), 2469–2476 (2013).
- 684 ¹⁴M. S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, “A tutorial on particle filters
685 for online nonlinear/non-Gaussian Bayesian tracking,” *IEEE Trans. Sign. Process.* **50**(2),
686 174–188 (2002).
- 687 ¹⁵S. M. Bozic, *Digital and Kalman filtering* (Edward Arnold Ltd., London, UK, 1979), p.
688 157.
- 689 ¹⁶R. Mahler, “”Statistics 101” for multisensor, multitarget data fusion,” *IEEE Aerosp. Elec-*
690 *tron. Syst. Mag.* **19**(1), 53–64 (2004).
- 691 ¹⁷R. P. Mahler, *Statistical multisource-multitarget information fusion* (Artech House, Inc.,
692 Norwood, MA, USA, 2007), p. 856.
- 693 ¹⁸R. Mahler, “A theoretical foundation for the Stein-Winter ”Probability Hypothesis Density
694 (PHD)” multitarget tracking approach,” in *Proceedings of the 2000 MSS National Sympo-*
695 *sium on Sensor and Data Fusion*, DTIC Document, San Antonio, Texas, US (2000), pp.
696 99–117.

- 697 ¹⁹R. Mahler, “Multitarget Bayes filtering via first-order multitarget moments,” IEEE Trans.
698 Aerosp. Electron. Syst. **39**(4), 1152–1178 (2003).
- 699 ²⁰D. E. Clark, I. Ruiz, Y. Petillot, and J. Bell, “Particle PHD filter multiple target tracking
700 in sonar image,” IEEE Trans. Aerosp. Electron. Syst. **1**(43), 409–416 (2007).
- 701 ²¹Y.-D. Wang, J.-K. Wu, A. A. Kassim, and W. Huang, “Data-driven probability hypothesis
702 density filter for visual tracking,” IEEE Trans. Circuits Syst. Video Technol. **18**(8), 1085–
703 1095 (2008).
- 704 ²²E. Maggio, M. Taj, and A. Cavallaro, “Efficient multitarget visual tracking using random
705 finite sets,” IEEE Trans. Circuits Syst. Video Technol. **18**(8), 1016–1027 (2008).
- 706 ²³T. M. Wood, C. A. Yates, D. A. Wilkinson, and G. Rosser, “Simplified multitarget tracking
707 using the PHD filter for microscopic video data,” IEEE Trans. Circuits Syst. Video Technol.
708 **22**(5), 702–713 (2012).
- 709 ²⁴B.-N. Vo and W.-K. Ma, “The Gaussian mixture probability hypothesis density filter,”
710 IEEE Trans. Sign. Process. **54**(11), 4091–4104 (2006).
- 711 ²⁵B.-N. Vo, S. Singh, and A. Doucet, “Sequential Monte Carlo implementation of the PHD
712 filter for multi-target tracking,” in *Proceedings on International Conference on Information*
713 *Fusion* (2003), pp. 792–799.
- 714 ²⁶T. Zajic and R. P. Mahler, “Particle-systems implementation of the PHD multitarget-
715 tracking filter,” in *Proceedings of SPIE* (2003), Vol. 5096, pp. 291–299.
- 716 ²⁷K. E. Frasier, E. Elizabeth Henderson, H. R. Bassett, and M. A. Roch, “Automated
717 identification and clustering of subunits within delphinid vocalizations,” *Marine Mammal*

718 Science **32**(3), 911–930 (2016).

719 ²⁸S. Baumann-Pickering, S. M. Wiggins, J. A. Hildebrand, M. A. Roch, and H.-U. Schnitz,
720 “Discriminating features of echolocation clicks of melon-headed whales (*Peponocephala*
721 *electra*), bottlenose dolphins (*Tursiops truncatus*), and Gray’s spinner dolphins (*Stenella*
722 *longirostris longirostris*),” J. Acoust. Soc. Am. **128**(4), 2212–2224 (2010).

723 ²⁹B.-N. Vo, M. Mallick, Y. Bar-Shalom, S. Coraluppi, R. Osborne III, R. Mahler, and B.-T.
724 Vo, “Multitarget tracking,” in *Wiley Encyclopedia of Electrical and Electronics Engineer-*
725 *ing* (John Wiley and Sons, Inc., 2015), pp. 1–23.

726 ³⁰B. Ristic, D. Clark, and B.-N. Vo, “Improved SMC implementation of the PHD filter,” in
727 *13th Conference on Information Fusion 2010*, IEEE (2010), pp. 1–8.

728 ³¹B. Ristic, D. Clark, B.-N. Vo, and B.-T. Vo, “Adaptive target birth intensity for PHD and
729 CPHD filters,” IEEE Trans. Aerosp. Electron. Syst. **48**(2), 1656–1668 (2012).

730 ³²B. Ristic, M. Beard, and C. Fantacci, “An overview of particle methods for random finite
731 set models,” Inf. Fusion **31**, 110–126 (2016).

732 ³³K. Panta, B.-N. Vo, and S. Singh, “Improved probability hypothesis density (PHD) filter
733 for multitarget tracking,” in *Third International Conference on Intelligent Sensing and*
734 *Information Processing, 2005. ICISIP 2005.*, IEEE (2005), pp. 213–218.

735 ³⁴D. E. Clark and J. Bell, “Multi-target state estimation and track continuity for the particle
736 PHD filter,” IEEE Trans. Aerosp. Electron. Syst. **43**(4), 1441 – 1453 (2007).

737 ³⁵K. Panta, B.-N. Vo, and S. Singh, “Novel data association schemes for the probability
738 hypothesis density filter,” IEEE Trans. Aerosp. Electron. Syst. **43**(2), 556–570 (2007).

- 739 ³⁶T. Li, M. Bolic, and P. M. Djuric, “Resampling methods for particle filtering: classification,
740 implementation, and strategies,” *IEEE Sign. Process. Mag.* **32**(3), 70–86 (2015).
- 741 ³⁷X. R. Li and V. P. Jilkov, “Survey of maneuvering target tracking. Part I: Dynamic
742 models,” *IEEE Trans. Aerosp. Electron. Syst.* **39**(4), 1333–1364 (2003).
- 743 ³⁸L. Wang, L. Zhang, and Z. Yi, “Trajectory predictor by using recurrent neural networks
744 in visual tracking,” *IEEE Trans. Cybern.* **47**(10), 3172–3183 (2017).
- 745 ³⁹C. M. Bishop, *Neural networks for pattern recognition* (Oxford University Press, Inc., New
746 York, NY, USA, 1995), pp. 164–191.
- 747 ⁴⁰D. Arthur and S. Vassilvitskii, “k-means++: The advantages of careful seeding,” in *Pro-
748 ceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, Society
749 for Industrial and Applied Mathematics (2007), pp. 1027–1035.
- 750 ⁴¹B.-N. Vo, S. Singh, and A. Doucet, “Sequential Monte Carlo methods for multitarget
751 filtering with random finite sets,” *IEEE Trans. Aerosp. Electron. Syst.* **41**(4), 1224–1245
752 (2005).
- 753 ⁴²C. M. Bishop, *Pattern recognition and machine learning* (Springer Science & Business
754 Media, New York, NY, USA, 2006), p. 738.
- 755 ⁴³C. M. Jarque and A. K. Bera, “A test for normality of observations and regression resid-
756 uals,” *Int. Stat. Rev.* **55**(2), 163–172 (1987).
- 757 ⁴⁴B. Ristic, “Efficient update of persistent particles in the SMC-PHD filter,” in *IEEE Inter-
758 national Conference on Acoustics, Speech and Signal Processing (ICASSP), 2015*, IEEE
759 (2015), pp. 4120–4124.

760 ⁴⁵R. P. Mahler, B.-T. Vo, and B.-N. Vo, “CPHD filtering with unknown clutter rate and
761 detection profile,” *IEEE Trans. Sign. Process.* **59**(8), 3497–3513 (2011).

762 ⁴⁶D. Clark, B. Ristic, B.-N. Vo, and B. T. Vo, “Bayesian multi-object filtering with amplitude
763 feature likelihood for unknown object SNR,” *IEEE Trans. Sign. Process.* **58**(1), 26–37
764 (2010).

Algorithm 1 Pseudo-code of the SMC-PHD filter for whistle contour tracking (adapted

 based on Ref.³²)

 1: Input $\mathcal{P}_{k-1} \equiv \{w_{k-1}^{(i)}, \mathbf{x}_{k-1}^{(i)}\}_{1 \leq i \leq N_{k-1}}; \mathbf{Z}_k$

 2: **Step 1 Prediction**

 3: Draw particles from proposal density to obtain $\mathbf{x}_{k|k-1}^{(i)}$ ▷ see Section II C 2

 4: Compute their weights: $w_{k|k-1}^{(i)} = p_S w_{k-1}^{(i)}$

 5: **Step 2 Update, Resampling, State Estimation**

 6: Partition $\{w_{k|k-1}^{(i)}, \mathbf{x}_{k|k-1}^{(i)}\}_{1 \leq i \leq N_{k-1}}$ to form clusters $C_{k|k-1}(z), z \in \mathbf{Z}_k \cup \emptyset$

 7: Initialize $\mathcal{P}_k = \emptyset, \hat{\mathbf{X}}_k = \emptyset$

 8: **for** every $z \in \mathbf{Z}_k$ **do**

 9: **if** $C_{k|k-1}(z) \neq \emptyset$, it consists of M weighted particles $\{w_{k|k-1}^{(m)}, \mathbf{x}_{k|k-1}^{(m)}\}_{1 \leq m \leq M}$ **then**

 10: Update their weights: $\hat{w}_k^{(m)} = \frac{p_D g_k(z | \mathbf{x}_{k|k-1}^{(m)}) w_{k|k-1}^{(m)}}{\kappa_k + p_D \sum_{n=1}^M g_k(z | \mathbf{x}_{k|k-1}^{(n)}) w_{k|k-1}^{(n)}}$

 11: Compute probability of cluster's existence: $p_e(z) = \sum_{m=1}^M \hat{w}_k^{(m)}$

 12: Resample based on \hat{w}_k to generate M_p particles $\mathbf{x}_k^{(l)}, l = 1, \dots, M_p$

 13: Set the resampled particle weights to $w_k^{(l)} = p_e(z) / M_p, l = 1, \dots, M_p$

 14: $\{w_k^{(l)}, \mathbf{x}_k^{(l)}\}_{1 \leq l \leq M_p}$ represent updated cluster $C_k(z)$, and $\mathcal{P}_k = \mathcal{P}_k \cup C_k(z)$

 15: **if** $p_e(z) > \eta$ **then** ▷ η is a threshold determined in Section II C 4

 16: Estimate whistle state $\hat{\mathbf{x}}_k$ from $C_k(z)$: $\hat{\mathbf{x}}_k = 1 / M_p \sum_{l=1}^{M_p} \mathbf{x}_k^{(l)}$

 17: $\hat{\mathbf{X}}_k = \hat{\mathbf{X}}_k \cup \{\hat{\mathbf{x}}_k\}$

 18: **end if**

 19: **end if**

 20: **end for**

21: **for** every pair $(w_{k|k-1}, \mathbf{x}_{k|k-1}) \in C_{k|k-1}(\emptyset)$ **do**

22: **if** $w_{k|k-1} > \xi$ **then** ▷ ξ is a threshold determined in Section II C 4

23: Update weights as: $w_k = (1 - p_D)w_{k|k-1}$

24: And add the weighted particles to \mathcal{P}_k

25: **end if**

26: **end for**

27: **Step 3 Whistle birth**

28: **for** every $z \in \mathbf{Z}_{b,k}$ **do**

29: Generate N_b particles and compute their weights ▷ see Section II C 3

30: Add the newborn weighted particles to \mathcal{P}_k

31: **end for**

32: Output: $\mathcal{P}_k \equiv \{w_k^{(i)}, \mathbf{x}_k^{(i)}\}_{1 \leq i \leq N_k}; \hat{\mathbf{X}}_k$

TABLE I. Summary of the parameters used in the SMC-PHD filter for dolphin whistle tracking. p_S and p_D denote the probabilities of survival and detection respectively; r denotes the average number of clutter measurements per time step; M_p and N_b denote the number of particles per persistent and newborn whistle respectively; η denotes state estimation threshold; ξ denotes particle elimination threshold, where M is the number of particles in cluster $C_{k|k-1}(\emptyset)$.

p_S	p_D	r	M_p	N_b	η	ξ
0.994	0.99	10	50	50	0.0005	1/M

765 **FIGURE CAPTIONS**

766 Fig.1 (Color online) The performance of the SMC-PHD using a linear and RBF motion
 767 models across a range of track length criteria (from 10- 28 time steps; 53 - 150 ms).
 768 The performance is computed across all ground truth whistles that met the criteria
 769 and is not the average of file or species performances. Each error bar indicates one
 770 standard deviation of a given metric.

771 Fig.2 (Color online) An example of the whistle tracking scenario. (A) Hand-annotated data
 772 (solid lines denote valid whistles, dashed lines not-valid whistles; for definition see
 773 Section IID). (B) Measurements (spectral peaks) - inputs for the SMC-PHD filters.
 774 Extracted whistles with the SMC-PHD filter that utilised linear motion model for
 775 track length 10 (C) and 28 (D) time steps. Extracted whistles with the SMC-PHD
 776 filter that utilised RBF motion model for track length 10 (E) and 28 (F) time steps.

777 Fig.3 (Color online) Pseudo-marginal distributions of the SMC-PHD parameters listed in
 778 Table I. Average F1 score per bin is shown, with the error bars indicating 1 SD. The
 779 values of the parameters that were used in Table I are denoted by “x”.

780 Fig.4 (Color online) An example of the false positive detection of a burst pulse. A spec-
 781 trogram of raw data is shown (left) and the measurements (dots) with detected false
 782 positives with the SMC-PHD filter (lines) are shown (right).