

## Efficient abstracting of dive profiles using a broken-stick model

Theoni Photopoulou<sup>1\*</sup>, Philip Lovell<sup>2</sup>, Michael A. Fedak<sup>2</sup>, Len Thomas<sup>3</sup> and Jason Matthiopoulos<sup>4</sup>

<sup>1</sup>Centre for Statistics in Ecology, Environment and Conservation, Department of Statistical Sciences, University of Cape Town, Rondebosch, Cape Town, 7701, South Africa; <sup>2</sup>Sea Mammal Research Unit, Scottish Oceans Institute, University of St Andrews, Scotland KY16 8LB, UK; <sup>3</sup>Centre for Research into Ecological and Environmental Modelling, The Observatory, University of St Andrews, Scotland KY16 9LZ, UK; and <sup>4</sup>Institute of Biodiversity, Animal Health and Comparative Medicine, Graham Kerr Building University of Glasgow, Glasgow, Scotland G12 8QQ, UK

### Summary

**1.** For diving animals, animal-borne sensors are used to collect time–depth information for studying behaviour, ranging patterns and foraging ecology. Often, this information needs to be compressed for storage or transmission. Widely used devices called conductivity-temperature-depth satellite relay data loggers (CTD-SRDLs) sample time and depth at high resolution during a dive and then abstract the time–depth trajectory using a broken-stick model (BSM). This approximation method can summarize efficiently the curvilinear shape of a dive, using a piecewise linear shape with a small, fixed number of vertices, or break points.

**2.** We present the process of abstracting dives using the BSM and quantify its performance, by measuring the uncertainty associated with the profiles it produces. We develop a method for obtaining a confidence zone and an index for the goodness-of-fit (dive zone index, DZI) for abstracted dive profiles. We validate our results with a case study using dives from elephant seals (*Mirounga* spp.). We use generalized additive models (GAMs) to determine whether the DZI can be used as a proxy for an absolute measure of fit and investigate the relationship between the DZI and the dive shape.

**3.** We found a strong correlation between the residual sum of squares (RSS) for the difference between the detailed and abstracted profiles, and the DZI and maximum residual (R<sub>4</sub>), for dives resulting from CTD-SRDLs (69% deviance explained). On its own, the DZI explained a lower percentage of deviance which was variable for abstracted dives with different numbers of break points. We also found evidence for systematic differences in the DZI for different dive shapes (65% deviance explained).

**4.** Although the proportional loss of information in the abstraction of time–depth dive profiles by BSM is high, what remains is sufficient to infer goodness-of-fit of the abstracted profile by reversing the abstraction process. Our results suggest that together the DZI and R<sub>4</sub> can be used as a proxy for the RSS, and we present the method for obtaining these metrics for BSM-abstracted profiles.

**Key-words:** animal telemetry, broken-stick model, CTD-SRDL, dive profile, dive type, dive zone index, elephant seal, data abstraction

### Introduction

#### DATA ABSTRACTION IN ANIMAL TELEMETRY: NEEDS AND CONSEQUENCES

One of the most effective ways to remotely study movement and behaviour in marine animals is to use animal-borne sensors. Satellite-linked and archival animal telemetry devices have developed rapidly, driven by questions about the behaviour and movement of large vertebrates at sea. A range of purpose-built hardware and software is widely available for deployment on animals. Although animal telemetry devices

are able to record information at high temporal and spatial resolution, in many cases, devices cannot be recovered, which means the data they collect have to be transmitted. Additionally, it is seldom possible to transmit all data that are collected during a deployment, because the quantity and resolution of the received telemetry data are constrained by several factors: the desired observation time, battery life of the device, bandwidth of the communication system used to relay data, behaviour of the animal and the software specifications. This means that not all information that is recorded can be recovered. Consequently, data abstraction (defined here as reduction in volume to a simplified representation of the original) is unavoidable for many types of telemetry device.

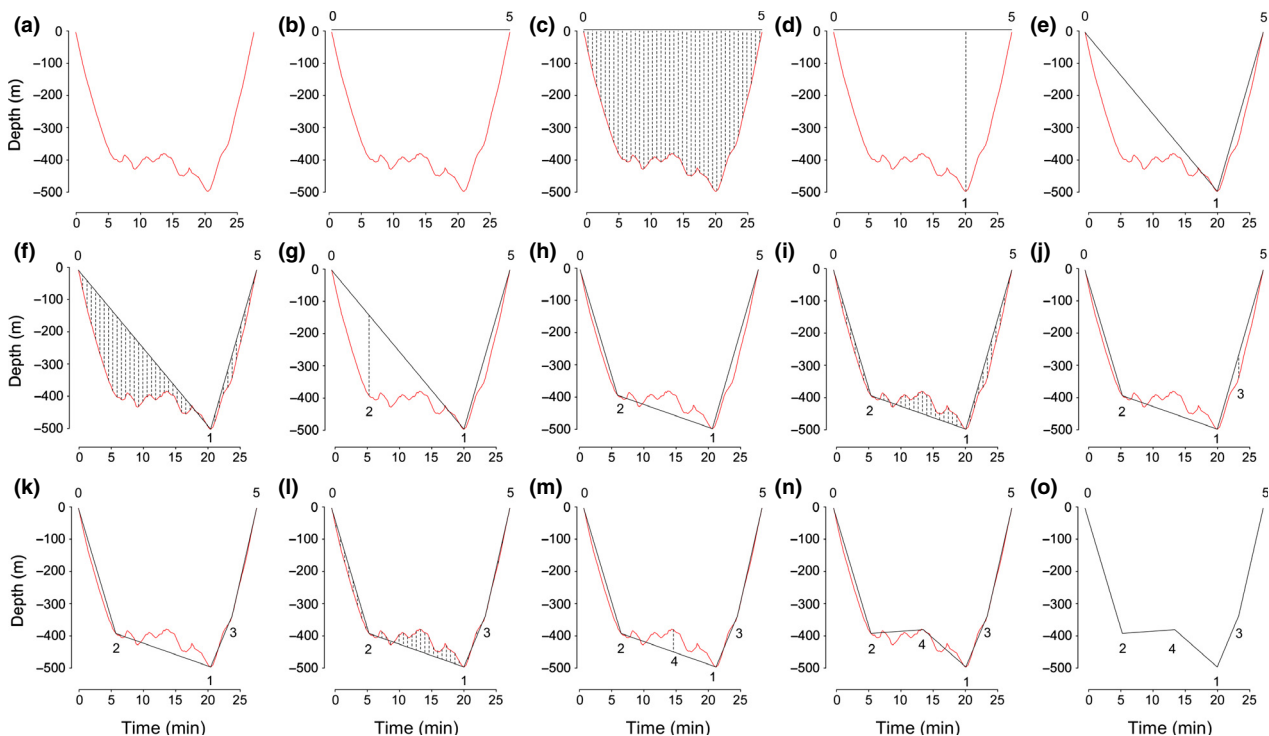
The trade-off between temporal data resolution (i.e. the rate of data sampling and delivery) and the operational lon-

\*Correspondence author. E-mail: theoni.photopoulou@uct.ac.za

geivity of the telemetry device has driven the development of efficient software and memory-saving processing algorithms. One of these is a broken-stick model (BSM), used for the abstraction of two-dimensional dive trajectories (time–depth dive profiles) on-board telemetry devices prior to transmission (Fedak *et al.* 2002). BSMs are change-point models, falling under piecewise linear approximation and are used to identify points of abrupt change in time series. The piecewise linear profile of a time series of depths, generated by an efficient linear abstraction method, should have low average deviation from the detailed dive profile. When processing takes place on a small device with limited power for processing and transmissions, it is advantageous to represent piecewise linear profiles using a fixed and small number of bits of information, particularly when using CLS Argos, where message size is fixed (Argos 2011). The algorithm should also be time efficient, its execution time should scale linearly with dive duration, it should consistently encode biologically relevant information that might enable inferences on dive function. The BSM fulfils these criteria and was adopted as the default dive abstraction algorithm on conductivity-temperature-depth satellite relay data loggers (CTD-SRDLs) in 2006 (pers. comm., Phil Lovell). The predecessor of this model placed break points at locations of maximum inflection in

the detailed profile, and while it performed equally well, processing required more time and energy (Fedak, Lovell & Grant 2001). The BSM was chosen empirically, because it was found to provide a highly satisfactory compromise between the priorities described above. However, its performance has not been formally tested, nor have the consequences of its performance on the biological and ecological conclusions that are drawn in studies using dive profiles abstracted with the BSM.

In ecology, change-point models have a long history (MacArthur 1957; MacArthur & MacArthur 1961) and have seen application in a range of fields, for example, in oceanography to reduce data volume (Rual 1989), to identify edge effects in plant communities (Toms & Lesperance 2003), to locate ontogenetic shifts in southern elephant seal diet using stable isotopes (Authier *et al.* 2012) and in a Bayesian context applied to allometric relationships between tree height and diameter (Beckage *et al.* 2007). We illustrate the use of the BSM for dive abstraction with a time–depth dive profile as an example (Fig. 1). We call *abstracted* those dive profiles that have been processed and reduced in resolution using this algorithm. We call *detailed* those dive profiles that are recorded at regular and frequent time intervals, at the sampling resolution of the device, before they are abstracted.



**Fig. 1.** Stepwise illustration of time–depth dive profile abstraction for a dive from elephant seal 10 943 (see Table 1) using a broken-stick model. With this abstraction regime, highly resolved time–depth measurements are processed to generate abstracted dive profiles made up of  $I + 1$  consecutive line segments that approximate the true, nonlinear time–depth trajectory travelled. This is achieved through  $I (=4, \text{ here})$  iterations of the algorithm. The red lines represent the true time–depth dive path, the solid black lines represent the abstracted dive path and the dashed lines represent the residuals for the abstracted dive path, at each iteration of the algorithm. The numbers represent the order in which points are added to the abstracted profile. Points 0 and 5 mark the beginning and end of the dive. The dashed black lines represent the residuals that are calculated at each time step between the true and abstracted dive.

# ASSESSING UNCERTAINTY IN ABSTRACTED DIVE PROFILES FROM ELEPHANT SEALS (*MIROUNGA* SP.) AS A CASE STUDY

The need for abstraction becomes critical for deployments on wide-ranging marine species, such as seals and turtles, when geographic and temporal data coverage is of interest, and when devices cannot be recovered. For elephant seals and other phocid seals, for example, complete time series of year-round locations and behaviour may be more biologically interesting than detailed information over short periods and more useful for understanding their life histories (McConnell, Chambers & Fedak 1992; Hebblewhite & Haydon 2010). Until now, the uncertainty associated with abstracted profiles has not been quantified, and the implications for ecological studies that use abstracted profiles have not been assessed.

Abstracted profiles are, by construction, information-poor versions of the detailed trajectories, but since the abstraction process is known in the case of CTD-SRDLs, it is possible to reverse the deterministic steps and retrieve some of the information. Historically, once the abstraction was completed, the high-resolution time–depth profile was overwritten, but current tags store all information that they record, and this can be accessed if the tag is retrieved. Here, we show that it is possible to construct a 100% confidence zone around an abstracted profile (i.e. upper and lower depth limits at each time point within which the true depth must lie) and compare the zone for different dives. This confidence zone is hereafter referred to as the *dive zone*, and the relative measure of maximum deviation of the abstracted profile, from the detailed profile, is referred to as the *dive zone index* (DZI).

The consequences of the abstraction regime by BSM are investigated here using detailed and abstracted dive profiles, from northern (*Mirounga angustirostris*) and southern elephant seals (*Mirounga leonina*), as a case study. Elephant seals are large-bodied, long-lived and abundant marine mammals. They spend many months at sea in the open ocean and in coastal or marginal ice zones (McConnell, Chambers & Fedak 1992; Jonker & Bester 1998; Bailleul et al. 2007; Campagna et al. 2007; Biuw et al. 2010). They frequently visit high latitudes for extended periods, diving deeply and almost continually, returning to land twice a year to breed and moult. CTD-SRDLs are regularly used in studies of their movement and diving behaviour and that of other wide-ranging phocid seals.

The characterization of dives into types, based on dive parameters, has been a popular approach to the study of diving behaviour (Hindell, Slip & Burton 1991; Schreer & Testa 1996; Schreer, Kovacs & O'Hara Hines 2001; Baechler 2002). In general, the identification of types or groups of behaviour is useful for comparing behavioural patterns and activity budgets between individuals and in different spatial and temporal contexts and is carried out using a wide range of methods including empirical methods, machine learning and state-space methods, to name a few (e.g. Fauchald & Tveraa 2003; Thums, Bradshaw & Hindell 2008; McKellar et al. 2014). BSM dives are used widely in studies of diving behaviour and physiology (e.g.

McConnell et al. 1999; Bailleul et al. 2007; Biuw et al. 2007, 2010), without considering or accounting for uncertainty in the abstracted dive profiles. Ignoring the uncertainty in BSM-derived time–depth profiles may lead to incorrect inferences if the BSM output has substantial error associated with it, and if dives with different shape characteristics differ systematically in the amount of error associated with them. This computationally expedient method has been thought to perform well at capturing biologically relevant aspects of time–depth dive profiles, but this impression has, to date, remained anecdotal.

## AIMS AND QUESTIONS

In this study, we aim to explain the BSM for dive profile abstraction, provide a method that extracts as much information as possible from abstracted dive data and improve the interpretability of abstracted dive profiles. To do this, we (i) present an overview of the process by which the BSM generates abstracted dive profiles; (ii) assess the performance of the BSM for dive profile representation, by comparing detailed and abstracted time–depth dive profiles from elephant seals, as a case study; (iii) present a three-step method for obtaining, *post hoc*, the depth limits on the detailed dive (i.e. the dive zone) based on its BSM-abstracted profile; (iv) develop an index of goodness-of-fit of abstracted dives (i.e. DZI) and use detailed dive profiles to validate it; and (v) use this index to determine whether there are systematic differences in the amount of error associated with different dive types, following (Hindell, Slip & Burton 1991).

We recast these five aims as research questions: (1) How does the BSM work for dive profile abstraction? How does the representation of the detailed dive change with increasing BSM points? (2) Is the sample of study dives representative? (3) What can we learn from abstracted dives? (4) How is the DZI derived? and (5) Can the DZI be used as a proxy for the residual sum of squares (RSS)? Does the DZI vary systematically between dive types?

## Materials and methods

### HOW DOES THE BSM WORK FOR DIVE PROFILE ABSTRACTION?

The BSM is an iterative process. For time–depth dive profiles, it is based on minimizing the vertical distance (i.e. difference in depth) between the detailed trajectory recorded by the tag and the abstracted dive profile being proposed, at the sampling resolution of the dive (for CTD-SRDLs 4s, 8s, 16s or 32s). We call these vertical distances *residuals* (Fig. 1). The basic principle of the model is that, at each iteration, the residuals are calculated, and the point with the biggest residual is added to the abstracted profile. At the first iteration, which we call step zero, the abstracted profile consists only of the start and endpoints of the dive, forming a straight line at what the tag perceives to be zero depth. This corresponds to a depth buffer at the surface (0–6 m), which is intended to exclude any less interesting shallow undulations from the dive record, because they would compete with regular deep dives for transmission. The distance from this straight line to the detailed profile is measured at each time point, and the point at which the piecewise

**Table 1.** Deployment information and morphometrics for the detailed data sets from one northern and three southern elephant seal (a sample of 60 study dives was taken from the dive record of each seal) and the abstracted data set(s) from 45 southern elephant seals. For 'ct' deployments, length is given as the mean (standard error) of all animals of each sex. For deployments 12 454, 12 453, 12 451, two length measurements were available for each animal, so we present the mean (standard error) of the two measurements. Only one length measurement was made during deployment 10 943

Deployment	Species	Deployment location	Period (UTC)	Sampling regime	Morphometrics	Sex	Dive duration
10 943	<i>Mirounga angustirostris</i> 1 adult male	Año Nuevo CA, USA	23/08/2008 to 16/02/2009	Time and depth every 4s *	Length: 350.0 cm Axial Girth: 295.0 cm	M	24.95 min (0.77)
12 454	<i>Mirounga leonina</i> 1 adult female	Kerguelen Islands	30/10/2012 to 15/02/2013	Time and depth every 4s	Length: 235.5.0 cm (6.5) Weight: 258 kg	F	19.18 min (0.63)
12 453	<i>M. leonina</i> 1 adult female	Kerguelen Islands	20/10/2012 to 11/02/2013	Time and depth every 4s	Length: 271.5 cm (6.5) Weight: 425 kg	F	23.48 cm (0.67)
12 451	<i>M. leonina</i> 1 adult female	Kerguelen Islands	30/10/2012 to 15/02/2013	Time and depth every 4s	Length: 234.0 cm (1.0) Weight: 275 kg	F	18.52 min (0.64)
ct1	<i>M. leonina</i> 6 adult females	Husvik, South Georgia Island	06/01/2004 to 25/08/2004	CTD_GEN_07B†	Length: 242.5 cm (6.7)	F	23.96 min (0.07)
ct8	<i>M. leonina</i> 7 adult females	Husvik, South Georgia Island	13/01/2005 to 31/10/2005	CTD_GEN_07B	Length: 255.7 cm (2.0) Length: 301.3 cm (8.5)	F M	36.59 min (0.25)
ct40	<i>M. leonina</i> 5 adult females	Husvik, South Georgia Island	28/01/2008 to 09/12/2008	CTD_GEN_07B	Length: 250.5 cm (13.0) Length: 343.6 cm (24.4)	F M	27.47 min (0.38)
ct45	<i>M. leonina</i> 10 adult females	Husvik, South Georgia Island	17/10/2008 to 01/02/2009	CTD_GEN_07B	Length: 248.7 cm (2.9)	F	19.77 min (0.18)
ct49	<i>M. leonina</i> 11 adult females	Husvik, South Georgia Island	28/01/2009 to 12/01/2010	CTD_GEN_07B	Length: 240.5 cm (4.0) Length: 237.00 cm	F M	25.92 min (0.35)
ct58	<i>M. leonina</i> 13 adult females	Husvik, South Georgia Island	22/10/2009 to 05/02/2010	CTD_GEN_07B	Length: 249.5 cm (3.2)	F	18.97 min (0.19)

\*This tag sampled time and depth at 1 Hz, but resolution was reduced to one sample every 4s for consistency with other time-depth data. The tag operated on a 3-day duty cycle (3 days on, 3 days off).

†CTD\_GEN\_07B is the parameter specification with which the instruments were programmed.



linear abstracted profile deviates most from the detailed profile is added to the abstracted profile, creating two new line segments. This is called a break point. This step creates a new piecewise linear profile comprising  $I + 2$  points, connected by linear segments, and completes one iteration of the model (Fig. 1).

The maximum residual ( $R_I$ ) is calculated for each of the resulting line segments, and the point with the greatest departure is selected as the next break point and added to the profile. This process is repeated until the desired number of break points is reached, and the resulting piecewise linear abstracted profile has been constructed. When the abstraction process is complete, the abstracted time–depth dive profile includes  $I + 2$  time points ( $T_1$  to  $T_I$ ) and the corresponding  $I$  depth points ( $D_1$  to  $D_I$ ). The first and last points ( $T_0, T_{I+1}$ ,  $D_0$  and  $D_{I+1}$ ) are not transmitted. At  $T_0$ , time is considered to be zero, and at  $T_{I+1}$ , the time elapsed since the beginning of the dive will equal the dive duration. Similarly, at  $D_0$  and  $D_{I+1}$ , the depth will both be 0–6 m. The order in which the time–depth points were selected is not stored or transmitted.

#### HOW DOES THE REPRESENTATION OF THE DETAILED DIVE CHANGE WITH INCREASING BSM POINTS?

Using detailed dives from four high-resolution data sets, each representing a continuous dive record from one individual (10 943, 12 454, 12 453 and 12 451 in Table 1), we estimated the proportion of high-resolution samples that was represented by the number of BSM points in the corresponding abstracted dive profile (Table 2). We generated an abstracted profile with 3–12 BSM points for each study dive. This resulted in a data set of 2400 proportions, from 240 dives.

#### IS THE SAMPLE OF STUDY DIVES REPRESENTATIVE?

Four iterations of the BSM are carried out on-board CTD-SRDLs, resulting in a dive profile consisting of six time–depth points: two at the surface, and four at depth, at irregular times, which vary from dive to dive. The number of iterations of the BSM algorithm to be carried out on CTD-SRDLs was chosen as the minimum sufficient number to con-

vey the geometric shape of the dive profile, while keeping the number of computations low (Fedak, Lovell & Grant 2001).

We confirmed that the sample of detailed dives from the four individuals was consistent with elephant seal diving behaviour in general. To do this, the DZI was calculated for a sample of 4000 abstracted dives from 45 southern elephant seals instrumented in four different field seasons (1000 dives each from two post-moult deployments and two post-breeding deployments; ct40, ct45, ct49 and ct58) over 2 years at the island of South Georgia, South Atlantic (Table 1). The resulting distribution of DZI was visually compared with the distribution of DZI for the detailed dives (Appendix S2, Figs S2.1 and S2.2).

#### WHAT CAN WE LEARN FROM ABSTRACTED DIVES?

The process by which abstracted dive profiles arise when they are collected by CTD-SRDL is known; therefore, it is possible to reverse the deterministic steps and obtain limits to the depth at which the trajectory could have passed, before it was abstracted. The information required to build the dive zone (the 100% confidence zone for depth) includes the (i) temporal resolution at which the dive was recorded by the tag, (ii) the  $I + 2$  locations, in time and depth (including maximum dive depth), (iii) the residual associated with the final,  $I + 1$ th break point, and critically, (iv) the order in which break points were selected during abstraction. The temporal resolution of the detailed dive data is known from the duration of the dive, and the locations are received in the satellite transmissions, but the residuals and the order in which the break points were added need to be determined. Constructing the dive zone involves three steps.

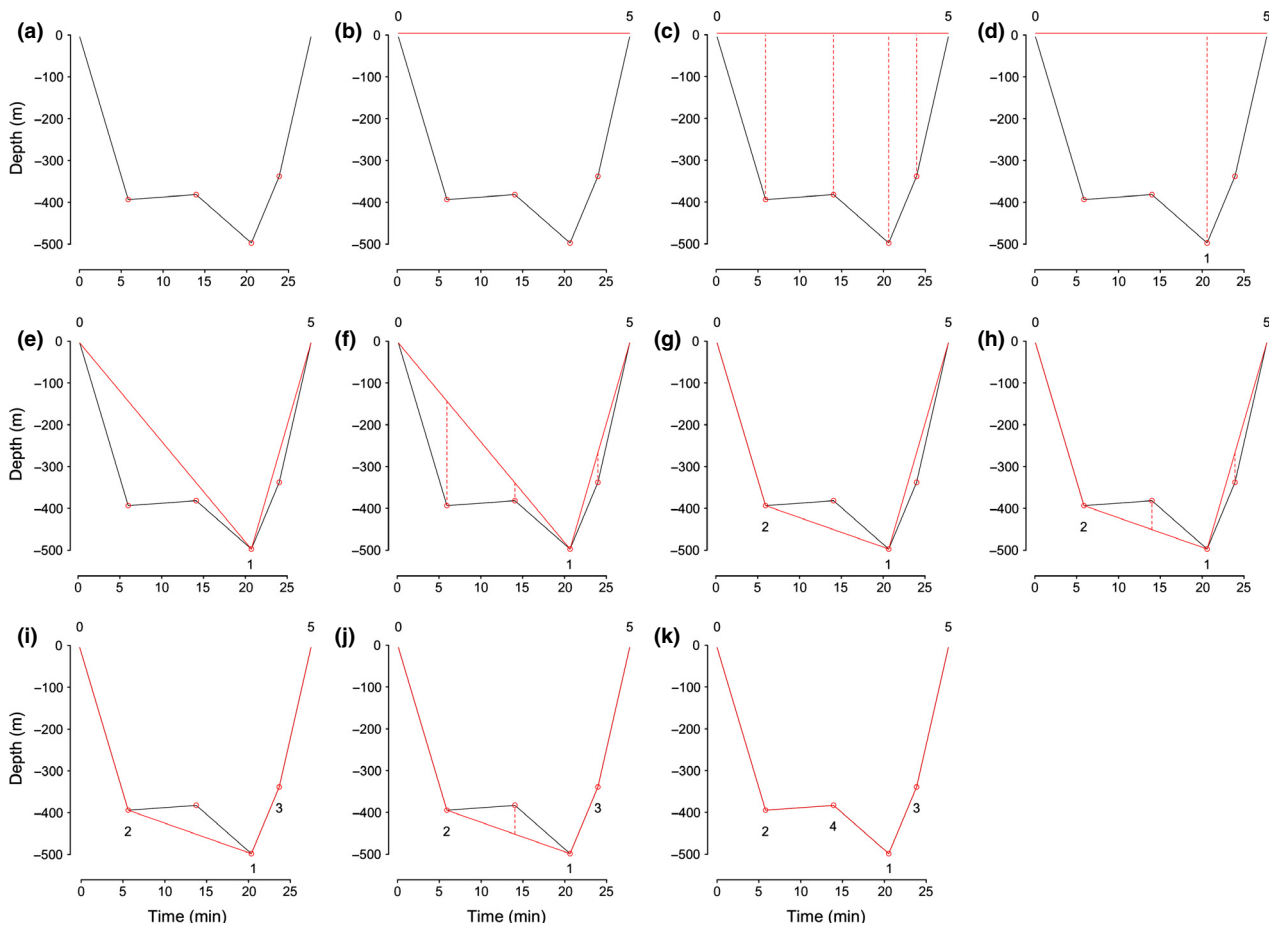
First, to find the order in which points were added to the profile, the BSM must be applied to the already abstracted dive trajectory (Fig. 2). At each of the  $I$  iterations, one of the break points is selected as the point of greatest deviation hence retrieving the order in which the break points were added.

Secondly, to determine the limits of the zone, the residuals corresponding to the break points need to be calculated. It is tempting to assume that the residual associated with the  $I$ th break point,  $R_I$ , applies to all segments in the final profile, determining the limits of the dive zone. This is not the case; instead, the limits of the dive zone at each iteration of the model interact with those from previous iterations. As break points are added to the profile, the dive zone changes shape and size. Although the dive zone will always get smaller with subsequent iterations of the algorithm, this geometric effect means that the dive zone needs to be constructed based on all iterations of the model, up to the last one. Furthermore, the resulting dive zone at the final iteration is not symmetric around all segments in a profile, and several break points can touch the limit of the dive zone (Fig. 3).

The depth points selected by the BSM are coded before transmission according to a pseudo-logarithmic mantissa and exponent representation. With this representation, resolution can be made proportional to the scale of the number being represented, making it useful for depth. Data are then truncated during the decoding process, once they are received. As a consequence, the received depth measurements are binned (i.e. each reported depth has an upper and lower bound) and bin width increases with depth. Bin width is usually smaller than the dive zone height at each break point. As a result, the dive zone is truncated near the deepest point(s) of the dive where the depths along the line segment near the bottom of the dive exceed the known accuracy of the greatest depth reached during the dive. This is justified since we know that no point in the true profile can be deeper than the maximum depth recorded by the tag. Depth is also truncated

**Table 2.** Observed mean proportion of detailed depth samples from individuals 10 943, 12 454, 12 453, 12 451 represented by an abstracted profile resulting from the BSM with a given number of break points (standard error, SE) and mean dive zone index (DZI) (standard error, SE). The grey row highlights the case of a dive profile with 6 points, produced by four iterations of the BSM

Number of BSM points	Mean proportion of depth samples represented by abstracted profile (SE)	Mean DZI (SE)
3	0.010 (0.0002)	0.993 (0.0002)
4	0.013 (0.0003)	0.549 (0.0156)
5	0.017 (0.0003)	0.303 (0.0103)
6	0.020 (0.0004)	0.158 (0.0048)
7	0.023 (0.0005)	0.110 (0.0036)
8	0.027 (0.0005)	0.085 (0.0030)
9	0.030 (0.0006)	0.071 (0.0026)
10	0.033 (0.0007)	0.060 (0.0023)
11	0.037 (0.0007)	0.051 (0.0019)
12	0.040 (0.0008)	0.044 (0.0017)



**Fig. 2.** Stepwise illustration of the calculation of the order in which break points were added to a time–depth dive profile abstracted onboard a CTD-SRDL telemetry device using the broken-stick model (BSM). The black line represents the abstracted time–depth dive profile received from the device. The red points represent known break points. The solid red lines represent the proposed BSM-abstracted dive path at each iteration of the model, and the dashed red lines represent the residuals that correspond to known break points.

at each break point, for the same reason, so the dive zone appears ‘pinched’ at each break point (Fig. 3).

When the  $i$ th iteration is complete, it is known with certainty that there were no points in the detailed trajectory that had a greater residual than the one corresponding to the last point (Fig. 2). All other depth points in the true trajectory will now have a smaller vertical distance to the abstracted profile.

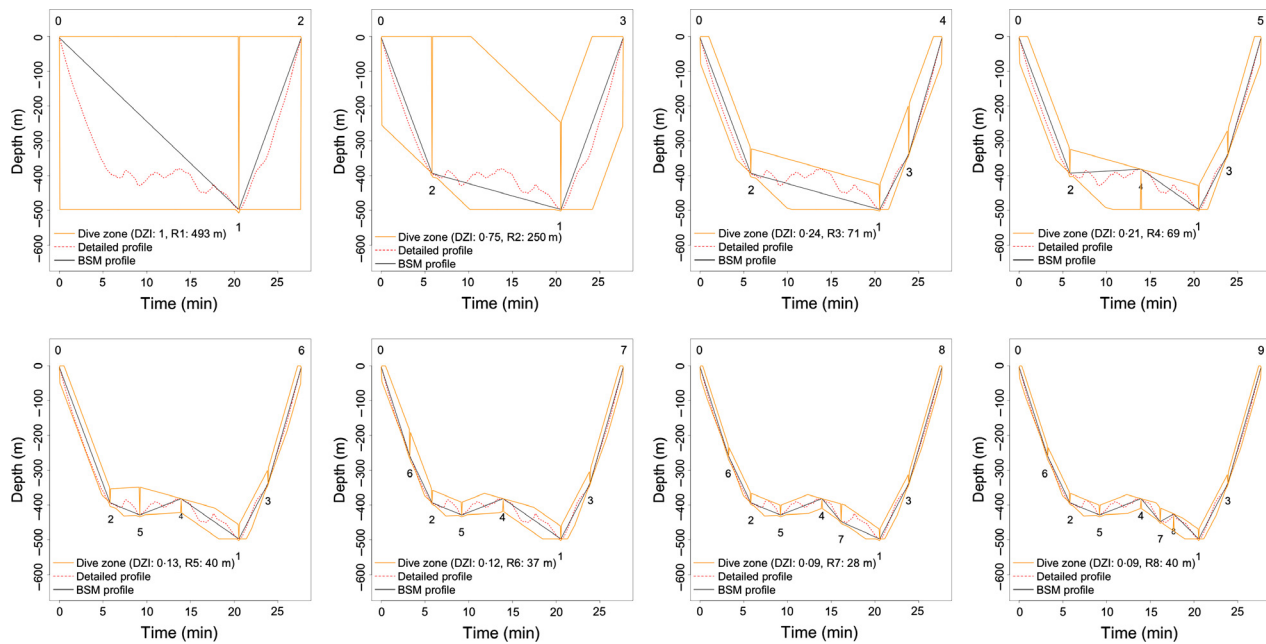
Thirdly, to construct the dive zone, a number of equally spaced time points need to be selected at which to sample vertical sections of the time–depth space. The resolution of time points should not exceed the resolution at which depth data were collected by the tag. At each of these time points ( $t, \dots, t_{\max}$ ), the estimated lower,  $L_t$ , and upper,  $U_t$ , depth bounds define the depth interval through which the true trajectory passed with 100% confidence (Fig. 3). Computer code for all algorithms described was written in R (R Core Team 2013) and can be found in the Supporting information (Appendix S1).

#### HOW IS THE DZI DERIVED?

The approximation of a nonlinear path will improve as the number of points ( $N$ ) that are used to approximate it increases. This should be reflected in a reduction in the size of the maximum residual at the last iteration, as the number of BSM iterations increases. In absolute terms, a small value for  $R_t$  means that the biggest vertical outlier in

the true path was a small distance away from the abstracted path, which might suggest a better fit. However,  $R_t$  is affected by the depth, duration and sinuosity of the detailed dive trajectory and does not follow a strictly decreasing relationship with the number of iterations of the model. When detailed dive data are not available, a reliable and unbiased way of measuring goodness-of-fit is required to assess the accuracy of abstracted profiles, irrespective of depth and duration. On its own,  $R_t$  does not provide an objective way of assessing goodness-of-fit over the whole dive, because it depends on the maximum depth of the dive and the slopes and lengths of the segments that make up the abstracted profile. The mean dive zone height should be a better measure of fit, because it does contain information on the whole dive and includes the geometric effects that result from the slopes and lengths of the segments that make up the abstracted dive, though importantly, it also does not include information on sinuosity. However, to be comparable between dives, it needs to be standardized by dive depth and dive duration. This is the basis for the construction of the DZI.

The DZI is calculated using the sum of the differences between the upper,  $U$ , and lower,  $L$ , limits of the dive zone at each time step,  $t$ , in a dive. The sum of these heights is divided by the product of the maximum dive depth,  $\max_{\text{dep}}$ , and the number of depth points that were recorded by the tag for the dive prior to abstraction,  $t_{\max}$ . This quotient ranges between 0 and 1, where values close to 0 indicate a narrow zone



**Fig. 3.** Stepwise illustration of the construction of the dive zone and its evolution as points are added to the abstracted dive profile. We use an example dive from individual 10 943 (see Table 1), with up to 10 iterations of the broken-stick model, for example, using 3–12 points, in total.

around an abstracted dive and are desirable, and values close to 1 indicate a wide zone around an abstracted dive and relatively low confidence in the abstracted dive profile.

$$DZI = \frac{\sum_{t=1}^{t_{\max}} U_t - L_t}{\maxdep * t_{\max}} \quad \text{eqn 1}$$

#### CAN THE DZI BE USED AS A PROXY FOR THE RSS?

When detailed dive data are available, goodness-of-fit can also be assessed using the sum of squared residuals (RSS) between the detailed and abstracted depths,

$$RSS = \sum_{t=1}^{t_{\max}} (T_t - A_t)^2 \quad \text{eqn 2}$$

where  $T_t$  is the  $t$ th depth in the detailed profile, and  $A_t$  is the  $t$ th depth in the linearly interpolated abstracted profile. The RSS, an absolute measure of fit for detailed dive profiles, was compared to the DZI, the relative measure of fit developed here for abstracted dives, and the biggest residual at the last iteration of the BSM,  $R_I$  to describe the relationship between them and determine whether they could be used as a proxy for the RSS.

To do this, the RSS, DZI and  $R_I$  were calculated for abstracted profiles with 3–12 BSM points for each of the 240 study dives ( $n = 2400$ ). RSS is a nonzero real number, so a generalized additive model (GAM) with a gamma distribution and a log link function were used, fitted with `mgcv` in R (Wood 2000, 2011). We wanted to know whether the DZI could be used as a proxy for the RSS in already abstracted dives received by CTD-SRDL, but we were also interested to know whether the amount of deviance explained by the DZI changed with increasing break points. Hence, we fitted a model to a subset of the data representing six break points and did model selection to find the best model, but also fitted a model to a subset of the data representing each number of break points explored in this study (3–12) including only DZI as a covariate without doing model selection.

In the first case, RSS was modelled with DZI and  $R_4$  as explanatory variables and individual seal as a random effect. In the second case, RSS was modelled with DZI as the only explanatory variable and individual seal as a random effect, as above. In both cases, the relationship between the RSS and the covariates was nonlinear, so they were fitted as smooth functions with a shrinkage smoother ('cs') as the basis function and  $k = 4$  knots. The number of knots was found to be sufficient using standard `mgcv` checks. This basis allows for the smooth coefficients to be shrunk to zero and effectively removed from the model when there is no relationship with the response. We specified a gamma parameter value of 1.4 to reduce the chance of overfitting. The random effect was fitted using the 're' smooth, as described above. We used restricted maximum likelihood (REML) as the fitting method (Wood 2011).

#### DOES THE DZI VARY SYSTEMATICALLY BETWEEN DIVE TYPES?

The availability of detailed dive data (depth sampled every 4 s) made it possible to investigate the effects of the abstraction process, develop methods to reverse the abstraction and quantify goodness-of-fit. Detailed data were taken from four high temporal resolution data sets, one recorded by a specially configured archival SRDL deployed on a northern elephant seal after the moult at Año Nuevo, California, and three from female southern elephant seals recorded by time-depth recorders deployed after the breeding season at Kerguelen Islands (Table 1). These dives were chosen by eye, as being representative of the six functional or behavioural characterizations sometimes used to classify elephant seal dives into types; U-shaped dives (U), V-shaped dives (V), square-bottom dives (SQ), wiggle dives (W), root-shaped dives (R) and drift dives (DR), after Hindell, Slip & Burton (1991). Ten dives of each type from each individual were included in the case study. A sample of over 22 000 abstracted dives that had already been classified into types was used to investigate systematic differences in goodness-of-fit between dive types (Table 1).

The classification of abstracted dives into types was done using the random forest tree-building method (Breiman 2001). Random forest is a machine learning tool; we used an implementation in the `randomForest` library in `R` (Liaw & Wiener 2002). A supervised version of the method was used to classify dives, whereby 3000 dives, 14% of the data set, were classified based on visual cues and used to train the remaining dives. The overall 'out-of-bag' error, an unbiased estimate of classification error, was 3.6%. This represents the aggregate of the prediction error rate at each bootstrap iteration (Liaw & Wiener 2002). The variables supplied to the function for classification were maximum dive depth, bathymetry and 15 dive parameters (Photopoulos 2007). This method has been found to work well for dive classification, using both detailed and abstracted dive data (Thums, Bradshaw & Hindell 2008).

The DZI and  $R_4$  were also calculated for each dive in the data set. The DZI used as the response variable in GAM with a quasibinomial distribution and a logit link function, dive type as a factor variable,  $R_4$  as a smooth covariate using 'cs' basis function ( $k = 4$  knots, checked as above) and individual animal as a simple random effect using the 're' smooth. The model was fitted using the REML method, and gamma was specified as 1.4, as above.

Model assessment was done based on the inspection of the residuals, the relationship between the observations and the fitted values for each model and the percentage of deviance explained.

## Results

In the introduction, we outlined five questions relating to the BSM for dive profile abstraction, the usefulness of the DZI as a goodness-of-fit measure and the validity of our assessment of it as such: (1) How does the BSM work for dive profile abstraction? How does the representation of the detailed dive change with increasing BSM points? (2) Is the sample of study dives representative? (3) What can we learn from abstracted dives? (4) How is the DZI derived? and (5). We present the results corresponding to these questions as a numbered list. Can the DZI be used as a proxy for the RSS? Does the DZI vary systematically between dive types?

1. The BSM algorithm for dive profile abstraction is illustrated in Fig. 1 and animated in Appendix S2 of the Supporting information. The proportion of detailed samples in a dive increases with the number of break points in the abstracted profile and depends on dive duration. Overall, the proportion of a dive represented by its abstracted profile is very low, starting at 1% with three break points and reaching 4% with twelve, and increases by a constant 0.33% for each break points added (Table 2). The mean relationship was similar for all animals, but there were differences in the variability in the relationship (Fig. 4).

2. Variability between the samples of abstracted and detailed dives was expected, due to individual variability and the different regions where data were collected. However, visual comparison of the distributions of DZI from the different samples did not suggest any striking differences (Appendix S2, Figs S2.1 and S2.2, Supporting information).

3. Even though the detailed trajectory cannot be recovered unless the device is physically recovered, the order in which break points were added to the profile and the 100% confi-

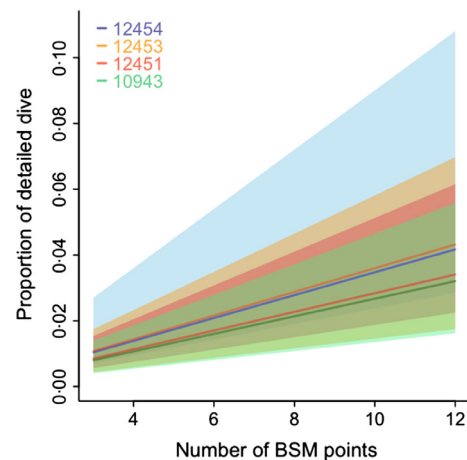


Fig. 4. The relationship between the proportion of high-resolution time-depth samples and the number of break points in the abstracted dive profiles of 240 case study dives from a northern elephant seal and three southern elephant seals (see Table 1). Abstracted profiles with 3–12 points in total were generated for each of the study dives and compared with the full-resolution profile. The coloured areas includes minimum and maximum observed range of the relationship for each individual.

dence limits to the detailed profile, which we call the dive zone, can be calculated from abstracted dives. This makes it possible to derive the DZI (Fig. 3).

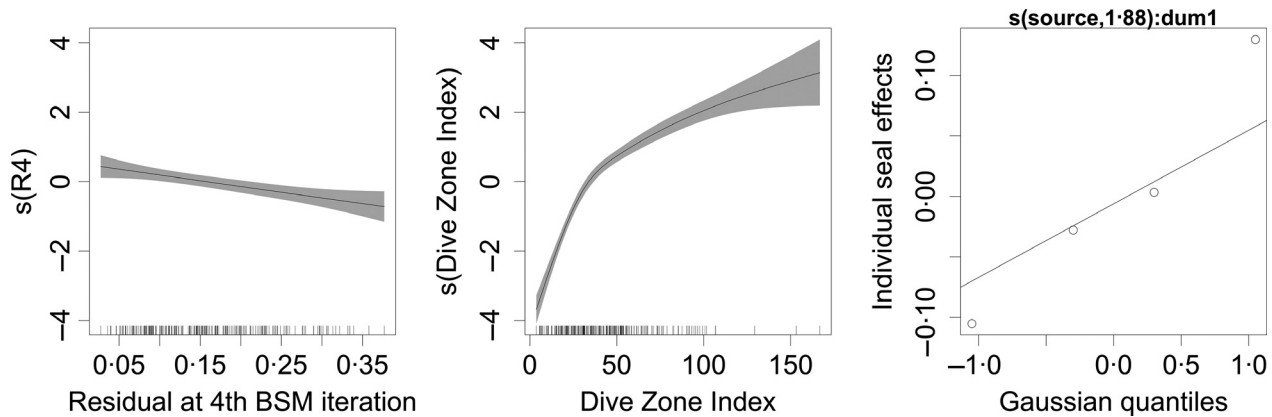
4. The derivation of the DZI is based on the maximum depth of the dive, its duration and the upper and lower limits to the dives zone.

5. There was a strong, positive relationship between the RSS and the DZI together with  $R_4$ . We found that the DZI,  $R_4$  and a random effect for individual explained 69% of the variability in the RSS for abstracted profiles with six break points (deviance explained) (Fig. 5). On its own, the DZI explained a variable proportion of deviance for abstracted dives with differing numbers of break points, but there was an overall positive relationship with increasing break points for dives with four or more break points (Fig. 6). The DZI varied substantially between dive type and had an increasing relationship with  $R_4$  (Fig. 7). The dive type associated with the biggest DZI values were square dives (SQ), and both V-shaped (V) and DR had the smallest DZI, under the model (Fig. 8). Dive type and  $R_4$  together explained 65% of the variability in the DZI (deviance explained), having accounted for individual variability by fitting a random effect for individual.

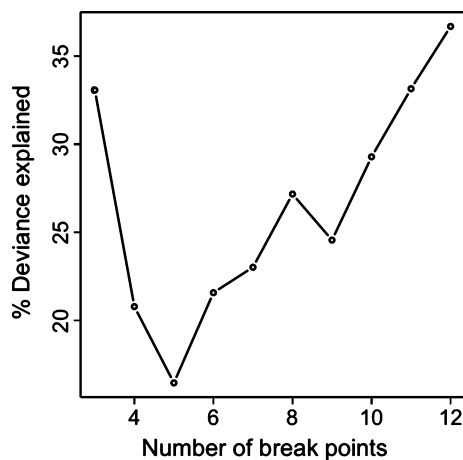
## Discussion

Change-point models are used in many fields, from neuroscience and epidemiology to genetics and finance, to identify the changes in time series. One of these, the BSM, was adopted on-board CTD-SRDLs, as a working solution to the problem of linearly approximating a nonlinear path in the vertical dimension with as little information as possible, while aiming to retain biologically relevant content. Until now, the





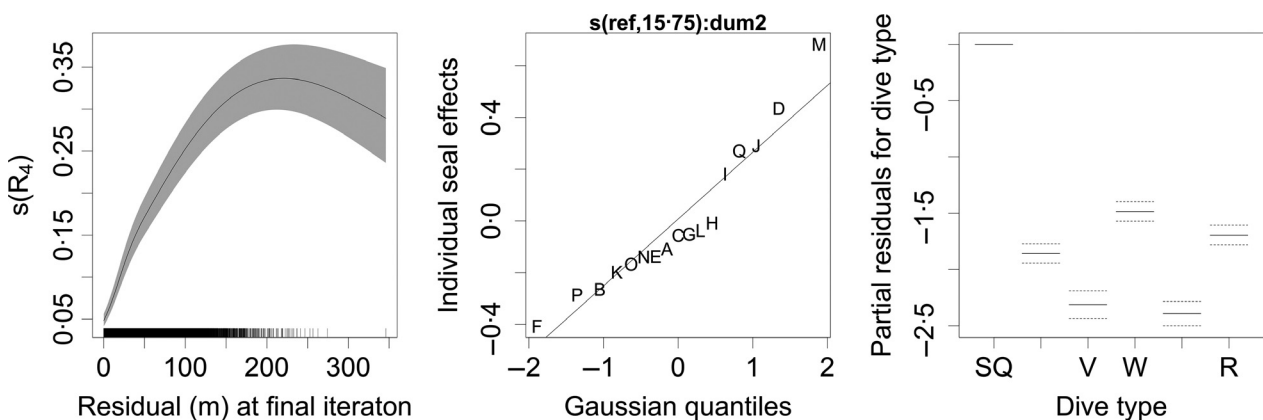
**Fig. 5.** Smooth functions of the covariates in a generalized additive model with  $\log(\text{RSS})$  as the response, dive zone index (DZI) a smooth covariate and individual dive as a random effect. The data used to fit this model were abstracted dive profiles of 4000 study dives from 45 southern elephant seals instrumented at South Georgia Island in 2008 and 2009. The grey area includes two standard errors for the fitted relationship. RSS, residual sum of squares.



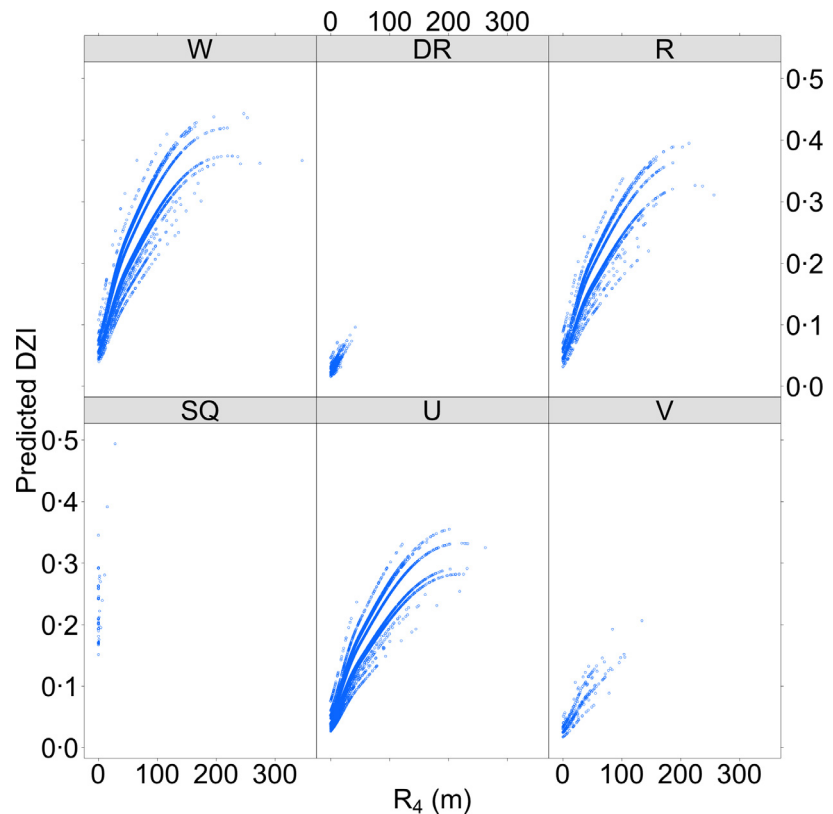
**Fig. 6.** The relationship between the percentage of deviance explained by a model for residual sum of squares (RSS) with dive zone index (DZI) as the explanatory variable and individual seal as a random effect, and the number of breakpoints in the dives being considered.

approximation error associated with abstracted profiles had not been investigated. We provide a way of calculating and summarizing the error associated with dive profiles derived using the BSM, in order to assess the information content of abstracted dives. Our results suggest that BSM-abstracted dive profiles do, in fact, retain enough information to estimate goodness-of-fit of the abstracted profile to the detailed profile. The strong, positive relationship between the DZI and  $R_4$ , and the RSS is evidence for this. It means that researchers using BSM-abstracted dive profiles to make inferences about animal behaviour and diving ecology can now calculate the DZI and  $R_4$  and incorporate a relative measure of error into their analyses.

With the BSM, as with other linear approximation methods, the number of iterations is critical to the quality of the abstracted dive. Our results confirm the findings of a preliminary investigation during tag development, regarding the adequacy of different numbers of iterations, which resulted in the standard use of four iterations of the algorithm for dives collected by CTD-SRDLs. That study found that the information



**Fig. 7.** Smooth functions of the covariates in a generalized additive model with dive zone index as the response, dive type as a factor variable, the residual as the final iteration of the broken-stick model (BSM) as a smooth covariate and individual seal as a random effect. The data used to fit this model were abstracted dive profiles of 22 305 study dives from 17 southern elephant seals instrumented at South Georgia Island in 2004 and 2005. The grey area includes two standard errors for the fitted relationship.



**Fig. 8.** The fitted relationship between the dive zone index (DZI) and the residual at the final interaction of the broken-stick model ( $R_4$ ) for each dive type. These are the predictions based on a generalized additive model, with  $R_4$  as smooth covariate, dive type as a factor variable and individual seal as a random effect. The data used to fit this model were abstracted dive profiles of 22 305 study dives from 17 southern elephant seals instrumented at South Georgia Island in 2004 and 2005.

gained in the transmission of a fifth break point is relatively small and that it would be more useful to receive a measure of the variance of the detailed dive, either for the whole dive or each segment, which is supported by the results presented in this study.

It is worth noting that although the BSM can efficiently summarize a curvilinear trajectory, using a piecewise linear shape, profiles resulting from low iteration numbers may mask important biological features even if they closely approximate the detailed trajectory. Our methods do not provide means for assessing the biological content of abstracted dives, since the detail lost, however, small, might be the most biologically interesting. For example, a useful feature of the BSM is that it is efficient at identifying long sections where the trajectory has low variability. In the case of dives, these are often the descent and ascent phases, leaving only two points to confer information about the bottom phase of the dive, which is arguably the most interesting biologically. Dives with low variability in change in depth in the bottom phase and also have the lowest DZI, as we found here for drift dives (DR type). However, when classifying abstracted dives with a method like a random forest algorithm, ancillary behavioural data are necessary for validating classes as being functionally distinct, in addition to being phenomenologically distinct.

Large numbers of dive profiles are collected using CTD-SRDLS from many different species (over 21 million profiles since 1991, SMRU 2012, unpublished data) and used to make inferences about the biology and behaviour of the instrumented animals. It seems essential that a method for assessing the accuracy of these abstracted dives, at least statistically if not

biologically, is made widely available. The methods we have presented here make that possible. They also provide a way of carrying out a 'pilot' analysis when detailed dive data are available. The differences in diving behaviour between species and habitats may render different numbers of break points appropriate for the questions being asked. When detailed dive data are available, the result of dive abstraction with different number of iterations can be investigated to achieve the best result for a specific study, prior to deployment. Together, these uses for our methods may help make more robust the behavioural conclusions we can draw from telemetry data.

More generally, through this work, we have developed a method for quantifying uncertainty in the fitted values for BSMs, when the original time series is no longer available. This method could be applied to any situation where the original time-series data are no longer available. This could be useful in situations where large amounts of data are being generated and cannot be stored or transmitted at the original resolution. We demonstrate that measures of fit like the DZI and  $R_4$  have a strong correspondence to the RSS and could therefore be used instead.

## Acknowledgements

This work was supported by SMRU Ltd (now SMRU Marine) in the form of a PhD fellowship (T.P.). Completion of the manuscript was supported by a National Research Foundation Scarce Skills Postdoctoral Fellowship at the University of Cape Town, South Africa (T.P.). The CTD-SRDLS data presented in this manuscript were collected as part of a project funded by the Natural Environment Research Council (NERC) grants NE/E018289/1 and NER/D/S/2002/00426. The authors are extremely grateful for access to three high-resolution data sets from Kerguelen Islands that were made available by Christophe Guinet and

collected under the SO-MEMO framework. Fieldwork was carried out according to the Animals (Scientific Procedures) Act 1986 guidelines and approved by the University of St Andrews Animal Welfare Committee. The authors thank Martin Biuw (Akvaplan-niva, Tromsø, Norway) and Lars Boehme (SMRU) for use of the CTD-SRDL data, Nora Hanson (SMRU) and Lars Boehme for useful comments on the manuscript, Samantha Gordine (SMRU) for testing the code and two anonymous reviewers for their constructive comments on the manuscript.

## Data accessibility

- R scripts: uploaded as supporting material
- Example data of detailed and abstracted dive data: uploaded as supporting material
- For access to the full data sets used in this study please contact the data owners directly. South Georgia Island and Año Nuevo deployments: Mike Fedak (maf3@st-andrews.ac.uk) and Lars Boehme (lb284@st-andrews.ac.uk). Kerguelen Islands deployments: Christophe Guinet (Christophe.GUINET@cebc.cnrs.fr).

## References

- Argos (2011) Argos User's Manual. URL [http://www.argos-system.org/files/pmedia/public/r363\\_9\\_argos\\_manual\\_en.pdf](http://www.argos-system.org/files/pmedia/public/r363_9_argos_manual_en.pdf).
- Authier, M., Martin, C., Ponchon, A., Steeland, S., Bentaleb, I. & Guinet, C. (2012) Breaking the sticks: a hierarchical change-point model for estimating ontogenetic shifts with stable isotope data. *Methods in Ecology and Evolution*, **3**, 281–290.
- Baechler, J. (2002) Dive shapes reveal temporal changes in the foraging behaviour of different age and sex classes of harbour seals (*Phoca vitulina*). *Canadian Journal of Zoology*, **80**, 1569–1577.
- Bailleul, F., Charrassin, J.-B., Monestiez, P., Roquet, F., Biuw, M. & Guinet, C. (2007) Successful foraging zones of southern elephant seals from the Kerguelen Islands in relation to oceanographic conditions. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **362**, 2169–2181.
- Beckage, B., Joseph, L., Belisle, P., Wolfson, D.B. & Platt, W.J. (2007) Bayesian change-point analyses in ecology. *The New Phytologist*, **174**, 456–467.
- Biuw, M., Boehme, L., Guinet, C., Hindell, M., Costa, D., Charrassin, J.-B. et al. (2007) Variations in behavior and condition of a Southern Ocean top predator in relation to in situ oceanographic conditions. *Proceedings of the National Academy of Sciences of the United States of America*, **104**, 13705–13710.
- Biuw, M., Nøst, O., Stien, A., Zhou, Q., Lydersen, C. & Kovacs, K. (2010) Effects of hydrographic variability on the spatial, seasonal and diel diving patterns of southern elephant seals in the eastern Weddell Sea. *PLoS One*, **5**, 1–14.
- Breiman, L. (2001) Random Forests. *Machine Learning*, **45**, 5–32.
- Campagna, C., Piola, A.R., Marin, M.R., Lewis, M., Zajackovski, U. & Fernández, T. (2007) Deep divers in shallow seas: Southern elephant seals on the Patagonian shelf. *Deep Sea Research Part I: Oceanographic Research Papers*, **54**, 1792–1814.
- Fauchald, P. & Tveraa, T. (2003) Using first-passage time in the analysis of area-restricted search and habitat selection. *Ecology*, **84**, 282–288.
- Fedak, M.A., Lovell, P. & Grant, S. (2001) Two approaches to compressing and interpreting time-depth information as collected by time-depth recorders and satellite-linked data recorders. *Marine Mammal Science*, **17**, 94–110.
- Fedak, M.A., Lovell, P., McConnell, B. & Hunter, C. (2002) Overcoming the constraints of long range radio telemetry from animals: getting more useful data from smaller packages. *Integrative and Comparative Biology*, **42**, 3–10.
- Hebblewhite, M. & Haydon, D.T. (2010) Distinguishing technology from biology: a critical review of the use of GPS telemetry data in ecology. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **365**, 2303–2312.
- Hindell, M., Slip, D. & Burton, H. (1991) The diving behaviour of adult male and female southern elephant seals, *Mirounga leonina* (Pinnipedia: Phocidae). *Australian Journal of Zoology*, **39**, 595–619.
- Jonker, F. & Bester, M. (1998) Seasonal movements and foraging areas of adult southern female elephant seals, *Mirounga leonina*, from Marion Island. *Antarctic Science*, **10**, 21–30.
- Liaw, A. & Wiener, M. (2002) Classification and regression by randomForest. *R News*, **2**, 18–22.
- MacArthur, R. (1957) On the relative abundance of bird species. *Proceedings of the National Academy of Sciences of the United States of America*, **43**, 293–295.
- MacArthur, R. & MacArthur, J. (1961) On bird species diversity. *Ecology*, **42**, 594–598.
- McConnell, B., Chambers, C. & Fedak, M.A. (1992) Foraging ecology of southern elephant seals in relation to the bathymetry and productivity of the Southern Ocean. *Antarctic Science*, **4**, 393–398.
- McConnell, B., Fedak, M.A., Lovell, P. & Hammond, P. (1999) Movements and foraging areas of grey seals in the North Sea. *Journal of Applied Ecology*, **36**, 573–590.
- McKellar, A.E., Langrock, R., Walters, J.R. & Kesler, D.C. (2014) Using mixed hidden Markov models to examine behavioral states in a cooperatively breeding bird. *Behavioral Ecology*, 1–10. doi: 10.1093/beheco/aru171.
- Photopoulou, T. (2007) *Behavioural Changes of a Long-Ranging Diver in Response to Oceanographic Conditions*. University of St Andrews, St Andrews, Scotland.
- R Core Team (2013) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rual, P. (1989) For a better XBT bathy-message on board quality control, plus a new data reduction method. *Western Pacific International Meeting and Workshop on Toga Coare* (eds J. Picaut, R. Lukas & T. Delcroix), pp. 823–833. Noumea, New Caledonia.
- Schreer, J., Kovacs, K. & O'Hara Hines, R. (2001) Comparative diving patterns of pinnipeds and seabirds. *Ecological Monographs*, **71**, 137–162.
- Schreer, J. & Testa, J. (1996) Classification of Weddell seal diving behavior. *Marine Mammal Science*, **12**, 227–250.
- Thums, M., Bradshaw, C.J.A. & Hindell, M.A. (2008) A validated approach for supervised dive classification in diving vertebrates. *Journal of Experimental Marine Biology and Ecology*, **363**, 75–83.
- Toms, J. & Lesperance, M. (2003) Piecewise regression: a tool for identifying ecological thresholds. *Ecology*, **84**, 2034–2041.
- Wood, S.N. (2000) Modelling and smoothing parameter estimation with multiple quadratic penalties. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **62**, 413–428.
- Wood, S.N. (2011) Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **73**, 3–36.

Received 13 June 2014; accepted 26 November 2014

Handling Editor: Luca Börger

## Supporting Information

Additional Supporting Information may be found in the online version of this article.

**Appendix S1.** R code and example data for BSM functions.

**Appendix S2.** Supplementary figures and animations.