
van der Vaart E, Prangle D, Sibly R. [Taking Error Into Account When Fitting Models Using Approximate Bayesian Computation](#). *Ecological Applications* 2018. DOI: 10.1002/eap.1656

Copyright:

Copyright by the Ecological Society of America. This is the peer reviewed version of the following article: [van der Vaart E, Prangle D, Sibly R. Taking Error Into Account When Fitting Models Using Approximate Bayesian Computation. *Ecological Applications* 2018. DOI: 10.1002/eap.1656], which has been published in final form at <https://doi.org/10.1002/eap.1656>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Self-Archiving.

DOI link to article:

<https://doi.org/10.1002/eap.1656>

Date deposited:

30/12/2017

1 **Running Head:**

2 TAKING ERROR INTO ACCOUNT WITH ABC

3

4 **Title:**

5 Taking Error Into Account When Fitting Models Using Approximate Bayesian Computation

6

7 **Authors:**

8 Elske van der Vaart^{a,b*}, vdrvaart@uvt.nl, +31 13 466 31 60

9 Dennis Prangle^c, dennis.prangle@newcastle.ac.uk

10 Richard M. Sibly^a, r.m.sibly@reading.ac.uk

11

12 ^aSchool of Biological Sciences, University of Reading, Harborne Building, University of
13 Reading, Whiteknights, Reading, Berkshire, RG6 6AS, United Kingdom

14 ^bCognitive Science and Artificial Intelligence, Tilburg University, School of Humanities, PO
15 Box 90153, 5000 LE Tilburg, the Netherlands

16 ^cSchool of Mathematics and Statistics, Newcastle University, Herschel Building, Newcastle
17 University, Newcastle upon Tyne, NE1 7RY, United Kingdom

18

19 **Abstract** (maximum 200 words)

20 Stochastic computer simulations are often the only practical way of answering questions relating
21 to ecological management. However, due to their complexity, such models are difficult to
22 calibrate and evaluate. Approximate Bayesian Computation (ABC) offers an increasingly
23 popular approach to this problem, widely applied across a variety of fields. However, ensuring

24 the accuracy of ABC's estimates has been difficult. Here, we obtain more accurate estimates by
25 incorporating estimation of error into the ABC protocol. We show how this can be done where
26 the data consist of repeated measures of the same quantity and errors may be assumed to be
27 normally distributed and independent. We then derive the correct acceptance probabilities for a
28 probabilistic ABC algorithm, and update the 'coverage test' with which accuracy is assessed. We
29 apply this method – which we call 'error-calibrated ABC' – to a toy example and a realistic 14-
30 parameter simulation model of earthworms that is used in environmental risk assessment. A
31 comparison with exact methods and the diagnostic 'coverage test' show that our approach
32 improves estimation of parameter values and their credible intervals for both models.

33

34 **Keywords**

35 ABC, IBM, approximate Bayesian computation, individual-based model, parameter estimation

36

37 **Introduction**

38 Stochastic computer simulations are increasingly used to make realistic predictions about real
39 world ecological processes (Hartig et al. 2011); from the survival of shorebirds (West et al.
40 2002) to the effects of climate change (Zurell et al. 2012) and the invasiveness of plants
41 (Nehrbass and Winkler 2007). Because such models attempt to simulate all relevant aspects of a
42 real physical system, they often involve many parameters, some of which will be difficult to set
43 correctly. Understanding the overall uncertainty introduced by these unknown parameter values
44 is crucial, especially when the final objective of these models is to assess the possible
45 consequences of management decisions, such as the translocation of vulnerable species
46 (Lethbridge and Strauss 2015) or the placement of wind turbines (Nabe-Nielsen et al. 2014).

47
48 Approximate Bayesian Computation, or ABC, is a promising technique for estimating parameter
49 values together with their credible intervals. Standard Bayesian methods explore properties of
50 the multivariate posterior distribution over the parameters (Gelman et al. 2013), often by
51 sampling parameter vectors from it. This posterior distribution specifies the degree of support for
52 different parameter vectors given the model, data and prior knowledge about the values the
53 parameters are likely to take. Sampling from the exact posterior is not always feasible, leading to
54 the development of approximate Bayesian methods, such as ABC.

55
56 Originally developed within population genetics (Tavaré et al. 1997, Pritchard et al. 1999,
57 Beaumont et al. 2002), ABC is now widely used, with recent applications to, for example, range
58 expansions (Rasmussen and Hamilton 2012), infectious diseases (Kosmala et al. 2016), and
59 forest dynamics (Lagarrigues et al. 2015). However, ensuring the accuracy of ABC's estimates
60 remains difficult. Here, we improve the estimation process for cases where the data consists of
61 repeated measures of the same quantity, such as a time series. We do this using Wilkinson
62 (2013)'s insight that accurate estimates can be obtained if the form of the error – the distribution
63 of the differences between model outputs and data – is incorporated into the ABC protocol.

64
65 Bayesian inference generally requires an analytical likelihood, expressing how the likelihood of
66 the data depends on the model parameters, but for mechanistic simulation models, this is often
67 not possible. Instead, ABC is based on simulations using the model. By repeatedly sampling
68 parameters from a model's prior, running the model, and then retaining the simulations closest
69 the data according to some distance function, ABC can approximate a model's posterior with an

70 accuracy that depends on the distance allowed between model outputs and data. This version of
71 ABC is referred to as ‘rejection ABC’. However in many cases even the best-fitting model will
72 not replicate the data exactly – even with the best parameters, there will always be some residual
73 distance between the model and the data, due to either model misspecification, observational
74 measurement error, or both. In these cases, taking error into account can greatly increase
75 posterior accuracy. Accounting for different types of error is well established in deterministic
76 modelling (e.g., Campbell 2006, Higdon et al. 2008, Goldstein and Rougier 2009), but Wilkinson
77 (2013) was the first to consider it in the context of stochastic computer simulations and ABC.

78
79 Wilkinson’s (2013) method assumes that the data measurements D can be considered as a
80 realization of the model η run with its input parameters θ set at their best values, $\hat{\theta}$, plus an
81 independent term ϵ representing error (Equation 1). If the distribution of ϵ is known, Equation 1
82 determines a probability distribution for D given the input value of $\hat{\theta}$. Therefore there is an
83 associated likelihood function. However, for most simulators $\eta(\hat{\theta})$ is extremely complicated, so
84 the likelihood function cannot be expressed as a simple mathematical formula. This means
85 standard Bayesian or maximum likelihood methods cannot be used.

$$D = \eta(\hat{\theta}) + \epsilon \quad \text{Equation 1}$$

86 The distribution of ϵ would ideally be based on a priori knowledge, with a principled
87 decomposition into model and measurement error. However for many ecological applications,
88 this is not practical. Any model concerned with the behavior of real organisms will have
89 structural inadequacies that are difficult to formally characterise, and many models are validated
90 against empirical data that was collected long ago, by other researchers, so that measurement

91 error is also unknown. In this paper, we present a simple approach to using Wilkinson's (2013)
92 method in cases where the empirical data consist of many data points of the same type.

93

94 Using the difference between the observations and the model at its best-fitting parameter values,
95 we parameterise a normally distributed estimate of the error, and then derive the corresponding
96 optimal acceptance probabilities for a new 'error-calibrated ABC' algorithm. We illustrate the
97 use of this new algorithm by analysing both a toy example and a complex computer simulation
98 of earthworms (Johnston et al. 2014), which was developed for the purpose of pesticide risk
99 assessment. This model was previously calibrated using 'rejection ABC' (van der Vaart et al.
100 2015), but a diagnostic 'coverage test' showed some inaccuracies in the posteriors. In this paper,
101 we update this diagnostic so that it also takes error into account, and show that 'error-calibrated
102 ABC' improves the estimation process for both the toy example and the earthworm simulation.

103 **Methods**

104 In previous work (van der Vaart et al. 2015) we implemented the most basic form of ABC,
105 'rejection ABC', using Algorithm 1. 'Rejection ABC' takes a sample of the parameter values
106 needed to run the model from a prior distribution which expresses existing knowledge about
107 what values each parameter is likely to take. The model is run with those parameter values, and
108 then the process is repeated thousands of times with different sets of parameter values randomly
109 drawn from the prior distribution. 'Rejection ABC' rejects all but the m best parameter values,
110 i.e., the m values that produce model outputs closest to the data points. These are samples from
111 an approximation to the Bayesian posterior distribution. The exact posterior distribution gives
112 the degree of support for each parameter vector, combining prior information and model
113 observations, and is used to produce univariate posterior distributions for each individual

114 parameter, as well as 95% credible intervals. The accuracy of the ABC approximation to the
 115 posterior can be assessed using ‘coverage tests’ (Prangle et al. 2013).

1. Repeat n times:
 - a. Draw $\theta^i \sim \pi(\theta)$ (the prior distribution)
 - b. Simulate $X^i \sim \eta(\theta^i)$ (the computer model)
2. Accept the m runs (θ^i, X^i) that minimise $\rho(X^i, D)$.

116 **Algorithm 1. Original ‘rejection ABC’ algorithm used in van der Vaart et al. (2015).**

117 The computer model is represented by $\eta(\theta)$, with output X and input parameters θ . This model is
 118 stochastic: repeated evaluations using the same input usually produce different outputs. Though
 119 our methods are also valid for deterministic models, better alternatives are available for those
 120 cases. X is a vector of model outputs which are to be compared with a data vector D . θ is a vector
 121 of model parameters, drawn from a prior distribution, $\pi(\theta)$. We often specify a prior distribution
 122 for each individual parameter and form the overall prior by an independence assumption. In total,
 123 n model runs are done, and ρ is the distance between the model output X and the data D . The m
 124 runs that minimise ρ are accepted and then the accepted (θ^i, X^i) pairs form a sample from an
 125 approximate posterior. Since both parameters and outputs are vectors, we use subscripts to
 126 denote particular components. For example, θ_j^i represents the j^{th} parameter for model run i , while
 127 X_j^i represents the model output corresponding to the j^{th} data point in model run i .

128 1.1. Coverage

129 Coverage tests were introduced by Prangle et al. (2013) to check the accuracy of estimated
 130 posterior distributions. The idea is to randomly draw a model output X^i from ABC’s sample of
 131 accepted runs as the ‘pseudo-data’ X^0 for a new round of ABC. This does not require further

132 simulation runs, as the original runs can be re-used. The output of this new round of ABC is a set
 133 of accepted runs associated with X^0 . Then, for each parameter j , we calculate the p_j^0 , the
 134 proportion of accepted parameter values smaller than that which produced X^0 . We then repeat the
 135 whole process many times, ending up with a sample of p_j^0 values for each parameter j . Intuitively
 136 these should be spread out between 0 and 1, and not ‘bunched up’ at either the middle or the
 137 extremes of the estimated posteriors. Ideally, the p_j^0 values have a Uniform(0,1) distribution
 138 (Prangle et al. 2013). Algorithm S1 in Appendix S1 gives the coverage algorithm that we first
 139 applied to our earthworm model (van der Vaart et al. 2015); unfortunately this produced non-
 140 uniform coverage for several parameters, motivating the work reported here.

141 1.2. Error-Calibrated ABC

142 In order to improve our estimation procedure, we used Wilkinson’s (2013) version of ABC,
 143 which provides inference for the model given by Equation 1. How to choose π_ϵ , the probability
 144 density function for the error ϵ , is discussed in the next section. Now the acceptance step (2) of
 145 Algorithm 1 is replaced by a probabilistic version, where each (θ^i, X^i) pair is accepted with
 146 probability $\frac{\pi_\epsilon(D-X^i)}{c}$, where $\pi_\epsilon(D - X^i)$ is the probability density function of ϵ evaluated at $D -$
 147 X^i , and c is a constant chosen as the maximum of $\pi_\epsilon(D - X^i)$ (Wilkinson 2013).

148 1.3. Error Estimation

149 If π_ϵ were known it would be straightforward to implement Wilkinson’s (2013) algorithm,
 150 though perhaps slow to produce adequate sample sizes, but in general π_ϵ is not known.
 151 However, if the data come from replicated experiments or time series it is possible to estimate
 152 π_ϵ from the differences between the data and the output of the best-fitting model.

153

154 To do this, we first find \hat{X} , the model output X^i which minimises $\rho(X^i, D)$. When all data are of
 155 the same type, $\rho(X^i, D)$ is the sum of all Euclidean distances between X^i and D . When the data
 156 are of k different types, all Euclidean distances are centered and scaled before summing. For
 157 example, in our earthworm model, where some data points concern growth and others concern
 158 reproduction, all Euclidean distances are centered and scaled by the mean and standard deviation
 159 of all Euclidean distances of that type. This ensures that the overall distance calculation is not
 160 dominated by scale differences between the data types.

161
 162 We then assume that, for each data type, the errors on data points are independent of each other
 163 and drawn from a normal distribution with mean 0, as in classical statistics; this is an assumption
 164 that we discuss in our conclusion. To estimate the standard deviation λ of this normal
 165 distribution, we take the standard deviation $\hat{\lambda}$ of all the $\hat{X}_j - D_j$ values that are of the same type.
 166 So, for example, for the earthworm model, $\hat{\lambda}_{growth}$ is equal to the standard deviation of all
 167 differences between the best-fitting model output \hat{X} and the data D for all data points concerning
 168 growth, and $\hat{\lambda}_{reproduction}$ is equal to the standard deviation of all differences between the best-
 169 fitting model output \hat{X} and the data D for all data points concerning reproduction.

170
 171 Then, under our assumption of independent, normally distributed errors, the probability density
 172 function $\pi_\epsilon(D - X^i) \propto \prod_{j=1}^l \pi_{N(0,1)}\left(\frac{X_j^i - D_j}{\hat{\lambda}_{\tau(j)}}\right)$, where l is the number of data points, $\tau(j)$ is the
 173 type of the j th data point, and $\hat{\lambda}_\tau$ is the standard deviation of data points of type τ . In other
 174 words, the overall acceptance probability of a specific model run i can be calculated by
 175 multiplying the probability densities of each of the simulated data points being produced from

176 the empirical data, given the assumed error distribution. This density is quicker to compute via a
 177 transformation, giving $\pi_\epsilon(D - X^i) \propto \pi_{\chi^2_l}(s) s^{1-\frac{l}{2}}$, where $s = \sum_{j=1}^l \left(\frac{X_j^i - D_j}{\hat{\lambda}_{\tau(j)}} \right)^2$; i.e., the density of a
 178 chi-square distribution with l degrees of freedom evaluated at s , the summed squares of all
 179 normalised errors multiplied by a Jacobian term, $s^{1-\frac{l}{2}}$. Algorithm 2 shows the overall procedure,
 180 which we call ‘error-calibrated ABC’.

1. Repeat n times:
 - a. Draw $\theta^i \sim \pi(\theta)$
 - b. Simulate $X^i \sim \eta(\theta^i)$
2. Find \hat{X} , the simulated value that minimises $\rho(X^i, D)$.
3. For each data type k , calculate $\hat{\lambda}_k$, the standard deviation of all corresponding $\hat{X}_j - D_j$.
4. Accept (θ^i, X^i) with probability $\frac{\pi_{\chi^2_l}(s) s^{1-\frac{l}{2}}}{c}$, where $s = \sum_{j=1}^l \left(\frac{X_j^i - D_j}{\hat{\lambda}_{\tau(j)}} \right)^2$ and c is equal to the maximum acceptance probability across all runs.

181 **Algorithm 2. New ‘error-calibrated ABC’ algorithm.**

182 1.4. Error-Calibrated Coverage

183 Finally, to assess the accuracy of this new algorithm, we update our coverage test, as shown in
 184 Algorithm 3, where $d = 200$, following Prangle et al. (2013). The main change is that the
 185 ‘pseudo-data’ is no longer directly equal to a best-fitting model runs but to a model run plus
 186 estimated noise π_ϵ making it more like the empirical data.

1. Add noise to all simulation results X^i , creating new pairs (θ^i, W^i) . In

particular, add $\mathbf{N}(\mathbf{0}, \hat{\lambda}_{\tau(j)}^2)$ noise to X_j^i to get W_j^i , where $\hat{\lambda}_{\tau(j)}$ is the standard deviation of data points of type $\tau(j)$, i.e. of the same type as the j th data point.

2. For each of the d noisy (θ^i, W^i) that minimise $\rho(W^i, D)$:
 - a. Label as (θ^0, W^0) and do ‘error-calibrated ABC’ with $D = W^0$, using all remaining non-noisy model runs as the simulations:
 - i. Accept each (θ^i, X^i) according to its acceptance probability, using the $\hat{\lambda}$ values calculated in the original analysis.
 - b. For each parameter j :
 - i. Calculate p_j^0 , the sum of all acceptance probabilities with $\theta_j^i \leq \theta_j^0$ divided by the sum of all acceptance probabilities.
3. Plot the distribution of all p_j^0 values, and check for uniformity.

187 **Algorithm 3. New coverage algorithm for ‘error-calibrated ABC’.**

188 1.5. Applications

189 To test this new ‘error-calibrated ABC’, we applied it first to a quadratic model where it is
 190 possible to calculate exact posteriors, and second to our earthworm simulation. In each case, we
 191 compared its results to those of ‘rejection ABC’, where we deterministically accepted the m runs
 192 with the highest acceptance probability according to Algorithm 1. For the quadratic model, the
 193 data consist of observations $D = \theta_1 + \theta_2 x + \theta_3 x^2$ plus noise $\epsilon = \mathbf{N}(\mathbf{0}, 100)$, evaluated for x
 194 values 1, 2, ..., 10 with the true $\theta_1 = -2$, $\theta_2 = 1$ and $\theta_3 = 2$. The simulator η has the same form
 195 without the error; to estimate the values of the θ_1 , θ_2 , and θ_3 parameters we took 10^5 samples of
 196 $[\theta_1, \theta_2, \theta_3]$ where each of θ_1 , θ_2 , and θ_3 were drawn from independent $\mathbf{N}(\mathbf{0}, 9)$ priors. This means

197 that, for this simple example, the simulator is deterministic rather than stochastic. Exact
198 posteriors were calculated using Bayesian regression; see Textbox S1.

199
200 For the earthworms, the observed data D consist of two types: 122 average body masses and 38
201 cocoon productions of earthworms living on experimental laboratory diets. In each case, five to
202 ten earthworms were placed in small containers filled with cattle manure for food (Reinecke and
203 Viljoen 1990, Gunadi et al. 2002, Gunadi and Edwards 2003). The model η is an individual-
204 based model (or IBM) that simulates the growth and reproduction of individual earthworms
205 according to established physiological principles (Sibly et al. 2013). Earthworms wriggle around
206 randomly as they forage, and allocate assimilated energy to maintenance, growth, reproduction,
207 and reserves, in a fixed order of priority; see Johnston et al. (2014). In total the model has
208 fourteen parameters θ , given in Table S1 in Appendix S1. The priors for all parameters were
209 lognormal, with means equal to previously determined literature values (see Johnston et al.
210 (2014)) and standard deviations equal to 0.3536. This produces samples where 95% of the values
211 lie between half and twice the literature values on the unlogged scale. We used ARCHER, the
212 UK's national supercomputing service, to do 10^6 runs; see van der Vaart et al. (2015) for details.

213 1.6. Implementation

214 All ABC code and the quadratic example were implemented in R (R Core Team 2015). The
215 earthworm model was built in NetLogo (Wilensky 1999), and $RNetLogo$ was used to run
216 NetLogo from R (Thiele et al. 2012). All statistical tests were corrected for multiple testing using
217 Holm's method, and all code and simulation results were deposited in a figshare repository.¹

¹ Link to be added before publication.

218 **Results**

219 For the quadratic example, ‘error-calibrated ABC’ estimated the standard deviation of the error,
 220 λ , to be 7.91, producing 330 acceptances. Figure 1A shows the model’s resulting fit (a ‘posterior
 221 predictive check’). The posteriors of all three parameters were not significantly different from
 222 those obtained by exact Bayesian regression (Figure 1B – D), and coverage plots were uniform
 223 (Figure 2A - C), suggesting accurate posteriors. By contrast, for ‘rejection ABC’ with $m = 330$
 224 acceptances, all three posteriors were significantly different from those obtained by exact
 225 Bayesian regression, and coverage plots were ‘U-shaped’, with an excess of p values at the
 226 extremes (Figure S2). After further varying m from 100 to 1000 to 10000, we found that
 227 ‘rejection ABC’ was only accurate for $m = 1000$ (Figure S1 & Figure S2); see Figure 2D – F.

228

229 For the earthworms, ‘error-calibrated ABC’ initially accepted only the best-fitting run, which is
 230 necessarily accepted (see Algorithm 2). Using this best-fitting run, we verified that the error
 231 distributions were normal for both masses and cocoons (Figure S3), and we estimated their
 232 standard deviations to be 0.08 and 10.4 respectively. To increase the number of acceptances, we
 233 fixed λ_{mass} and $\lambda_{cocoons}$ at their original values, but otherwise reduced the data set to every 6th
 234 point; see the Discussion for rationale. Now, 108 runs were accepted, giving the posterior
 235 predictive check of Figure 3. Relative to the priors, 4 out of 14 posteriors were significantly
 236 narrowed (Figure S4), and coverage was uniform for all 14 (Figure S5).

237

238 In comparison, ‘rejection ABC’ with $m = 100$ acceptances narrowed five posteriors (Figure S6).
 239 For h , the half saturation coefficient, IG_m , the maximum ingestion rate, and M_m , the maximum
 240 mass, these posteriors were significantly different from those of ‘error-calibrated ABC’ (

241 Figure 4). IG_m and M_m , along with three other parameters, also produced non-uniform coverage,
242 Figure S7. After varying m from 100 to 10^3 , 10^4 and 10^5 , we found that ‘rejection ABC’ never
243 produced uniform coverage for all parameters at once (Figure S8), with E , the activation energy,
244 for example, varying from ‘U-shaped’ at $m = 100$ to ‘mountain shaped’ at $m = 10^5$.

245 **Discussion**

246 We have shown how incorporating estimation of error into the ABC protocol can improve
247 estimates of parameter values and their credible intervals. To do this we specified ABC
248 acceptance probabilities for the case that errors are normally distributed and independent. Our
249 ‘error-calibrated ABC’ implements a general methodology introduced by Wilkinson (2013). To
250 diagnose the accuracy of our method, we updated Prangle et al.’s (2013) coverage test by adding
251 the estimated error to the simulation runs used as ‘pseudo-data’, improving their realism.

252
253 For our two example models, ‘error-calibrated ABC’ appears to have improved posterior
254 accuracy: Coverage plots were uniform for all parameters, and for the quadratic case, results
255 were indistinguishable from those of exact Bayesian regression. In both cases, ‘rejection ABC’
256 with an equivalent number of acceptances was demonstrably inaccurate. For the quadratic model,
257 this could be corrected by accepting more runs, but for the earthworm IBM, ‘rejection ABC’
258 never produced uniform coverage for all parameters simultaneously. Thus, we conclude that
259 ‘error-calibrated ABC’ offers a real improvement with respect to model calibration.

260
261 In essence, coverage checks for inaccuracies in ABC’s posteriors by repeatedly applying the
262 ABC protocol to ‘pseudo data’ for which the correct parameter values are known. Typically, a
263 lack of uniformity can then be due to either error or inadequacy in the ABC protocol; most

264 notably, an incorrect acceptance rate. A standard coverage test assumes that the model is perfect,
265 and that calibrated correctly, it can replicate the data exactly. However, our updated coverage
266 test drops this assumption, by adding ‘noise’ drawn from the error model to all data points before
267 using them as ‘pseudo-data’. In a coverage test, surpluses in the tails of the coverage distribution,
268 as in Figure 2D, imply that posteriors are too narrow, with too few runs accepted. At the other
269 extreme, deficits in the tails of the coverage distribution, as in Figure 2F, imply that posteriors
270 are too wide, with too many runs accepted. For this polynomial example, we know the error
271 model is correct, so any lack of uniformity must be due to problems with the acceptance criteria.

272
273 For the earthworm model, the use of ‘error-calibrated ABC’ required two approximations:
274 Firstly, it seems unlikely that its errors really are independent across observations and normally
275 distributed. However, this assumption has often been made by ecologists deploying regression
276 models, and would seem as justifiable here. For the future, it would be interesting to explore
277 methods that incorporate correlations between successive errors, since these could reduce the
278 degrees of freedom and so increase acceptance rates. Currently, as our second approximation, we
279 had to remedy a lack of acceptances by reducing the data set to every 6th data point. As the error
280 distribution π_{ϵ} is multivariate normal with dimension equal to the number of data points,
281 acceptance falls off exponentially as the number of data points increases. “Too much data” is a
282 common problem in ABC, known as ‘the curse of dimensionality’. It is generally addressed by
283 summarizing data sets into as few as one or two ‘summary statistics’ (see, e.g., Blum et al.
284 2013). Addressing the issue by ‘thinning out’ a time series, as here, is not an established
285 technique but has the same fundamental justification. While simple, it appears to work well in
286 this case; visually the accepted runs still mimic the full data set (Figure 3), and for ‘rejection

287 ABC', the posteriors estimated with the full and reduced data sets are similar (Figure S9).

288

289 Our overall approach is relatively simple, and does not make use of various sophistications
290 already present in the literature. These include techniques for 'correcting' accepted parameter
291 values on the basis of the resulting model fit, for example using regression (Beaumont et al.
292 2002), by estimating the error simultaneously with a model's parameters, as in $ABC\mu$ (Ratmann
293 et al. 2009), by analysing time series data sequentially (Jasra 2015), or by sampling a model's
294 parameters more efficiently, as in MCMC-ABC (Marjoram et al. 2003) or SMC-ABC (Sisson et
295 al. 2007). We see the simplicity of 'error-calibrated ABC' as an attraction; more efficient
296 sampling schemes are harder to implement and make it impossible to re-use runs for the purpose
297 of calculating coverage. In these cases, 'error-calibrated ABC' offers an accessible approach to
298 improving models' posteriors, with the additional benefit of explicitly accounting for error.

299

300 **Acknowledgements**

301 This work used the ARCHER UK National Supercomputing Service (<http://www.archer.ac.uk>)
302 and was supported by NERC grant number NE/K006282/1. DP was supported by a Richard
303 Rado postdoctoral fellowship from the University of Reading during much of this project. We
304 thank A Meade for computational assistance, S Watson and the University of Reading's
305 Bayesian reading group for discussion, and PJ van Leeuwen, R Everitt, M Beaumont, M
306 Kosmala, J Barber and two anonymous referees for very helpful comments on the manuscript.

307 **References**

308 Beaumont, M. A., W. Zhang, and D. J. Balding. 2002. Approximate Bayesian computation in
309 population genetics. *Genetics* **162**:2025 - 2035.

- 310 Blum, M. G. B., M. A. Nunes, D. Prangle, and S. A. Sisson. 2013. A comparative review of
311 dimension reduction methods in approximate Bayesian computation. *Statistical Science*
312 **28**:189-208.
- 313 Campbell, K. 2006. Statistical calibration of computer simulations. *Reliability Engineering &*
314 *System Safety* **91**:1358-1363.
- 315 Gelman, A., C. J.B., H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. 2013. *Bayesian*
316 *Data Analysis*. 3rd edition. Chapman & Hall/CRC.
- 317 Goldstein, M., and J. Rougier. 2009. Reified Bayesian modelling and inference for physical
318 systems. *Journal of Statistical Planning and Inference* **139**:1221-1239.
- 319 Gunadi, B., C. Blount, and C. A. Edwards. 2002. The growth and fecundity of *Eisenia fetida*
320 (Savigny) in cattle solids pre-composted for different periods. *Pedobiologia* **46**:15-23.
- 321 Gunadi, B., and C. A. Edwards. 2003. The effects of multiple applications of different organic
322 wastes on the growth, fecundity and survival of *Eisenia fetida* (Savigny) (Lumbricidae).
323 *Pedobiologia* **47**:321-329.
- 324 Hartig, F., J. M. Calabrese, B. Reineking, T. Wiegand, and A. Huth. 2011. Statistical inference
325 for stochastic simulation models - theory and application. *Ecology Letters* **14**:816-827.
- 326 Higdon, D., J. Gattiker, B. Williams, and M. Rightley. 2008. Computer model calibration using
327 high-dimensional output. *Journal of the American Statistical Association* **103**:570-583.
- 328 Jasra, A. 2015. Approximate Bayesian Computation for a class of time series models.
329 *International Statistical Review* **83**:405-435.
- 330 Johnston, A. S. A., M. E. Hodson, P. Thorbek, T. Alvarez, and R. M. Sibly. 2014. An energy
331 budget agent-based model of earthworm populations and its application to study the
332 effects of pesticides. *Ecological Modelling* **280**:5-17.

- 333 Kosmala, M., P. Miller, S. Ferreira, P. Funston, D. Keet, and C. Packer. 2016. Estimating
334 wildlife disease dynamics in complex systems using an Approximate Bayesian
335 Computation framework. *Ecological Applications* **26**:295-308.
- 336 Lagarrigues, G., F. Jabot, V. Lafond, and B. Courbaud. 2015. Approximate Bayesian
337 computation to recalibrate individual-based models with population data: Illustration with
338 a forest simulation model. *Ecological Modelling* **306**:278-286.
- 339 Lethbridge, M. R., and J. C. Strauss. 2015. A novel dispersal algorithm in individual-based,
340 spatially explicit Population Viability Analysis: A new role for genetic measures in
341 model testing? *Environmental Modelling & Software* **68**:83-97.
- 342 Marjoram, P., J. Molitor, V. Plagnol, and S. Tavaré. 2003. Markov chain Monte Carlo without
343 likelihoods. *Proceedings of the National Academy of Sciences of the United States of*
344 *America* **100**:15324-15238.
- 345 Nabe-Nielsen, J., R. M. Sibly, J. Tougaard, J. Teilmann, and S. Sveegard. 2014. Effects of noise
346 and by-catch on a Danish harbour porpoise population. *Ecological Modelling* **272**:242-
347 251.
- 348 Nehrbass, N., and E. Winkler. 2007. Is the Giant Hogweed still a threat? An individual-based
349 modelling approach for local invasion dynamics of *Heracleum mantegazzianum*.
350 *Ecological Modelling* **201**:377-384.
- 351 Prangle, D., M. G. B. Blum, G. Popovic, and S. A. Sisson. 2013. Diagnostic tools for
352 approximate Bayesian computation using the coverage property. *Australian & New*
353 *Zealand Journal of Statistics* **56**:309-329.

- 354 Pritchard, J. K., M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman. 1999. Population growth
355 of human Y chromosomes: A study of Y chromosome microsatellites. *Molecular Biology*
356 and Evolution **16**:1791-1798.
- 357 R Core Team. 2015. R: A Language and Environment for Statistical Computing. The R
358 Foundation for Statistical Computing, Vienna, Austria.
- 359 Rasmussen, R., and G. Hamilton. 2012. An approximate Bayesian computation approach for
360 estimating parameters of complex environmental processes in a cellular automata.
361 *Environmental Modelling & Software* **29**:1-10.
- 362 Ratmann, O., C. Andrieu, C. Wiuf, and D. M. Richardson. 2009. Model criticism based on
363 likelihood-free inference, with an application to protein network evolution. *Proceedings*
364 of the National Academy of Sciences **106**:10576-10581.
- 365 Reinecke, A. J., and S. A. Viljoen. 1990. The influence of feeding patterns on growth and
366 reproduction of the vermicomposting earthworm *Eisenia fetida* (Oligochaeta). *Biology*
367 and Fertility of Soils **10**:184-187.
- 368 Sibly, R. M., V. Grimm, B. T. Martin, A. S. A. Johnston, K. Kułakowska, C. J. Topping, P.
369 Calow, J. Nabe-Nielsen, P. Thorbek, and D. L. DeAngelis. 2013. Representing the
370 acquisition and use of energy by individuals in agent-based models of animal
371 populations. *Methods in Ecology and Evolution* **4**:151-161.
- 372 Sisson, S. A., Y. Fan, and M. A. Tanaka. 2007. Sequential Monte Carlo without likelihoods.
373 *Proceedings of the National Academy of Sciences of the United States of America*
374 **104**:1760-1765.
- 375 Tavaré, S., D. J. Balding, R. C. Griffiths, and P. Donnelly. 1997. Inferring coalescence times
376 from DNA sequence data. *Genetics* **145**:505-518.

- 377 Thiele, J. C., W. Kurth, and V. Grimm. 2012. RNetLogo: An R package for running and
378 exploring individual-based models implemented in NetLogo. *Methods in Ecology and*
379 *Evolution* **3**:480-483.
- 380 van der Vaart, E., M. A. Beaumont, A. Johnston, and R. M. Sibly. 2015. Calibration and
381 evaluation of individual-based models using Approximate Bayesian Computation.
382 *Ecological Modelling* **312**:182-190.
- 383 West, A. D., J. Goss-Custard, R. A. Stillman, R. W. G. Caldow, S. E. A. L. D. Durell, and S.
384 McGrorty. 2002. Predicting the impacts of disturbance on shorebird mortality using a
385 behaviour-based model. *Biological Conservation* **106**:319-328.
- 386 Wilensky, U. 1999. NetLogo. Center for Connected Learning and Computer-Based Modeling,
387 Northwestern University, Evanston, IL.
- 388 Wilkinson, R. D. 2013. Approximate Bayesian Computation (ABC) gives exact results under the
389 assumption of model error. *Statistical Applications in Genetics and Molecular Biology*
390 **12**:129-141.
- 391 Zurell, D., V. Grimm, E. Rossmannith, N. Zbinden, N. E. Zimmerman, and B. Schröder. 2012.
392 Uncertainty in predictions of range dynamics: Black grouse climbing the Swiss Alps.
393 *Ecography*:590-603.

394

395 **Figure Legends**

396 **Figure 1. Results for the quadratic example.** A: Posterior check. Black points represent the
397 data, the result of $\theta_1 + \theta_2 x + \theta_3 x^2$ plus $N(\mathbf{0}, 100)$ noise, and the semi-transparent grey lines are
398 the ‘posterior predictive check’, i.e., 100 random samples from runs accepted by ‘error-calibrated

399 ABC'. B – D: Posterior distributions. Bars are 'error-calibrated ABC', lines are exact Bayesian
 400 regression, all differences nonsignificant (Kolmogorov-Smirnov, $p > 0.01$). The true θ_1 , θ_2 and θ_3
 401 were -2, 1 and 2, respectively, marked on the x-axes by arrows; posteriors are centred differently
 402 because of the added noise and the priors that were used. On the horizontal axes, ticks are placed
 403 at the mean of the exact posterior density and three standard deviations above and below.

404

405 **Figure 2. Coverage for the quadratic example.** A - C: 'Error-calibrated ABC'. D – F:
 406 'Rejection ABC' for parameter θ_3 at different acceptance rates m . Asterisks mark significant
 407 departures from uniformity (Kolmogorov-Smirnov, $p < 0.01$).

408

409 **Figure 3. Body masses and cocoon productions in the earthworm experiments.** The black
 410 lines show the empirical data (Reinecke and Viljoen 1990, Gunadi et al. 2002, Gunadi and
 411 Edwards 2003), the thick grey line is the 'best-fitting run' and the semi-transparent grey lines are
 412 the 'posterior predictive check', i.e., the output of 100 new simulations using random samples
 413 from runs accepted by 'error-calibrated ABC'. Only every 6th data point, marked by a circle, was
 414 used in the analysis; those marked by a cross were removed to improve acceptance rates. Arrows
 415 indicate when food was added (\uparrow) or removed (\downarrow). See van der Vaart et al. (2015) for details.

416

417 **Figure 4. Posterior distributions for the earthworm model.** Black lines show 'error-calibrated
 418 ABC' accepting 108 runs; grey lines 'rejection ABC' accepting 100. Circles represent medians,
 419 whiskers 95% credible intervals. Asterisks mark significant differences (Kolmogorov-Smirnov, p
 420 < 0.01). All parameter values were scaled by dividing by the corresponding literature value.