

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

Testing peatland testate amoeba transfer functions: appropriate methods for clustered training-sets.

Richard J. Payne^{1,2*}, Richard J. Telford^{3*}, Jeffrey J. Blackford², Antony Blundell⁴, Robert K. Booth⁵, Dan J. Charman⁶, Łukasz Lamentowicz⁷, Mariusz Lamentowicz⁸, Edward A.D. Mitchell⁹, Genevieve Potts², Graeme T. Swindles⁴, Barry G. Warner¹⁰ and Wendy Woodland¹¹

¹School of Science and the Environment, Manchester Metropolitan University, Chester Street, Manchester, M1 5GD, UK.

²Geography, School of Environment and Development, The University of Manchester, Oxford Road, Manchester, M13 9PL, UK.

³Bjerknes Centre for Climate Change and Department of Biology, University of Bergen, Post box 7820, N-5020 Bergen, Norway.

⁴School of Geography, University of Leeds, Leeds, LS2 9JT, UK.

⁵Earth and Environmental Sciences, Lehigh University, 31 Williams Drive, Bethlehem, PA 18015, USA

⁶Department of Geography, College of Life and Environmental Sciences, University of Exeter, Exeter, Devon, EX4 4RJ, UK

⁷Department of Hydrobiology, Adam Mickiewicz University, Umultowska 89, 61-614 Poznan, Poland.

⁸Dept. of Biogeography & Palaeoecology, Adam Mickiewicz University, Dziegielowa 27, PL-61-680 Poznan, Poland.

33 ⁹Laboratoire de biologie du sol, Université de Neuchâtel, Rue Emile-Argand 11, Case
34 postale 158, 2009 Neuchâtel, Switzerland.

35

36 ¹⁰ Department of Earth and Environmental Sciences, University of Waterloo, 200
37 University Avenue West, Waterloo, Ontario, N2L 3G1, Canada

38

39 ¹¹ Department of Geography & Environmental Management, School of the Built and
40 Natural Environment, University of the West of England, Coldharbour Lane, Bristol,
41 BS16 1QY, UK.

42

43 * Joint first authors, e-mails: r.payne@mmu.ac.uk, richard.telford@bio.uib.no

44

45 ABSTRACT

46

47 Transfer functions are widely used in palaeoecology to infer past environmental
48 conditions from fossil remains of many groups of organisms. In contrast to
49 traditional training-set design with one observation per site, some training sets,
50 including those for peatland testate amoeba-hydrology transfer functions, have a
51 clustered structure with many observations from each site. Here we show that this
52 clustered design causes standard performance statistics to be overly optimistic.
53 Model performance when applied to independent data sets is considerably weaker
54 than suggested by statistical cross-validation. We discuss the reasons for these
55 problems and describe leave-one-site-out cross-validation and the cluster bootstrap
56 as appropriate methods for clustered training sets. Using these methods we show
57 that the performance of most testate amoeba-hydrology transfer functions is worse
58 than previously assumed and reconstructions are more uncertain.

59

60 KEYWORDS: Transfer functions; Palaeoclimate; Clustered data; Leave-one-site-out
61 cross-validation, Cluster bootstrap.

62

63

64 Transfer functions are widely used to generate quantitative environmental
65 reconstructions in palaeoecology. Traditional training-set design (e.g. Birks et al.
66 1990) has one observation per site. An alternative design with many observations at
67 each site is used for some training-sets, including those for chironomid-lake depth
68 (Kurek and Cwynar 2009); coastal diatom-water chemistry (Saunders et al. 2008);
69 diatom- and foraminifera-sea level (Massey et al. 2006; Zong & Horton 1999; Leorri
70 et al. 2008); and testate amoeba-hydrology transfer functions (Charman 2001,
71 Mitchell et al. 2008). Although the implications of, and methods for, such clustered
72 data are well known in other branches of statistics (Walsh 1947), the implications of
73 this design have been neglected for transfer functions.

74 One motivation for developing clustered training-sets is the presence within
75 each site of substantial environmental gradients, which may be large relative to the
76 differences between sites. This contrasts with the traditional one observation per
77 site training-set where typically the environmental variable (e.g. lake-pH) is assumed
78 to be spatially homogeneous at each site. Standard methods for assessing the
79 performance of transfer functions assume that the observations are independent
80 and are thus inappropriate for clustered data. Lack of independence between
81 observations, either because of spatial autocorrelation or a clustered design, will
82 cause performance statistics to be over-optimistic (Telford and Birks, 2005). Telford
83 and Birks (2009) have developed cross-validation methods appropriate for spatially
84 autocorrelated training sets; here we consider the problem of clustered training sets
85 and develop appropriate cross-validation methods. We focus on testate amoeba-
86 hydrology transfer functions from peatlands, which have become increasingly
87 important in shaping our understanding of Holocene climatic change (Charman et al.
88 2004, 2006).

89

90 Indications that standard tools are misleading

91 Training sets for peatland testate amoebae transfer functions have a highly
92 uneven spatial structure, with samples from individual sites often only separated by
93 a few metres, while sites may be separated by tens or hundreds of kilometres.
94 Ordinations of testate amoeba data frequently show distinct clustering of
95 observations from the same bog (e.g. Charman et al. 2007, Swindles et al. 2009) and

96 site identity typically explains a large proportion of variance in constrained
97 ordinations (Fig. 1).

98 To provide an independent estimate of transfer function performance, we
99 apply five transfer functions to all comparable independent datasets with
100 appropriate corrections for taxonomic and methodological differences (Appendix I).
101 Table 1 shows that most transfer functions perform worse than suggested by leave-
102 one-out (LOO) cross-validation when applied to independent data. Methodological
103 explanations for the poor model performance can largely be excluded. Differences in
104 time-discrete water-table measurements cannot explain the differences in rank-
105 order shown by Spearman's ρ . Any differences in sample preparation and analysis, or
106 residual taxonomic biases cannot explain poor performance where these are closely
107 harmonised (e.g. Polish data). Performance is particularly poor for two datasets from
108 Scotland (Payne 2010a; Potts & Blackford unpublished data); in the case of the Moss
109 of Achnacree, this is likely to be due to the limited WTD range in a site which has
110 experienced hydrological modification. As previously presented tests with transfer
111 functions from different regions have frequently (Charman et al. 2007; Booth et al.
112 2008; Payne 2011), but not universally (e.g. Swindles et al. 2009), shown
113 performance poorer than LOO cross-validation we conclude that model performance
114 *in praxis* appears to be weaker than suggested by conventional cross-validation.

115

116 Appropriate cross-validation methods for clustered data

117 Typically, transfer function model performance is assessed by either leave-
118 one-out (LOO) or bootstrap cross-validation. In LOO, one observation at a time is
119 omitted from the training-set of size n and the environmental value predicted using
120 the remaining $n-1$ observations. For clustered data, this can be extended to leave-
121 one-site-out cross-validation (LOSO), where data from one site is omitted from the
122 training set, and data from the remaining $m-1$ sites used to predict it. LOSO is also
123 known as leave-one-cluster-out cross-validation and sometimes as leave-one-group-
124 out cross-validation (confusingly, this latter term is also used to refer to k -fold cross-
125 validation in which k groups are created at random).

126 In standard bootstrap cross-validation, n observations are selected from the
127 training set with replacement, and used to predict the remaining observations and

128 new observations. There are several possible bootstrap schemes available for
129 clustered data including the cluster bootstrap, where m clusters are selected at
130 random with replacement, and the two-level bootstrap where m clusters are
131 selected at random and observations are selected at random from within each
132 cluster (Field and Welsh 2007). Here we use the cluster bootstrap following the
133 findings of Field and Welsh (2007) that the two-level bootstrap and the related
134 reverse-two-level bootstrap generate excessive variability.

135

136 Application to Testate Amoeba Training sets

137 We determine the performance of 14 published testate amoeba transfer
138 functions for water-table depth (WTD) using both robust cross-validation methods
139 and standard methods. In the case of the Jura training set (Mitchell et al. 1999) we
140 omit samples with estimated rather than measured water-table depths. For all
141 training sets, we use weighted averaging with inverse deshrinking as this transfer
142 function method is fairly robust to spatial autocorrelation (Telford and Birks, 2005)
143 and so should also be fairly robust to clustered data. Assemblage data were square
144 root transformed prior to analysis. All analyses were carried out in R (R Development
145 Core Team 2010) with the rioja library (Juggins 2010).

146 While differences are not always great, all transfer functions except for one
147 exhibit worse performance with LOSO than LOO cross-validation (Table 2). One
148 transfer function has an LOSO RMSEP greater than the standard deviation of WTD.
149 There are several possible reasons for this deterioration in performance. It could be
150 simply an artefact because the estimates are based on fewer observations as more
151 observations are omitted during LOSO than LOO. We tested for the importance of
152 this factor by running a modified cross-validation scheme termed leave-many-out
153 (LMO) that omits as many observations as LOSO when making each prediction but
154 with the observations chosen at random rather than being from the same site. We
155 repeated this analysis 100 times to get a distribution of performance statistics and
156 tested if the observed LOSO RMSEP is worse than the 95th percentile of the leave-
157 many-out RMSEP. Only the Poland (Lamentowicz & Mitchell 2005) training set had a
158 LOSO performance that was not statistically significantly worse than expected from
159 leaving out so many observations during cross-validation.

160 LOSO performance would be worse than LOO performance if each site only
161 covered part of the environmental gradient. This factor is likely to be of minor
162 importance, except in the Greece training set as all the other training sets have
163 replication along the WTD gradient and variance partitioning shows only a small
164 covariance between WTD and site for most of the training sets (Figure 1).

165 As for most training sets the WTD measurements are based on one-time spot
166 measurements, there may be site-specific errors in the WTD measurements if heavy-
167 rainfall or prolonged drought occurs between sampling the first and last bog. Most
168 training sets were collected within a short period of time, so major changes in WTD
169 are unlikely to have occurred however a few training sets were acquired over a
170 longer period of time and this may be an important factor (Charman et al. 2007;
171 Lamentowicz et al. 2008b).

172 There are likely to be important non-hydrological controls on amoebae which
173 differ between sites such as pollutant loading with recent studies showing sulphur
174 (Payne et al. 2010), reactive nitrogen (Nguyen-Viet et al. 2004; Mitchell 2004), heavy
175 metals (Nguyen-Viet et al. 2007; 2008) and particulate matter (Meyer et al. 2010) to
176 be important. Many transfer function studies have included sites of differing pH and
177 trophic status, and there is evidence for differences in amoeba communities and
178 their hydrological responses between fens and bogs (Payne 2011; Jassey et al. 2011).
179 Plant communities, which differ between sites in many studies, shape both the
180 physical and biotic environment of amoebae through processes such as root
181 exudation and allelopathy, particularly the production of phenolic compounds
182 (Jassey et al. 2011). The fundamental hydrological controls on amoeba communities
183 are poorly understood, while water table depth consistently explains the largest
184 proportion of variance in gradient studies it is clearly not water table depth *per se*
185 which is important to amoebae usually living well above the water table. Water table
186 depth is simply a robust measurement, which serves as a proxy for the hydrological
187 variables which do affect amoebae such as water film thickness and variability in the
188 top few cm of moss where amoebae live (Sullivan et al. 2011). These variables may
189 be controlled by fine-scale structural details of the peat and plant communities.

190

191 Predictors of LOSO relative performance

192 In an attempt to understand the attributes of training sets that have a large
193 decrease in performance with LOSO cross-validation, we regress the decrease in
194 performance, standardised by dividing by the standard deviation of WTD, against the
195 number of sites and observations, the proportion of variance explained by WTD, site,
196 and the covariance between WTD and site (Fig. 2). Of these predictors, only the
197 proportion of variance explained by WTD is a statistically significant predictor of the
198 deterioration in performance. Although the regression is not statistically significant,
199 there appears to be an increased risk of a large reduction in performance for training
200 sets with few sites.

201

202 Error decomposition

203 The magnitude of the RMSEP is not necessarily a good guide to the utility of a
204 transfer function. If, as is usually the case in testate amoeba palaeoecology, one is
205 interested only in identifying relatively wet and dry phases, then the absolute value
206 of the reconstruction is not very important. Thus, even transfer functions with a
207 large RMSEP could potentially have utility.

208 For each site in the clustered training-set, we can decompose the total sum of
209 squares of residuals into the proportion explained by site-specific offsets or biases
210 and the residual variation. Table 3 shows that when LOSO is used instead of LOO, the
211 site specific offset increases much more than the residual variation in both absolute
212 and relative terms. This suggests that the absolute values of reconstructions are
213 much more uncertain, but the relative values are only slightly more uncertain than
214 LOO suggests.

215

216 Reconstruction errors

217 Sample-specific (s_1 ; Birks et al., 1990; Birks, 1995) bootstrap errors for the
218 cluster bootstrap will always be larger than those from the standard bootstrap. Fig. 3
219 shows the WTD reconstruction for Jelenia Wyspa, Poland (Lamentowicz et al. 2007b)
220 using the Poland 2008 training set, with sample-specific bootstrap errors using both
221 bootstrap techniques. Bootstrap errors vary by sample but are in all cases greater
222 when using the cluster bootstrap and for some samples the errors are more than
223 double.

224

225 Recommendations

226 Given our results, improvements can be made in both the generation and
227 application of clustered training sets. We make four recommendations for
228 generating new training sets, which should be followed where it is practical to do so
229 and may not be possible to satisfy simultaneously. First, efforts should be made to
230 sample the full environmental gradient at each site, or at least to ensure that all
231 parts of the gradient are replicated in several sites. Ideally, the gradients should be
232 uniformly sampled at each site (Telford and Birks 2011). Second, approximately the
233 same number of observations should be made at each site, so that in LOSO cross-
234 validation the number of observations omitted is close to constant. Third, a large
235 number of sites should be sampled, as the cluster bootstrap is not appropriate for
236 datasets with few clusters. Finally, the sites should be similar to each other with
237 respect to, for example, vegetation and climate, with the proviso that care is taken
238 to include sufficient diversity of sites to ensure that all fossil samples have good
239 analogues in the training set.

240 We recommend that the robust cross-validation methods developed here are
241 used when testing the performance of clustered training sets. We anticipate that the
242 performance statistics of transfer function methods robust to autocorrelation (e.g.,
243 WA) will deteriorate less with robust cross-validation than methods more sensitive
244 to autocorrelation (e.g., WAPLS with several components). If there is a choice of
245 training set that could be applied to the fossil data, we recommend, all else being
246 equal, using the training set with the smallest loss of performance when robust
247 cross-validation is used. Single-site training sets (e.g. Booth et al. 2008; Payne et al.
248 2008) will be immune to cluster problems but this may be offset by poor
249 reconstructive ability. As always in quantitative palaeoecology, caution should be
250 used in interpreting small changes in reconstructions and replication using multi-
251 core, multi-proxy and multi-site records is desirable.

252

253 Conclusions

254 Published performance statistics of testate amoeba transfer functions are
255 over-optimistic due to the clustered design of the training sets. LOO cross-validation

256 is biased by the lack of independence of the observations. As amoeba communities
257 in a sample tend to be more similar to other samples from the same site than to
258 samples from different sites, if samples from the same site remain in the training set
259 during cross-validation, then the model will generate unrealistically accurate
260 predictions of water-table depth in the training set.

261

262

263

264

265

266

267

268

269

270

271

272

273

274 ACKNOWLEDGEMENTS

275

276 RJP was supported by a Humanities Research Fellowship from the University of
277 Manchester and a Study Grant from the British Institute at Ankara. Norwegian
278 Research Council projects ARCTREC and PES helped support RJT. We thank H.J.B.
279 Birks for his comments on this manuscript. R-code for leave-one-site-out and cluster
280 bootstrap cross-validation has been implemented in the rioja library. This is
281 publication no. A358 from the Bjerknes Centre for Climate Research

282 Author contributions:

283 RJP conceived and coordinated the project, compiled the data and carried out the
284 tests with independent data-sets. RJT devised and implemented the cross-validation
285 procedures. RJP and RJT wrote the paper. Other authors contributed data, discussed
286 the taxonomic harmonisation issues and commented on the interpretation of the
287 results and manuscript.

289

290 TABLES

291 Table 1. Transfer function performance for five training sets tested by leave-one-out
 292 (LOO) cross-validation and application to independent test-sets, showing transfer
 293 function method used, number of samples (*n*), root mean squared error of
 294 prediction (RMSEP), R^2 , and Spearman's ρ . Some values differ from previously
 295 published values due to minor variation in sample selection and taxonomic
 296 harmonisation. Values in round brackets show performance when small taxa are
 297 excluded to account for differences in the use of back-sieving (Appendix 1). R^2 and ρ
 298 values in square brackets denote negative correlations.

299

Training-set	Transfer function	Test-set	Peatland type(s)	N	RMSEP (cm)	R^2	P
European (Charman et al. 2007)	2 component WA-PLS	LOO cross-validation	-	119	5.63 (5.80)	0.71 (0.69)	0.90 (0.89)
		All test data	Bogs	200	5.51	0.18	0.67
		Blythermo (Potts & Blackford, unpublished) ²	Bog	9	11.40	0.37	0.66
		Loonan (Potts & Blackford, unpublished) ²	Bog	11	13.02	[0.12]	[-0.38]
		Moss of Achnacree (Payne 2010a) ^{1,2}	Bog	30	6.65	[0.01]	[-0.01]
		Moidach More (Payne et al. 2010b) ¹	Bog	150	4.38	0.53	0.75
UK (Woodland et al. 1998)	WA-Tol (inverse deshrinking)	LOO cross-validation	-	160	3.94 (3.91)	0.29 (0.30)	0.64 (0.64)
		All test data	Bogs	200	6.71	0.25	0.60
		Blythermo (Potts & Blackford, unpublished) ²	Bog	9	13.18	0.56	0.82
		Loonan (Potts & Blackford, unpublished) ²	Bog	11	17.05	[0.13]	[-0.21]
		Moss of Achnacree (Payne 2010a) ^{1,2}	Bog	30	10.19	0.01	0.11
		Moidach More (Payne et al. 2010b) ¹	Bog	150	4.86	0.23	0.42
Alaska (Payne et al. 2006)	2 component WA-PLS	LOO cross-validation	-	91	9.99	0.53	0.81
		Alaska (Markel et al. 2010)	Various	126	16.52	0.42	0.61
Alaska (Markel et al. 2010)	2 component WA-PLS	LOO cross-validation	-	126	8.50	0.63	0.84
		Alaska (Payne et al. 2006)	Various	91	16.94	0.42	0.69
Poland (Lamentowicz & Mitchell 2005)	WA-Tol (inverse deshrinking)	LOO cross-validation	-	36	7.75	0.72	0.94
		All test data	Various	213	11.23	0.20	0.48
		Jedwabna (Lamentowicz et al. 2008b)	Poor fen	10	5.77	0.17	0.53
		Mietlica (Lamentowicz et al. 2008b)	Poor fen	12	7.86	0.85	0.77
		Ostrowite (Lamentowicz et al. 2008b)	Bog	7	13.41	0.82	0.85
		Rybie Oko (Lamentowicz et al. 2008b)	Bog	16	6.35	0.80	0.84
		Skrzynka (Lamentowicz et al. 2008b)	Poor fen	12	4.13	0.55	0.60
		Stawek (Lamentowicz et al. 2008b)	Poor fen	9	8.69	0.52	0.39
		Stężki (Lamentowicz et al. 2008b)	Moderately rich fen	10	7.89	0.51	0.71

		Żabieniec (Lamentowicz et al. 2008b)	Schwingmoor	8	3.83	0.76	0.96
		Chlebowo (Lamentowicz et al. 2007a, 2008a)	Poor fen	27	5.96	0.27	0.54
		Linje (Lamentowicz et al. 2008b)	Bog and poor fen	46	12.07	0.52	0.55
		Słowińskie Błota (Lamentowicz et al. 2008b)	Bog	25	29.58	0.24	0.73
		Jeziorka Kozie (Lamentowicz et al. 2008b)	Poor fen	31	11.34	0.00	0.27

300

¹Back-sieving not used so small taxa excluded.

301

²Lower counts of around 100 tests.

302

303 Table 2. Root mean squared error of prediction for 14 published training sets
 304 calculated with leave-one-out (LOO), leave-one-site-out (LOSO), and leave-many-out
 305 (LMO) cross-validation. The 95th percentile of the LMO distribution is shown. Results
 306 are based on weighted averaging with inverse deshrinking on square root
 307 transformed data. Also shown are the DWT range (cm), number of sites (*m*) and
 308 observations (*n*), and the standard deviation of WTD (sd).

	Range (cm)	<i>m</i>	<i>n</i>	LOO	LOSO	LMO 95%	sd
Europe (Charman et al. 2007)	-3-35	7	119	6.2	6.9	6.3	10.5
Alaska 2006 (Payne et al. 2006)	7-67	8	91	10.8	14.0	11.1	14.6
Alaska 2010 (Markel et al. 2010)	-18-46	12	126	8.6	9.3	8.8	14.0
Engadine (Lamentowicz et al. 2010)	-20-76	6	84	9.8	11.0	10.3	16.1
Greece (Payne and Mitchell 2007)	-1-14.5	4	57	2.2	3.3	2.2	4.1
Jura (Mitchell et al. 1999)	3-53	4	36	9.5	12.4	10.4	13.4
Minnesota/Ontario (Warner and Charman	0-100	10	49	20.1	22.7	20.8	26.2
Newfoundland (Charman and Warner 1997)	-4-46	6	57	7.2	8.1	7.6	11.8
Northern Ireland (Swindles et al. 2009)	-10-38	3	81	5.3	6.0	5.6	12.2
Rockies (Booth and Zygmunt 2005)	-5-50	14	139	7.5	8.0	7.6	16.1
UK (Woodland et al. 1998)	0-19	9	160	4.0	4.8	4.1	4.7
North America (Booth 2008)	-13-75	31	403	8.1	8.2	8.2	17.1
Poland 2008 (Lamentowicz et al. 2008b)	-25-84	15	249	14.0	16.3	14.1	17.8
Poland 2005 (Lamentowicz and Mitchell 2005)	-3-55	3	36	9.6	9.3	11.8	14.7

309

310

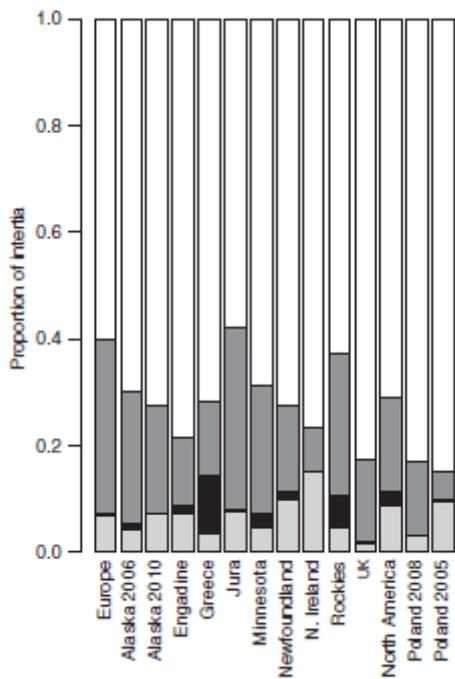
311

312 Table 3. Decomposition of the mean total sum of squares of the transfer function
313 residuals into the portion explained by site-specific offsets and the residual variation
314 for both LOO and LOSO cross-validation, and the ratio of the LOSO and LOO results.

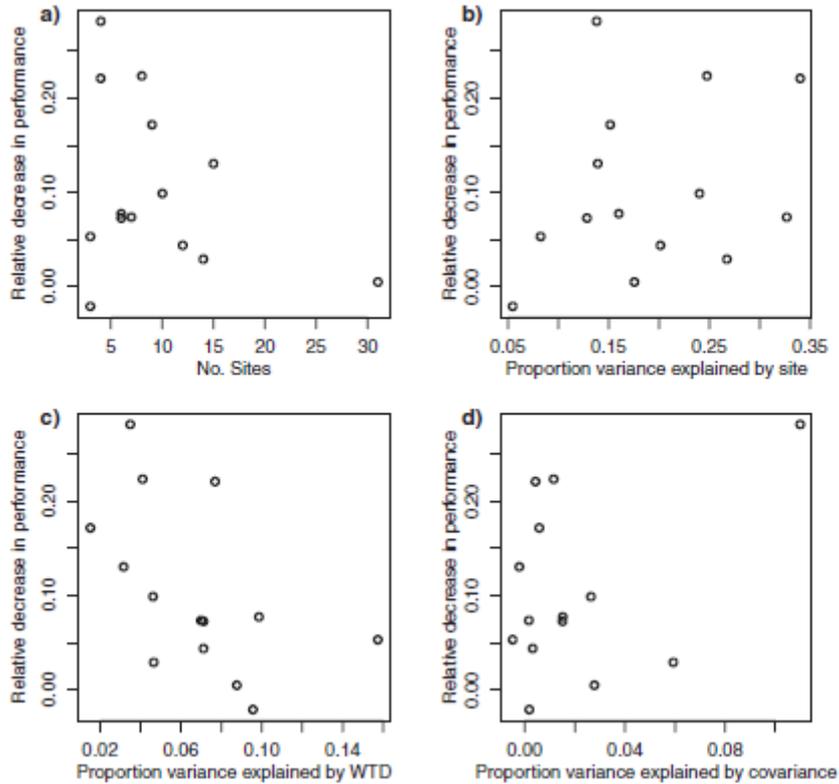
	LOO			LOSO			LOSO/LOO		
	Total	Site	Residual	total	Site	Residual	total	Site	Residual
Europe	38	9	29	48	16	32	1.26	1.89	1.08
Alaska 2006	116	53	63	197	121	75	1.69	2.28	1.19
Alaska 2010	75	13	61	86	25	60	1.14	1.88	0.98
Engadine	96	17	79	120	30	90	1.25	1.72	1.15
Greece	5	2	2	11	8	3	2.35	3.56	1.22
Jura	90	8	82	154	69	85	1.71	8.93	1.04
Minnesota/Ontario	405	177	228	516	250	266	1.27	1.41	1.17
Newfoundland	52	15	37	66	29	37	1.26	1.87	1.01
Northern Ireland	28	5	24	35	9	26	1.25	2.04	1.10
Rockies	57	8	48	64	16	48	1.12	1.95	0.98
UK	16	4	12	23	11	11	1.44	2.74	0.98
North America	66	12	54	68	13	54	1.02	1.12	1.00
Poland 2008	196	72	124	266	134	133	1.36	1.85	1.07
Poland 2005	91	11	80	84	13	71	0.92	1.18	0.88

315

316 Figure 1.
 317 Variance partitioning of the inertia in the different data-sets into components
 318 explained by water table depth (light grey), site (dark grey), covariance between site
 319 and water table depth (black). Unexplained inertia is shown in white. See Table 2 for
 320 data sources. Site is a statistically significant predictor for all training sets except
 321 Poland 2005.

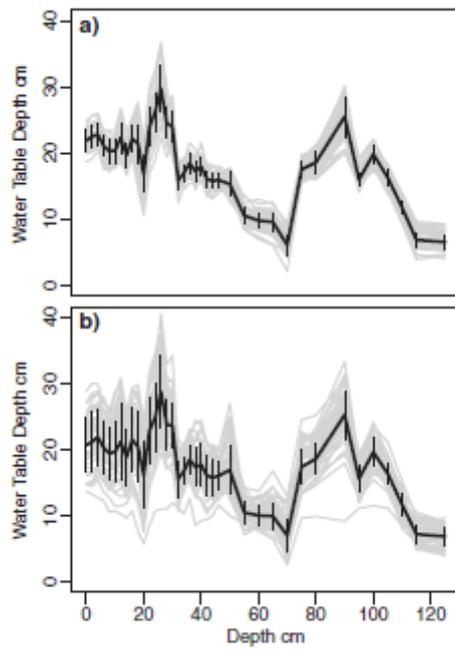


322
 323 Figure 2. Scatter plots of the relative decrease in performance against different
 324 predictors: a) number of sites; and proportion of variance explained by b) site, c)
 325 water table depth and d) covariance between water table depth and site in a CCA.



326

327 Fig. 3. Water table reconstruction from Jelenia Wyspa, Poland (Lamentowicz et al.
 328 2007b) calculated using weighted averaging with inverse deshrinking on square root
 329 transformed data with the expanded Polish training set (Lamentowicz et al. 2008b).
 330 Reconstructions (black) are based on 1000 bootstrap predictions (50 of which are
 331 shown in grey) for a) conventional bootstrap and b) cluster bootstrap. The standard
 332 deviation of the bootstrap predictions (error component s_1) is shown with vertical
 333 black lines).



334

335

336

337

338

339 REFERENCES

340

341 Birks HJB, Line JM, Juggins S, Stevenson AC, ter Braak CJF (1990) Diatoms and pH
342 reconstruction. *Philosophical Transactions of the Royal Society London B* 327:
343 263-278.

344 Birks HJB (1995) Quantitative palaeoenvironmental reconstructions, In: Maddy D and
345 Brew JS (eds.) *Statistical modelling of quaternary science data*. Technical Guide 5.
346 Cambridge: Quaternary Research Association.

347 Booth RK (2008) Testate amoebae as proxies for mean annual water-table depth in
348 *Sphagnum*-dominated peatlands of North America. *Journal of Quaternary Science*
349 23: 43-57.

350 Booth RK, Sullivan ME, Sousa VA (2008) Ecology of testate amoebae in a North
351 Carolina pocosin and their potential use as environmental and
352 palaeoenvironmental indicators. *Ecoscience* 15: 277-289.

353 Booth RK, Zygmunt JR (2005) Biogeography and comparative ecology of testate
354 amoebae inhabiting *Sphagnum*-dominated peatlands in the Great Lakes and
355 Rocky Mountain regions of North America. *Diversity and Distributions* 11: 577-
356 590.

357 Charman DJ (2001) Biostratigraphic and palaeoenvironmental applications of testate
358 amoebae. *Quaternary Science Reviews* 20, 1753-1764.

359 Charman DJ, Hendon D (2000) Long-term changes in soil water tables over the past
360 4500 years: relationships with climate and North Atlantic atmospheric circulation
361 and sea surface temperature. *Climatic Change* 47: 45-59

362 Charman DJ, Warner BG (1992) Relationship between testate amoebae (Protozoa:
363 Rhizopoda) and microenvironmental parameters on a forested peatland in
364 northeastern Ontario. *Canadian Journal of Zoology* 70: 2474-2482.

365 Charman DJ, Warner BG (1997) The ecology of testate amoebae (Protozoa:
366 Rhizopoda) and microenvironmental parameters in Newfoundland, Canada:
367 modeling hydrological relationships for palaeoenvironmental reconstruction.
368 *Ecoscience* 4: 555-562.

369 Charman DJ, Brown A, Hendon D, Karofeld E (2004) Testing the relationship between
370 Holocene peatland palaeoclimate reconstructions and instrumental data at two
371 European sites. *Quaternary Science Reviews* 23: 137-143.

372 Charman DJ, Blundell A, Chiverrell RC, Hendon D, Langdon PG (2006) Compilation of
373 non-annually resolved Holocene proxy climate records: stacked Holocene
374 peatland palaeo-water table reconstructions from northern Britain. *Quaternary
375 Science Reviews* 25: 336–350.

376 Charman DJ, Blundell A, ACCROTELM members (2007) A new-European testate
377 amoebae transfer function for palaeohydrological reconstruction on
378 ombrotrophic peatlands. *Journal of Quaternary Science* 22: 209-221.

379 Field CA, Welsh AH (2007) Bootstrapping clustered data. *Journal of the Royal
380 Statistical Society B* 69: 369-390.

381 Hendon D, Charman DJ (2004) High-resolution peatland water-table changes for the
382 past 200 years: the influence of climate and implications for management. *The
383 Holocene* 14: 125-134.

384 Jassey VEJ, Chiapusio G, Mitchell EAD, Binet P, Toussaint M-L, Gilbert D (2011) Fine-
385 scale horizontal and vertical micro-distribution patterns of testate amoebae
386 along a narrow fen/bog gradient. *Microbial Ecology* 61: 374-385.

387 Juggins S (2003) *C2 user guide. Software for ecological and palaeoecological data
388 analysis and visualisation*. University of Newcastle, Newcastle Upon Tyne.

389 Kurek J, Cwynar LC (2009) The potential of site-specific and local chironomid-based
390 inference models for reconstructing past lake levels. *Journal of Paleolimnology
391* 42: 37-50.

392 Lamentowicz M, Mitchell EAD (2005) The ecology of testate amoebae (Protists) in
393 *Sphagnum* in north-west Poland in relation to peatland ecology. *Microbial
394 Ecology* 50: 48-63.

395 Lamentowicz Ł, Gabka M, Lamentowicz M (2007a) Species composition of testate
396 amoebae (Protists) and environmental parameters in a *Sphagnum* peatland.
397 *Polish Journal of Ecology* 55: 749-759.

398 Lamentowicz M, Tobolski K, Mitchell EAD (2007b) Palaeoecological evidence for
399 anthropogenic acidification of a kettle-hole peatland in northern Poland. *The
400 Holocene* 17: 1185-1196.

401 Lamentowicz Ł, Lamentowicz M, Gabka M (2008a) Testate amoebae ecology and a
402 local transfer function from a peatland in western Poland. *Wetlands* 28: 164-175.

403 Lamentowicz M, Obremaska M, Mitchell EAD (2008b) Autogenic succession, land-use
404 change, and climatic influences on the Holocene development of a kettle-hole
405 mire in Northern Poland. *Review of Palaeobotany and Palynology* 151: 21-40.

406 Lamentowicz M, Lamentowicz Ł, van der Knaap WO, Gąbka M, Mitchell EAD (2010)
407 Contrasting species—environment relationships in communities of testate
408 amoebae, bryophytes and vascular plants along the fen—bog gradient. *Microbial
409 Ecology* 59: 499-510.

410 Leorri E, Horton BP, Cearreta A (2008) Development of a foraminifera-based transfer
411 function in the Basque marshes, N. Spain: Implications for sea-level studies in the
412 Bay of Biscay. *Marine Geology* 251: 60-74.

413 Markel ER, Booth RK, Qin Y (2010) Testate amoebae and $\delta^{13}\text{C}$ of *Sphagnum* as
414 surface-moisture proxies in Alaskan peatlands. *The Holocene* 20: 463-475.

415 Massey A, Gehrels WR, Charman DJ, White SV (2006) An intertidal foraminifera-
416 based transfer function for reconstructing Holocene sea-level change in
417 southwest England. *Journal of Foraminiferal Research* 36: 215-232.

418 Mitchell EAD (2004) Response of testate amoebae (Protozoa) to N and P fertilization
419 in an Arctic wet sedge tundra. *Arctic Antarctic and Alpine Research* 36: 77–82.

420 Mitchell EAD, Buttler AJ, Warner BG, Gobat JM (1999) Ecology of testate amoebae
421 (Protozoa: Rhizopoda) in *Sphagnum* peatlands in the Jura mountains, Switzerland
422 and France. *Ecoscience* 6: 565-576.

423 Mitchell EAD, Charman DJ, Warner BG (2008) Testate amoebae analysis in ecological
424 and paleoecological studies of wetlands: past, present and future. *Biodiversity
425 and Conservation* 17: 2115-2137.

426 Meyer C, Bernard N, Moskura M, Toussaint ML, Denayer F, Gilbert D (2010) Effects of
427 urban particulate deposition on microbial communities living in bryophytes: An
428 experimental study. *Ecotoxicology and Environmental Safety* 73: 1776-1784.

429 Nguyen-Viet H, Gilbert D, Bernard N, Mitchell EAD, Badot PM (2004) Relationship
430 between atmospheric pollution characterized by NO_2 concentrations and testate
431 amoebae density and diversity. *Acta Protozoologica* 43: 233–239.

432 Nguyen-Viet H, Bernard N, Mitchell EAD, Cortet J, Badot PM, Gilbert D (2007)
433 Relationship between testate amoeba (Protist) communities and atmospheric
434 heavy metals accumulated in *Barbula indica* (Bryophyta) in Vietnam. *Microbial*
435 *Ecology* 53: 53–65.

436 Nguyen-Viet H, Bernard N, Mitchell EAD, Badot PM, Gilbert D (2008) Effect of lead
437 pollution on testate amoebae communities living in *Sphagnum fallax*: an
438 experimental study. *Ecotoxicology and Environmental Safety* 69: 130–138.

439 Payne RJ (2009) The standard preparation method for testate amoebae leads to
440 selective loss of the smallest shells. *Quaternary Newsletter* 119: 16-20.

441 Payne RJ (2010) Testate amoeba response to acid deposition in a Scottish peatland.
442 *Aquatic Ecology* 44, 373-385.

443 Payne RJ (2011) Can testate amoeba-based palaeohydrology be extended to fens?
444 *Journal of Quaternary Science* 26: 15-27.

445 Payne RJ, Mitchell EAD (2007) Ecology of testate amoebae from mires in the Central
446 Rhodope Mountains, Greece and development of a transfer function for
447 paleohydrological reconstruction. *Protist* 158: 159-171.

448 Payne RJ, Mitchell EAD (2009) How many is enough? Determining optimal count
449 totals for ecological and palaeoecological studies of testate amoebae. *Journal of*
450 *Paleolimnology* 42: 483-495.

451 Payne RJ, Kishaba K, Blackford JJ, Mitchell EAD (2006) The ecology of testate
452 amoebae in southcentral Alaskan peatlands: Building transfer function models for
453 palaeoenvironmental inference. *The Holocene* 16: 403-414.

454 Payne RJ, Charman DJ, Matthews S, Eastwood W (2008) Testate amoebae as
455 palaeoclimate proxies in Sürmene Ağaçbaşı Yaylasi peatland (Northeast Turkey).
456 *Wetlands* 28: 311-323.

457 Payne RJ, Charman, DJ, Gauci V (2010): The impact of simulated sulfate deposition
458 on peatland testate amoebae. *Microbial Ecology* 59: 76-83.

459 Payne RJ, Lamentowicz M, Mitchell EAD (2011) The perils of taxonomic
460 inconsistency in quantitative palaeoecology: simulations with testate amoeba
461 data. *Boreas* 40: 15-27.

462 R Development Core Team (2010) *R: A language and environment for statistical*
463 *computing*. R Foundation for Statistical Computing, Vienna, Austria.

464 Saunders KM, Hodgson DA, Harrison J, McMinn A (2008) Palaeoecological tools for
465 improving the management of coastal ecosystems: a case study from Lake King
466 (Gippsland Lakes) Australia. *Journal of Paleolimnology* 40: 33-47.

467 Sullivan, M.E. & Booth, R.K. (2011) The potential influence of short-term environmental
468 variability on the composition of testate amoeba communities in *Sphagnum* peatlands.
469 *Microbial Ecology* 62: 80-93.

470 Swindles GT, Charman DJ, Roe HM, Sansum PA (2009) Environmental controls on
471 peatland testate amoebae (Protozoa: Rhizopoda) in the North of Ireland:
472 Implications for Holocene palaeoclimate studies. *Journal of Paleolimnology* 42:
473 123–140.

474 Telford RJ, Birks HJB (2005) The secret assumptions of transfer functions: problems
475 with spatial autocorrelation in evaluating model performance. *Quaternary*
476 *Science Reviews* 24: 2173-2179.

477 Telford RJ, Birks HJB (2009) Evaluation of transfer functions in spatially structured
478 environments. *Quaternary Science Reviews* 28: 1309-1316.

479 Telford, R.J. Birks, H.J.B. (2011) Effect of uneven sampling along an environmental
480 gradient on transfer-function performance. *Journal of Paleolimnology* 46: 99-106.

481 Walsh JE (1947) Concerning the effect of intraclass correlation on certain significance
482 tests. *Annals of Mathematical Statistics* 18: 88-96.

483 Warner BG, Charman DJ (1994) Holocene changes on a peatland interpreted from
484 testate amoebae (Protozoa) analysis. *Boreas* 23: 270-280.

485 Woodland W, Charman DJ, Simms P (1998) Quantitative estimates of water tables
486 and soil moisture in Holocene peatlands from testate amoebae. *The Holocene* 8:
487 261-273.

488 Zong Y, Horton BP (1999) Diatom-based tidal-level transfer functions as an aid in
489 reconstructing Quaternary history of sea-level movements in the UK. *Journal of*
490 *Quaternary Science* 14: 153-167.

491

492

493 Appendix 1. Details of taxonomic harmonisation showing groupings and
 494 nomenclatural changes made to the original data. In addition to these changes small
 495 taxa (*Corythion* spp., *Trinema* spp., *Euglypha rotunda* type, *Euglypha cristata*,
 496 *Cryptodiffugia oviformis*, *Diffugia pulex* type and *Pseudodiffugia fulva* type) were
 497 eliminated where there was a difference in preparation method between training
 498 and test sets.

499

Dataset	Taxa in original data	Taxa here
Moss of Achnacree (Payne 2010a)	<i>Centropyxis aerophila</i> type <i>Phryganella acropodia</i> type <i>Corythion dubium</i> , <i>Trinema complanatum</i>	<i>Centropyxis cassis</i> type <i>Cyclopyxis arcelloides</i> type <i>Corythion-Trinema</i> type
Moidach More (Payne et al. 2010b)	<i>Phryganella acropodia</i> type <i>Corythion dubium</i> , <i>Trinema complanatum</i>	<i>Cyclopyxis arcelloides</i> type <i>Corythion-Trinema</i> type
UK (Woodland et al. 1998; Charman et al. 2007; Potts & Blackford unpublished data)	<i>Nebela minor</i> , <i>Nebela tincta</i> , <i>Nebela parvula</i>	<i>Nebela tincta</i> type
Alaska (Payne et al. 2006; Markel et al. 2010)	<i>Arcella arenaria</i> type, <i>A. catinus</i> type <i>Centropyxis aerophila</i> s.l., <i>C. cassis</i> type <i>Centropyxis laevis</i> , <i>C. ecornis</i> , <i>C. ecornis</i> type <i>Cyclopyxis arcelloides</i> type, <i>Phryganella acropodia</i> type, <i>P. acropodia</i> s.l. <i>Nebela dentistoma</i> , <i>N. vitraea</i> <i>Euglypha ciliata</i> , <i>E. compressa</i> , <i>E. strigosa</i> , <i>E. rotunda</i> s.l., <i>E. tuberculata</i> type, <i>E. strigosa</i> type, <i>E. rotunda</i> type <i>Nebela tincta</i> s.l., <i>N. tincta</i> , <i>N. parvula</i> <i>Placocista spinosa</i> s.l., <i>P. lens</i> , <i>P. spinosa</i> <i>Trigonopyxis arcula</i> , <i>T. minuta</i> <i>Trinema</i> spp., <i>T. lineare</i>	<i>Arcella catinus</i> type <i>Centropyxis aerophila</i> type <i>Centropyxis ecornis</i> type <i>Cyclopyxis arcelloides</i> type <i>Argynnia dentistoma</i> type <i>Euglypha</i> spp. <i>Nebela tincta</i> type <i>Placocista spinosa</i> type <i>Trigonopyxis arcula</i> type <i>Trinema</i> spp.

500