

1 Lake surface-sediment pollen dataset for the alpine meadow vegetation type from 2 the eastern Tibetan Plateau and its potential in past climate reconstructions

3 Xianyong Cao^{1,2*}, Fang ~~Tian~~³Tian², Kai ~~Li~~⁴Li³, Jian ~~Ni~~⁴Ni³, Xiaoshan Yu¹, Lina Liu¹, Nannan
4 Wang¹

5 ¹ Alpine Paleoecology and Human Adaptation Group (ALPHA), State Key Laboratory of ~~Alpine Ecology, Tibetan~~
6 Plateau Earth System, Resources and Environment (TPESRE), Institute of Tibetan Plateau Research, Chinese
7 Academy of Sciences, Beijing 100101, China

8 ~~² CAS Center for Excellence in Tibetan Plateau Earth Sciences, Institute of Tibetan Plateau Research, Chinese~~
9 ~~Academy of Sciences (CAS), Beijing 100101, China~~

10 ~~³ Beijing Key Laboratory of Resource Environment and GIS,~~² College of Resource Environment and Tourism,
11 Capital Normal University, Beijing, 100048, China

12 ⁴ College of Chemistry and Life Sciences, Zhejiang Normal University, Jinhua, 321004, China

13 Correspondence: Xianyong Cao (xcao@itpcas.ac.cn)

14

15

16 Abstract

17 A modern pollen dataset with an even distribution of sites is essential for pollen-based
18 past vegetation and climate estimations. As there were geographical gaps in previous
19 datasets covering the central and eastern Tibetan Plateau, lake surface-sediment
20 samples (n=117) were collected from the alpine meadow region on the Tibetan Plateau
21 between elevations of 3720 and 5170 m a.s.l. Pollen identification and counting were
22 based on standard approaches, and modern climate data were interpolated from a robust
23 modern meteorological dataset. A series of numerical analyses revealed that
24 precipitation is the main climatic determinant of pollen spatial distribution; Cyperaceae,
25 Ranunculaceae, Rosaceae, and *Salix* indicate wet climatic conditions, while Poaceae,

26 *Artemisia*, and Chenopodiaceae represent drought. Model performance of both
27 weighted-averaging partial least squares (WA-PLS) and the random forest (RF)
28 algorithm suggest that this modern pollen dataset has good predictive power in
29 estimating the past precipitation ~~for~~from pollen spectra from the eastern Tibetan Plateau.
30 In addition, a comprehensive modern pollen dataset can be established by combining
31 our modern pollen dataset with previous datasets, which will be essential for the
32 reconstruction of vegetation and climatic signals for fossil pollen ~~spectra~~spectra on the
33 Tibetan Plateau. Pollen datasets including both pollen counts and percentages for each
34 sample together with their site location and climatic data are available at the National
35 Tibetan Plateau Data Center (TPDC; DOI: 10.11888/Paleoenv.tpdc.271191).

36

37 **1 Introduction**

38 The relationship between modern pollen and climate, and its representation of
39 vegetation, is the basis for explaining and reconstructing past climate and vegetation
40 qualitatively or quantitatively (Juggins and Birks, 2012), so improving the quality of
41 the modern pollen dataset is a primary step for an objective investigation of the modern
42 relationship and to ensure reliable climate and vegetation reconstructions (Cao et al.,
43 2018). To make the pollen-source area and taphonomy as compatible as possible,
44 modern pollen assemblages should be retrieved from the same type of sedimentary
45 environment as the fossil pollen spectra (Birks et al., 2010). Hence, to reconstruct past
46 climate and vegetation from fossil pollen extracted from a lacustrine sediment, a
47 corresponding modern pollen dataset of samples collected from lake surface-sediments
48 is necessary. Although there are some modern pollen datasets for the Tibetan Plateau,
49 established to investigate the relationships between pollen and climate or vegetation
50 (Shen et al., 2006; Herzschuh et al., 2010; Ma et al., 2017), there are geographical gaps
51 (e.g. the central and eastern Tibetan Plateau) in the sampled lakes which may bias
52 interpretations.

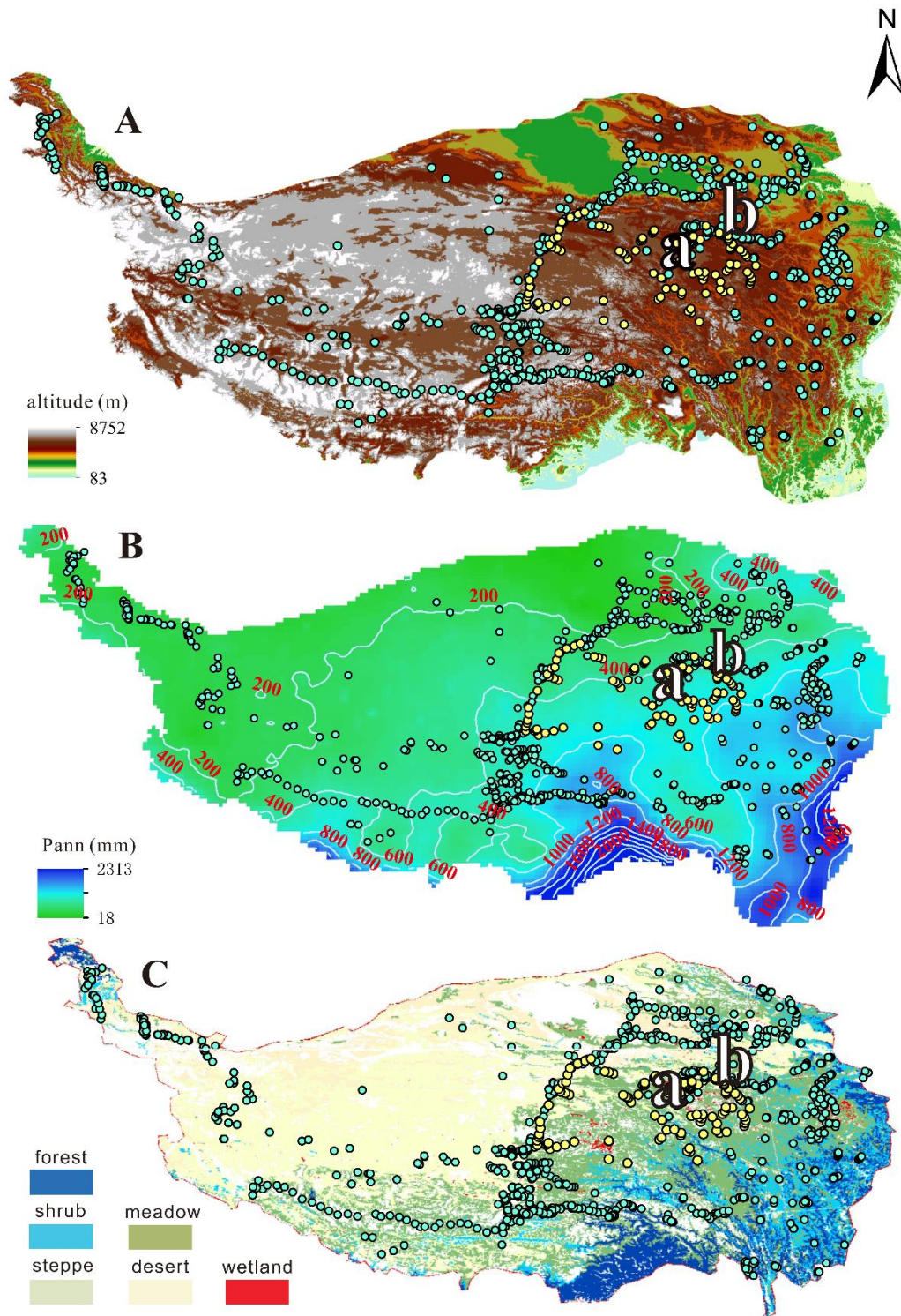
53 The available modern pollen datasets reveal that pollen assemblages on the Tibetan
54 Plateau are generally simple with Cyperaceae, *Artemisia*, Poaceae, and
55 Chenopodiaceae as the dominant taxa (e.g. Herzsuh et al., 2010; Cao et al., 2014),
56 with arboreal pollen taxa becoming more influential in the marginal areas (e.g. Ma et
57 al., 2017; Li et al., 2020). It is essential to identify the climatic indicators of the modern
58 pollen taxa (particular for the four dominant taxa) on the Tibetan Plateau, because the
59 climatic indicators derived from modern pollen datasets from the surrounding lowland
60 cannot be directly employed on the Tibetan Plateau. With our current modern pollen
61 dataset extracted from lake surface-sediments we aim to 1) fill a geographical gap and
62 thus establish a comprehensive modern pollen dataset covering the entire Tibetan
63 Plateau; 2) determine the climatic indicators for common pollen taxa from the alpine
64 meadow ecosystem; and 3) evaluate the predictive power of the modern dataset to
65 reconstruct past climate and assess the reliability of the random forest algorithm in
66 calibrating the pollen-climate relationship.

67

68 **2 Study area**

69 The elevation range of the lakes sampled for our pollen dataset is between 3720 and
70 5170 m a.s.l. with a median of 4420 m a.s.l. (the 25% quantile is 4230 m a.s.l and the
71 75% quantile is 4550 m a.s.l.; Figure 1). Climate of this region is controlled by the
72 Asian Summer Monsoon in summer with warm and wet climatic conditions, and by
73 westerlies in winter with cold and dry conditions (Wang, 2006). The eastern and central
74 Tibetan Plateau containing these sampled lakes (with >4000 m a.s.l elevation) is
75 covered by alpine meadow with sporadic patches of subalpine shrub. The plant
76 communities of the alpine meadow are dominated by *Kobresia* species (Cyperaceae)
77 generally, with Ranunculaceae, Asteraceae, *Polygonum* (Polygonaceae), *Potentilla*
78 (Rosaceae), Fabaceae, and Caryophyllaceae as the common taxa. The subalpine shrub
79 is generally distributed on the northern slopes of mountains with *Salix oritrepha* and
80 *Potentilla fruticosa* as the main shrub components, while the herbaceous taxa

81 mentioned above are also common (Wu, 1995; Herzschuh et al., 2010; unpublished
 82 vegetation survey).



83

84 **Figure 1** Spatial distribution of modern pollen samples (yellow dots: the 117 sampled

85 lakes; bluish green dots: previous samples (surface-soils and lake surface-sediments)

86 included in the dataset of Cao et al., 2014). A: Digital Elevation Model; B: isohyet map
 87 (mm); C: vegetation map. “a” and “b” indicate the locations of Koucha Lake and
 88 Xingxinghai Lake.

89 3 Materials and methods

90 3.1 Sample collecting and pollen processing

91 To ensure the even distribution of the representative lakes, we travelled not only along
 92 the hardened roads but also the dirt roads to collect samples from the alpine meadow
 93 on the eastern and central Tibetan Plateau, in July and August 2018. ~~Generally To~~
 94 ~~reduce the influence of long-distance pollen grains transported by wind and rivers,~~
 95 small and shallow ~~unnamed~~-lakes (or pools) with less than 100-m radius and without
 96 long inflow rivers (n=117) (locally sourced pollen grains are the dominant components
 97 for small lakes; Sugita, 1993) were selected to ~~reduce the influence of long distance~~
 98 ~~pollen transported by wind or rivers~~ collect pollen samples (Figure 1). To reduce the
 99 influence of the ~~local lake-shore~~ lake-shore vegetation component ~~from the lake shore~~, the lake
 100 surface-sediment samples were collected from the central part of each lake, with the
 101 top 2 cm of lake sediment forming the sample. ~~(Tian et al., 2008).~~ Although the selected
 102 lakes generally have an even distribution, there is still a gap in the south-west part of
 103 study area because of a lack of lake and road access (Figure 1).

104 For pollen extraction, approximately 10 g (wet untreated sediment) per sample were
 105 sub-sampled. Pollen samples were processed using standard acid-alkali-acid
 106 procedures (including 10% HCl, 10% KOH, 40% HF and 9:1 mixture of acetic
 107 anhydride and sulphuric acid successively; Fægri and Iversen, 1975) followed by 7-
 108 µm-mesh sieving. A tablet with *Lycopodium* spores (27560 grains/tablet) was added to
 109 each sample prior to pollen extraction as tracers (Maher, 1981). Pollen grains were
 110 identified with the aid of modern pollen reference slides collected from the eastern and
 111 central Tibetan Plateau (including 401 common species of alpine meadow; Cao et al.,
 112 2020) and published atlases for pollen and spores (Wang et al., 1995; Tang et al., 2017).

113 More than 500 terrestrial pollen grains were counted for each sample, and more than
114 200 *Lycopodium* spores were counted for most of the samples (mean=270 grains;
115 median=480 grains), both of which ensure a reliable representation of the entire pollen
116 assemblage by the counted pollen data.

117 3.2 Data processing

118 To obtain modern climatic data for the sampled lakes, the Chinese Meteorological
119 Forcing Dataset (CMFD; gridded near-surface meteorological dataset) with a temporal
120 resolution of three hours and a spatial resolution of 0.1° was employed (He et al., 2020).
121 The CMFD is made through the fusion of remote-sensing products, reanalysis datasets,
122 and *in situ* station data between January 1979 and December 2018, and its high
123 reliability has already ~~been~~ confirmed for western China including the Tibetan
124 Plateau (He et al., 2020). Geographical distances of each sampled lake to each pixel in
125 the CMFD were calculated based on their longitude/latitude coordinates using the
126 *rdist.earth* function in the *fields* package version 9.6.1 (Nychka, et al., 2019) for R
127 (version 3.6.0; R Core Team, 2019), and the ~~climatic data~~ meteorological data (three-
128 hour resolution between January 1979 and December 2018) of the nearest pixel to a
129 sampled lake were assigned to represent the climatic conditions of that lake. Finally,
130 the mean annual precipitation (P_{ann} ; mm), mean annual temperature (T_{ann} ; °C), and
131 mean temperature of the coldest month (Mt_{co} ; °C) and warmest month (Mt_{wa} ; °C) were
132 calculated for each sampled lake based on the long-term continuous meteorological
133 data.

134 To visualize the relationships between modern pollen assemblages and climatic
135 variables, ordination techniques were employed based on the square-root transformed
136 pollen data of 19 taxa (those present in at least 3 samples and with a $\geq 3\%$ maximum)
137 to stabilize variances and optimize the signal-to-noise ratio (Prentice, 1980). Detrended
138 correspondence analysis (DCA; Hill and Gauch, 1980) revealed that the length of the
139 first axis of the pollen data was 1.44 SD (standard deviation units), indicating a linear
140 response model is suitable for our pollen dataset (ter Braak and Verdonschot, 1995).

141 We performed redundancy analysis (RDA) to visualize the distribution of pollen
142 species and sampling sites along the climatic gradients, selecting the minimal adequate
143 model using forward selection and checking the variance inflation factors (VIF) at each
144 step. If VIF values were higher than 20, which ~~indicate~~indicates that some variables in
145 the model are co-linear, we stopped adding variables (ter Braak and Prentice, 1988).
146 These ordinations were performed using the *decorana* and *rda* functions in the *vegan*
147 package version 2.5-4 (Oksanen et al., 2019) for R.

148 Boosted regression tree (BRT) analysis was applied to determine how strongly the
149 climatic variables influence the distribution of each individual pollen taxon, using
150 square-root transformed pollen percentages. A BRT model was generated using the
151 *gbm.step* function in the *dismo* package 1.0-12 version (Hijmans et al., 2015) for R with
152 a Gaussian error distribution.

153 The basic assumption of pollen-based past climate reconstruction assumes that pollen
154 taxa recorded in the modern calibration-set have similar ecological requirements as
155 those in the fossil spectra (Juggins and Birks, 2012); in other words, the modern
156 vegetation-climate relationship is assumed to be stable temporally through the target
157 period for reconstruction. To evaluate the potential of the pollen dataset for past climate
158 reconstruction, both the traditional method of weighted-averaging partial least squares
159 (WA-PLS) and a new approach using the random forest (RF) algorithm were run. WA-
160 PLS was performed using the *WAPLS* function in the *rioja* package version 0.7-3
161 (Juggins, 2012) for R using leave-one-out cross-validation, pollen percentages of the
162 19 selected pollen taxa were square-root transformed, and the number of WA-PLS
163 components used was selected using a randomization *t*-test (Juggins and Birks, 2012).
164 We performed the RF algorithm with the *randomForest* package (version 4.6-14; Liaw,
165 2018) in R. RF is an algorithm that integrates multiple decision trees, and the
166 importance of each explanatory variable is measured as the percentage increase in the
167 residual sum of squares after randomly shuffling the order of the variables to determine
168 which explanatory variable can be added to the model. In our study, the importance of
169 all pollen taxa on the spatial distribution of P_{ann} was estimated and the model

170 systematically optimized by a stepwise reduction in variables by deleting the least
171 important one. Our final RF model includes 19 pollen taxa (Appendix [2B](#)), which all
172 make a positive contribution to the precipitation distribution. To assess the predictive
173 power of our pollen dataset, pollen spectra from Koucha Lake (covering the last 16 cal
174 ka BP; [\(calibrated thousand years before 1950 CE\)](#); 34.0°N; 97.2°E, 4540 m a.s.l.;
175 ~~Herzschuh et al., 2009; cal ka BP; calibrated thousand year before 1950 AD~~) and
176 Xingxinghai Lake (covering the last 7.5 cal ka BP; 34.8°N, 98.1°E, 4228 m a.s.l.; Zhang
177 et al., unpublished) were selected as the target fossil pollen datasets for quantitative
178 reconstruction. A statistical significance test for all reconstructions was performed
179 following the methods described in Telford and Birks (2011) using the *randomTF*
180 function in the *palaeoSig* package version 1.1.2 for both WA-PLS and RF
181 reconstruction methods separately (Telford, 2013).

182 [3.34](#) Data description

183 Pollen assemblages of the dataset from alpine ~~meadow~~[meadows](#) are dominated by
184 Cyperaceae (mean 68.4%, maximum 95.9%), with other herbaceous pollen taxa
185 common including Poaceae (mean 10.3%, maximum 87.7%), Ranunculaceae (mean
186 4.8%, maximum 33.6%), *Artemisia* (mean 3.7%, maximum 24.5%), and Asteraceae
187 (mean 2.1%, maximum 33.6%). *Salix* (mean 0.4%, maximum 5.3%) is the major shrub
188 taxon in these pollen assemblages, while arboreal taxa occur with low percentages
189 generally (mean total arboreal percentage 0.9%, maximum 5.8%), mainly comprising
190 *Pinus* (mean 0.3%, maximum 1.8%), *Betula* (mean 0.1%, maximum 0.9%), and *Alnus*
191 (mean 0.1%, maximum 0.7%). [These Published vegetation data \(e.g. Wu, 1995;](#)
192 [Herzschuh et al., 2010\) and our vegetation survey reveal that trees are absent from the](#)
193 [alpine meadow communities within the study area, thus we believe the arboreal pollen](#)
194 [with low abundances in the dataset will have been transported by wind from adjacent](#)
195 [regions to the south and east. Generally, these](#) pollen assemblages represent well the
196 plant components in the alpine meadow communities, although they are influenced

197 slightly by long-distance pollen transported by wind ~~or rivers (such as the arboreal~~
 198 ~~pollen taxa;~~ (Figure 2).

199

200

201 **Table 1** Summary statistics for parameters in the pollen dataset. Min.: minimum; Med.:
 202 median; Max.: maximum. Units for ~~Longitude~~longitude and ~~Latitude~~latitude are ~~degree,~~
 203 ~~for Altitude~~degrees, ~~elevation~~ is ~~in m a.s.l., for~~above sea level, Mt_{co}, Mt_{wa} and T_{ann}
 204 are °C, ~~for~~P_{ann} is mm, ~~while for~~and pollen ~~data~~data are %.

Parameter	Min.	Med.	Max.	Mean	Pollen taxa	Min.	Med.	Max.	Mean
Longitude	91.80	97.20	99.79	96.42	<i>Ilex</i>	0.00	0.00	0.18	0.00
Latitude	31.59	34.02	35.52	33.74	<i>Nitraria</i>	0.00	0.00	0.51	0.01
Elevation	3717	4422	5168	4399	Rosaceae	0.00	0.76	12.74	1.15
Mt _{co}	-19.21	-15.61	-7.41	-15.09	Tamaricaceae	0.00	0.00	0.75	0.03
Mt _{wa}	3.71	6.90	11.41	7.15	Apiaceae	0.00	0.16	3.98	0.32
T _{ann}	-7.27	-3.72	2.27	-3.39	<i>Artemisia</i>	0.19	2.43	24.51	3.68
P _{ann}	226	491	689	471	Asteraceae	0.00	1.46	33.56	2.09
Pollen taxa	Min.	Med.	Max.	Mean	Brassicaceae	0.00	0.36	28.17	1.22
<i>Abies</i>	0.00	0.00	0.38	0.01	Caryophyllaceae	0.00	0.16	2.26	0.23
<i>Cedrus</i>	0.00	0.00	0.19	0.00	Cyperaceae	4.84	76.24	95.91	68.67
<i>Picea</i>	0.00	0.00	2.52	0.10	Balsaminaceae	0.00	0.00	0.14	0.00
<i>Pinus</i>	0.00	0.18	1.76	0.32	Urticaceae	0.00	0.00	3.87	0.08
<i>Alnus</i>	0.00	0.00	0.67	0.11	Gentianaceae	0.00	0.16	4.85	0.40
<i>Betula</i>	0.00	0.00	0.94	0.11	Lamiaceae	0.00	0.00	1.05	0.12
<i>Carpinus</i>	0.00	0.00	0.63	0.06	Liliaceae	0.00	0.00	0.50	0.04
<i>Castanea</i>	0.00	0.00	2.44	0.06	Plantaginaceae	0.00	0.00	0.88	0.03
<i>Corylus</i>	0.00	0.00	1.88	0.07	Onagraceae	0.00	0.00	0.34	0.00
<i>Juglans</i>	0.00	0.00	0.82	0.01	Papaveraceae	0.00	0.00	0.82	0.03
Oleaceae	0.00	0.00	0.16	0.00	Poaceae	0.39	4.90	87.74	10.28
<i>Quercus</i>	0.00	0.00	2.00	0.06	Polemoniaceae	0.00	0.00	15.21	0.34
<i>Salix</i>	0.00	0.18	5.35	0.45	<i>Polygonum</i>	0.00	0.49	20.50	1.47
<i>Ulmus</i>	0.00	0.00	0.16	0.00	<i>Rumex</i>	0.00	0.00	1.64	0.03
Chenopodiaceae	0.00	0.48	15.44	0.86	<i>Koenigia</i>	0.00	0.00	2.96	0.39
<i>Ephedra</i>	0.00	0.00	1.66	0.12	Primulaceae	0.00	0.00	0.56	0.03
Ericaceae	0.00	0.00	0.19	0.01	Ranunculaceae	0.00	3.47	33.62	4.88
Euphorbiaceae	0.00	0.00	0.19	0.00	Saxifragaceae	0.00	0.00	4.69	0.10
Fabaceae	0.00	0.16	3.07	0.28	Scrophulariaceae	0.00	0.00	0.71	0.01
<i>Hippophaë</i>	0.00	0.00	5.62	0.27	Solanaceae	0.00	0.00	0.69	0.01
Rhamnaceae	0.00	0.00	0.17	0.00	<i>Thalictrum</i>	0.00	0.98	12.05	1.45

205

206 The region covered by these modern pollen samples has a P_{ann} gradient from 226 to 689
 207 mm, and cold thermal conditions with low T_{ann} (-7.3 to 2.3 °C) and Mt_{co} (-19.2 to
 208 -7.4 °C). A series of RDAs reveals that, relative to Mt_{co} and Mt_{wa} , P_{ann} explains more
 209 pollen assemblage variation (10.8% as a sole predictor in RDA) in the dataset (Table
 210 2). A biplot of the RDA shows that the direction of the P_{ann} vector has a smaller angle
 211 with the positive direction of Axis 1 (captures 43.2% of total inertia in the dataset)
 212 than with the positive direction of Axis 2 (10.3%), indicating that the major
 213 component of Axis 1 should be moisture. ~~The RDA separates axis 1, which is~~
 214 highly correlated with P_{ann} , divides the pollen taxa into two groups generally;
 215 Cyperaceae, Ranunculaceae, Rosaceae, and *Salix* indicating wet climatic conditions;
 216 (located along the positive direction of P_{ann}), while Poaceae, *Artemisia*, and
 217 Chenopodiaceae represent drought (located along the negative direction of P_{ann} ; Figure
 218 3). ~~Since the~~ Axis 2 is highly correlated with the two temperature variables; however
 219 these dominant pollen taxa have insignificant distributions along the axis, hence
 220 temperature is the secondary climatic variable for the pollen dataset relative to
 221 precipitation (Figure 3). Because of low occurrences and abundances for some rare
 222 pollen taxa, BRT models are only performed ~~successfully~~ for only 14 dominant or
 223 common pollen taxa. BRT modelling results ~~also~~ suggest that P_{ann} is the main climatic
 224 determinant for 9 out of 10 of the major pollen taxa with >0.6 prevalence, ~~while~~ with
 225 Asteraceae ~~is an~~ exception ~~with~~ having Mt_{co} as its main climatic determinant (68%;
 226 Table 3). BRT results reveal that pollen abundances of Cyperaceae, Ranunculaceae,
 227 and *Salix* are positively related to P_{ann} , while those of Poaceae, *Artemisia*, and
 228 Chenopodiaceae have a negative relationship with P_{ann} , ~~which are~~ consistent with the
 229 RDA results (Figure 3 and 4; Appendix 1).

230

231 5 Potential use of the modern pollen dataset

232 ~~Numerical~~ Numerical analyses reveal that P_{ann} is the most important climatic determinant
 233 of pollen distribution in the eastern Tibetan Plateau, hence, P_{ann} is selected as the target

234 variable in the calibration-set to assess the predictive power of this pollen dataset. Both
 235 approaches (WA-PLS, RF) perform well with low RMSEP values (the root mean square
 236 error of prediction) and high r^2 values (coefficient of determination between observed
 237 and predicted climatic variables; Figure 5). However, the plots of observed vs. predicted
 238 P_{ann} show a overestimate of P_{ann} for arid sites and an underestimate for wet sites (Figure
 239 5). Hence, the inevitable “edge effects” should be treated with caution. Nevertheless,
 240 ~~the reconstruction with~~ reconstructions covering ca. 400–500 mm P_{ann} should be reliable
 241 because of the low bias in the central part of the P_{ann} gradient (Figure 5).

242 Although the model performance of RF is not any better than that of WA-PLS, the
 243 reconstruction produced by RF might be more reliable as suggested by the statistical
 244 significance testing and comparison with modern observed P_{ann} for the two lakes
 245 (Koucha Lake and Xingxinghai Lake). Statistical significance testing ~~reveals~~ shows that
 246 ~~reconstructions based on the proportion of variance in the fossil data explained by the~~
 247 WA-PLS ~~explain~~ reconstruction is less ~~proportion~~ than the 95% quantile of the
 248 ~~proportion of~~ variance explained by a reconstruction based on random environmental
 249 variables (999 ~~timestrials~~) for the two lakes, while reconstructions ~~produced~~ produced
 250 by RF explain a higher proportion ~~than the 95% quantile~~ (Figure 6). In other words,
 251 reconstructions produced by RF might be controlled by the major pollen components,
 252 because the explained proportion of variance in the fossil pollen spectra is closer to that
 253 explained by the first PCA axis, while reconstructions by WA-PLS could be influenced
 254 more by the pollen taxa with low abundances (Figure 6). The hypothesis that WA-PLS
 255 is ~~more~~ influenced more by low-abundance pollen taxa is supported by the high-
 256 variation in reconstructed P_{ann} among the fossil pollen samples (Figure 7). Relative to
 257 reconstructions of WA-PLS, results of RF have lower temporal variation and fewer
 258 outliers, and the predicted P_{ann} by RF is closer to the observed P_{ann} for the two lakes
 259 (Koucha Lake, 500 mm; Xingxinghai Lake, 350 mm) than that by WA-PLS.

260

261 **Table 2** Summary statistics of redundancy analysis (RDA) of 19 pollen species and
 262 four climatic variables. VIF: variance inflation factor; P_{ann} : mean annual precipitation

263 (mm); Mt_{co} : mean temperature of the coldest month ($^{\circ}C$); Mt_{wa} : mean temperature of
 264 the warmest month ($^{\circ}C$); T_{ann} : annual mean temperature ($^{\circ}C$).

Climatic variables	VIF (without T_{ann})	VIF (with T_{ann})	Climatic variables as sole predictor	Marginal contribution based on climatic variables	
			Explained variance (%)	Explained variance (%)	<i>p</i> -value
P_{ann}	1.6	2.9	10.8	14.7	0.001
Mt_{co}	4.8	161.4	2.6	4.8	0.001
Mt_{wa}	3.8	83.9	1.6	1.3	0.100
T_{ann}	-	447.8	-	-	-

265

266

267

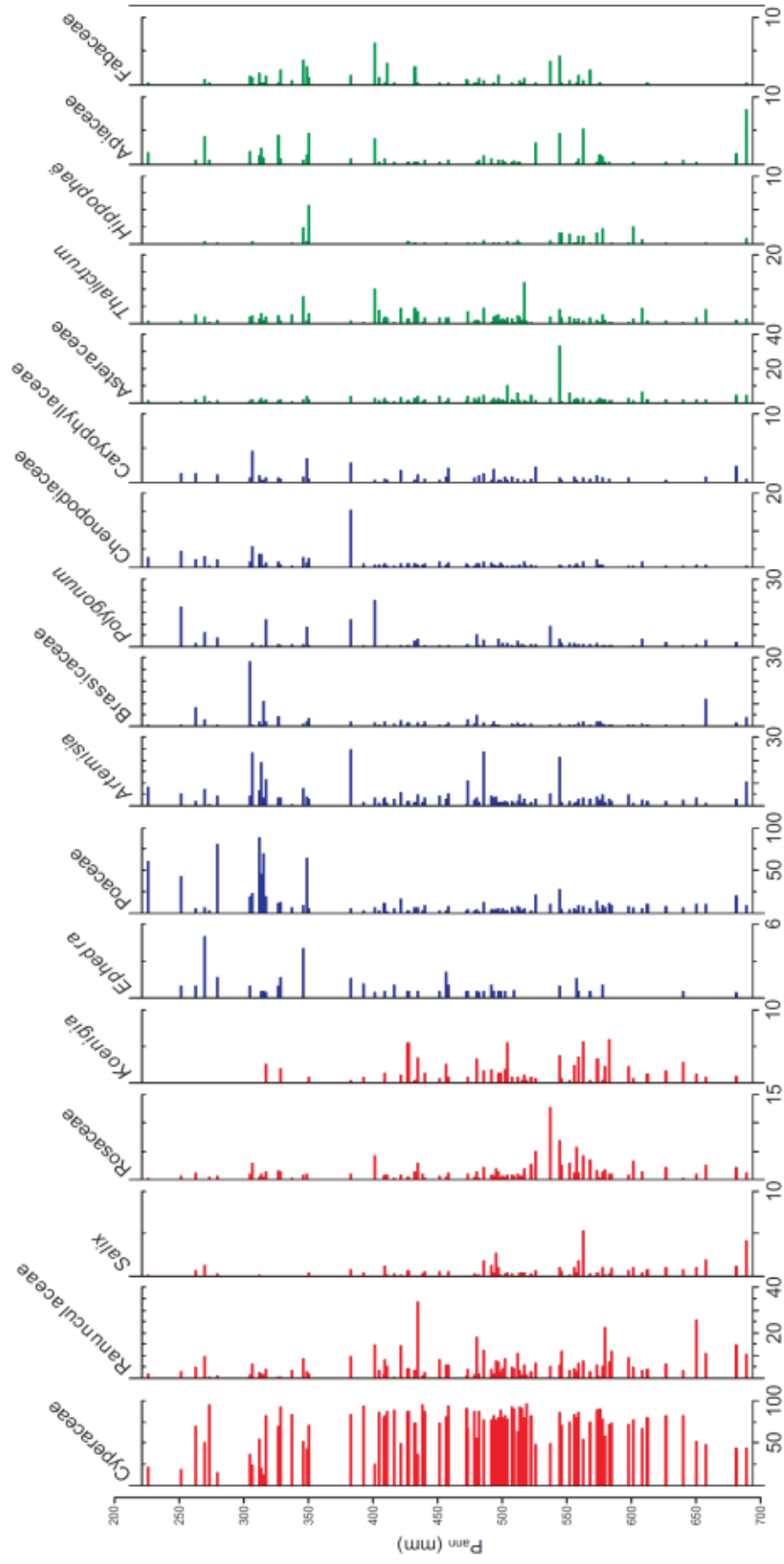
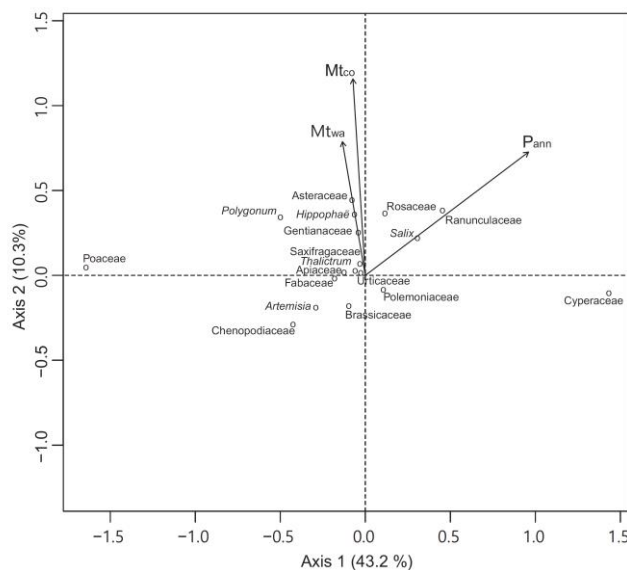


Figure 2 Pollen diagram showing the major taxa (percentage; %) of the 117 samples arranged by mean annual precipitation (P_{ann} ; mm). Pollen taxa with red bars are positively related to P_{ann} , those with blue bars are negatively related to P_{ann} , while the relationship is insignificant for those with green bars.

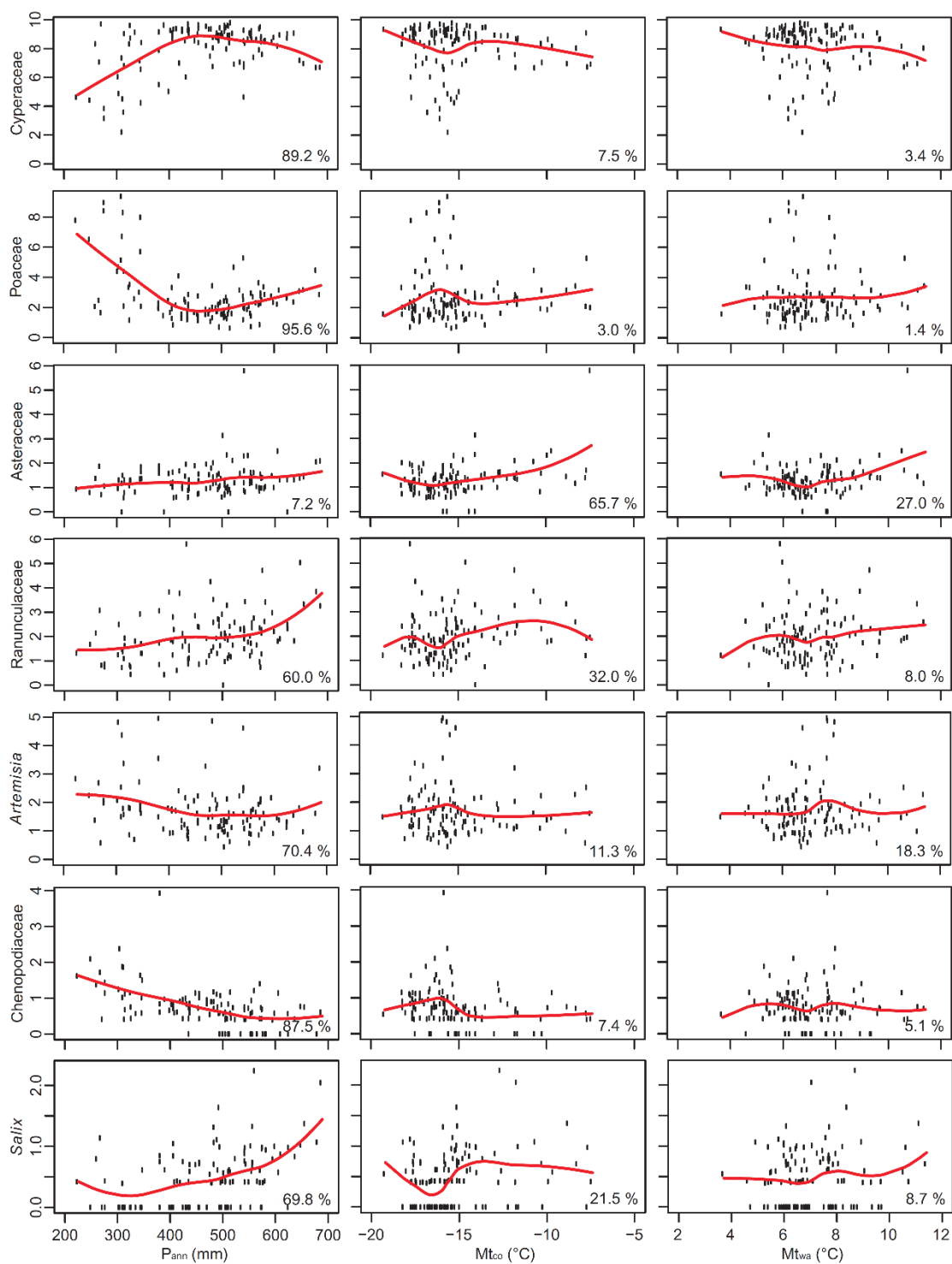


269

270 **Figure 3** Plot of the first two redundancy analysis (RDA) axes showing the
 271 relationships between 18 pollen taxa (circles) and 3 climatic variables (arrows). P_{ann}:
 272 mean annual precipitation (mm); Mt_{co}: mean temperature of the coldest month (°C);
 273 Mt_{wa}: mean temperature of the warmest month (°C).

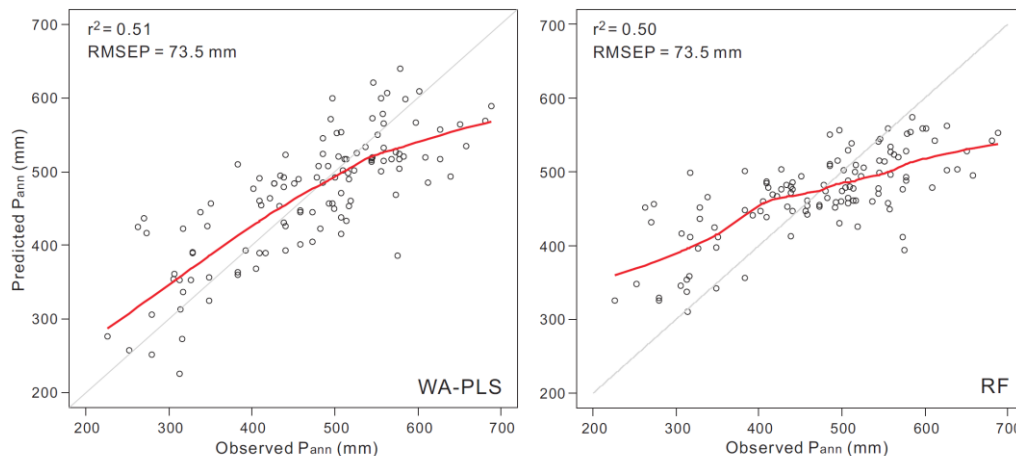
274 **Table 3** Relative influence of climatic variables to the spatial distributions of 14 pollen
 275 taxa based on boosted regression tree (BRT) models. For each variable, the relative
 276 influence is expressed as a percentage among the three variables. Pollen taxa are
 277 ordered by decreasing prevalence (the proportion of sites in which each taxon is
 278 present).

Taxa	Prevalence	P _{ann}	Mt _{co}	Mt _{wa}
Cyperaceae	1.00	89.3%	7.5%	3.2%
Poaceae	1.00	95.1%	3.3%	1.5%
<i>Artemisia</i>	1.00	69.3%	12.9%	17.8%
Ranunculaceae	0.99	56.9%	33.7%	9.4%
Asteraceae	0.97	7.2%	68.0%	24.8%
Rosaceae	0.90	32.2%	52.7%	15.1%
Chenopodiaceae	0.85	89.1%	5.8%	5.1%
Brassicaceae	0.81	49.6%	37.4%	13.0%
<i>Polygonum</i>	0.75	42.8%	31.9%	25.3%
<i>Salix</i>	0.63	71.2%	21.7%	7.1%
Fabaceae	0.54	79.3%	11.0%	9.6%
Gentianaceae	0.54	10.5%	63.1%	26.4%
Apiaceae	0.53	33.6%	30.5%	35.9%
<i>Hippophaë</i>	0.37	9.6%	77.6%	12.9%
Number of > 50% relative influence:		7	3	0



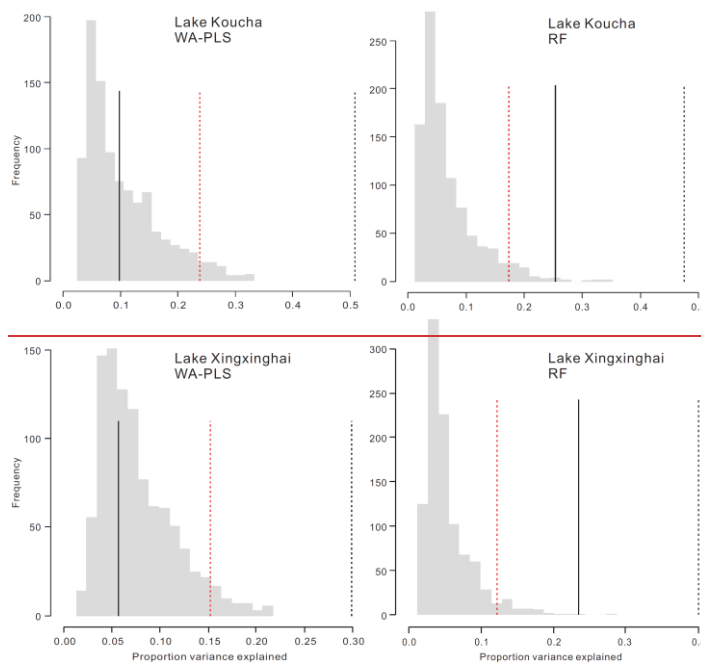
279

280 **Figure 4** Boosted regression tree (BRT) modelled climate influences on pollen (seven
 281 dominant or major taxa) percentages. The pollen responses to three climatic variables
 282 (red curves) are fitted with local polynomial regression (LOESS).

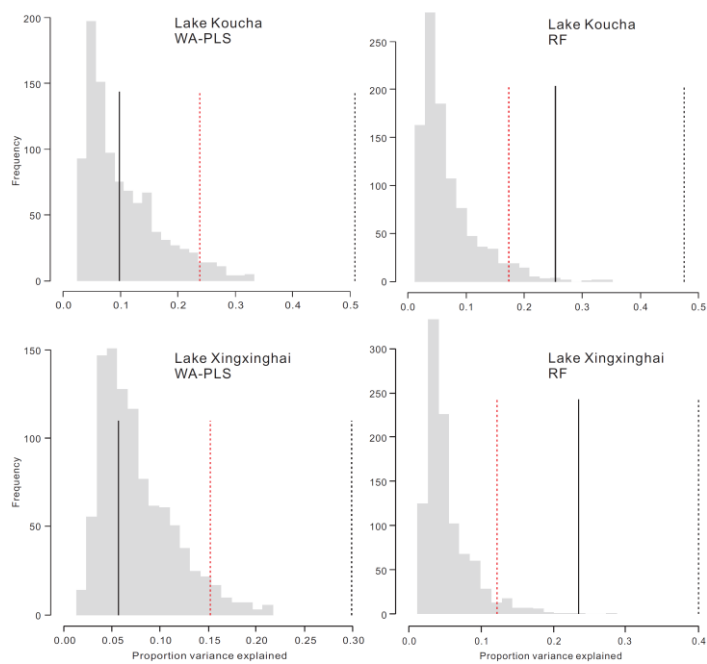


283

284 **Figure 5** Scatter plots of observed annual precipitation (P_{ann}) vs. predicted P_{ann} by
 285 weighted averaging partial least squares regression (WA-PLS) and random forest
 286 algorithm (RF).

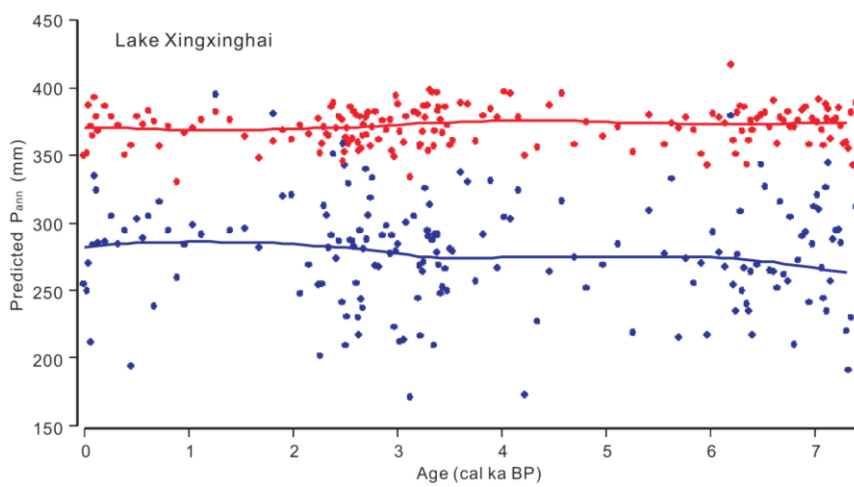
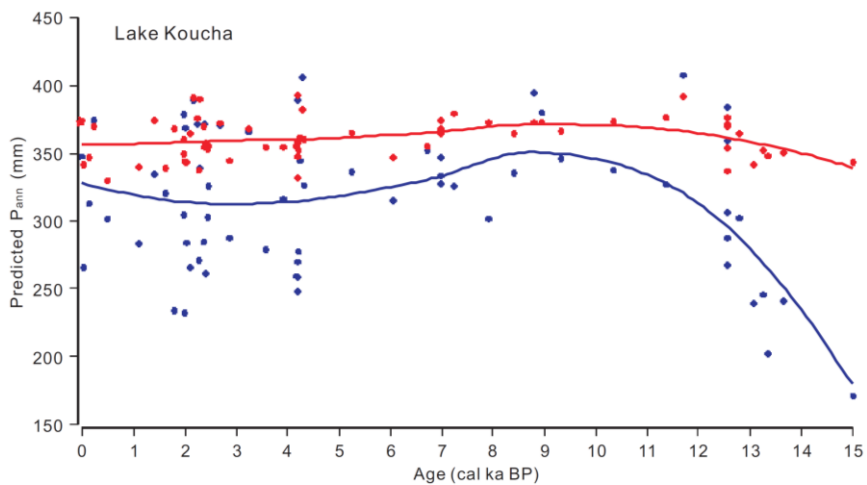


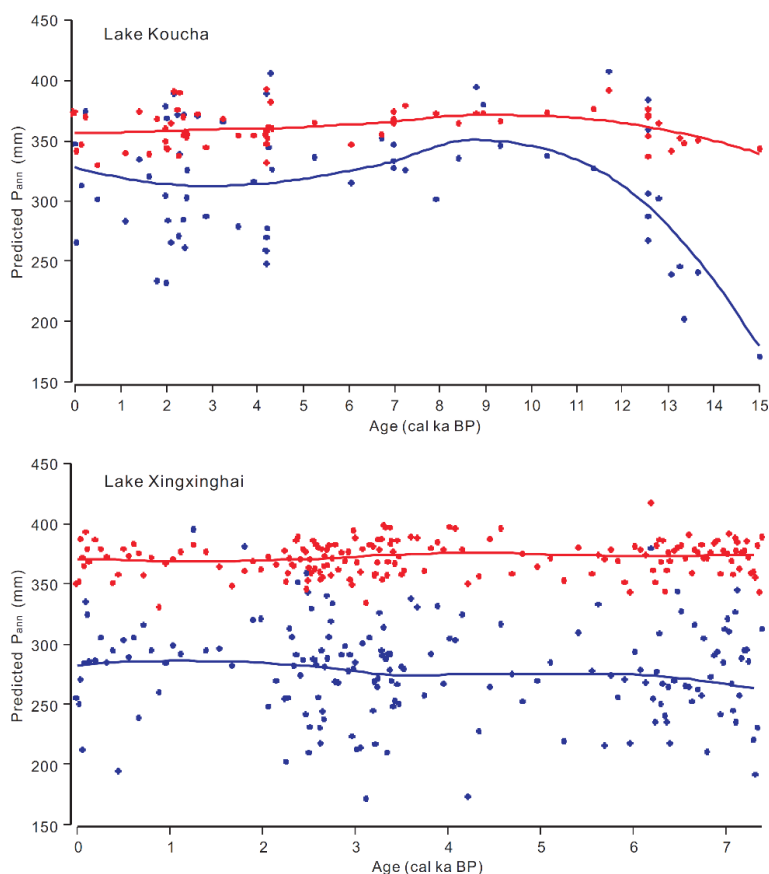
287



288

289 **Figure 6** Statistical significance test of P_{ann} ~~reconstruction~~reconstructions from two
 290 lakes using weighted-averaging partial least squares regression (WA-PLS) and the
 291 random forest (RF) algorithm. Grey histograms indicate the proportion of variance in
 292 the fossil pollen spectra explained by random variables (999 times) and the red dotted
 293 line is the 95% quantile, the black dotted line is the variance in the pollen explained by
 294 the first PCA axis, and the black solid line is the explanation by the reconstructed P_{ann} .





296

297 **Figure 7** Annual precipitation (P_{ann} ; mm) reconstructions for two Tibetan lakes using
 298 the weighted-averaging partial least squares regression (blue) and random forest
 299 algorithm (red). The curves are fitted by local polynomial regression (LOESS).

300 46 Data availability

301 Pollen datasets including both pollen counts and percentages for each sample together
 302 with their locations and climatic data are available at the National Tibetan Plateau Data
 303 Center (TPDC; DOI: [10.11888/Paleoenv.tpdc.271191](https://doi.org/10.11888/Paleoenv.tpdc.271191)).

304 7 Summary

305 We present a regional modern pollen dataset extracted from lake surface-sediments
 306 from the alpine meadow vegetation type on the Tibetan Plateau (eastern Tibetan Plateau,
 307 91.8° – 99.8° E and 31.6° – 35.5° N), including pollen counts and pollen percentages
 308 together with their positions and climatic data. Numerical analyses reveal that P_{ann} is

309 the most important climatic determinant for pollen distribution in the dataset, and our
310 dataset behaves reliably and has good predictive power for past moisture reconstruction,
311 and the random forest algorithm is a potentially ~~robust~~reliable approach in pollen-based
312 past environment reconstruction.

313 In addition, our open-access dataset can fill the ~~geographic~~geographical gap left by the
314 two previous modern pollen datasets (lake surface-sediments; Shen et al., 2006;
315 Herzschuh et al., 2010) on the eastern Tibetan Plateau. By combining our dataset here
316 with the previous ones (e.g. Herzschuh et al., 2019), a comprehensive modern pollen
317 dataset is created covering vegetation types from the alpine forest to alpine steppe on
318 the Tibetan Plateau, and will greatly improve the reliability of past vegetation
319 reconstructions and climate estimations.

320 ~~5 Data availability~~

321 ~~Pollen datasets including both pollen counts and percentages for each sample together~~
322 ~~with their locations and climatic data are available at the National Tibetan Plateau Data~~
323 ~~Center (TPDC; DOI: 10.11888/Paleoenv.tpdc.271191).~~

324 **Author contributions.** XC and JN designed the pollen dataset. XC and KL collected
325 pollen samples. XY and FT compiled the pollen identification and counting. XC and
326 FT performed numerical analyses and organized the manuscript, LL and NW prepared
327 the figures. All authors discussed the results and contributed to the final paper.

328 **Acknowledgements**

329 The sample collection and research were supported by the National Natural Science
330 Foundation of China (Grant No. 41877459 and 41930323), CAS Pioneer Hundred
331 Talents Program (Xianyong Cao) and Pan-Third Pole Environment Study for a Green
332 Silk Road of CAS Strategic Priority Research Program (XDA20090000).

333 **References**

334 Birks, H.J.B., Heiri, O., Seppä, H. and Bjune, A.E.: Strengths and weaknesses of
335 quantitative climate reconstructions based on late-Quaternary biological proxies,
336 *Open Ecol. J.* [3](#), 68–110, 2010.

337 Cao, X., Tian, F. and Ding, W.: Improving the quality of pollen-climate calibration-sets
338 is the primary step for ensuring reliable climate reconstructions, *Sci. Bull.* [63](#),
339 1317–1318, 2018.

340 Cao, X., Tian, F., Li, K. and Ni, J.: Atlas of pollen and spores for common plants from
341 the east Tibetan Plateau. National Tibetan Plateau Data Center, DOI:
342 10.11888/Paleoenv.tpdc.270735, 2020.

343 Cao, X.Y., Herzschuh, U., Telford, R.J. and Ni, J.: A modern pollen-climate dataset
344 from China and Mongolia: assessing its potential for climate reconstruction, *Rev.*
345 *Palaeobot. Palynol.* [211](#), 87–96, 2014.

346 Fægri, K. and Iversen, J.: Textbook of pollen analysis, Munksgaard, Copenhagen, 1975.

347 He, J., Yang, K., Tang, W., Lu, H., Qin, J., Chen, Y. and Li, X.: The first high-resolution
348 meteorological forcing dataset for land process studies over China, *Sci. Data*, [7](#),
349 25, DOI: 10.1038/s41597-020-0369-y, 2020.

350 Herzschuh, U., Birks, H.J.B., Mischke, S., Zhang, C. and Böhner, J.: A modern pollen-
351 climate calibration set based on lake sediments from the Tibetan Plateau and its
352 application to a Late Quaternary pollen record from the Qilian Mountains, *J.*
353 *Biogeogr.* [37](#), 752–766, 2010.

354 Herzschuh, U., Cao, X., Laepple, T., Dallmeyer, A., Telford, R., Ni, J., Chen, F., Kong,
355 Z., Liu, G., Liu, K.-B., Liu, X., Stebich, M., Tang, L., Tian, F., Wang, Y.,
356 Wischnewski, J., Xu, Q., Yan, S., Yang, Z., Yu, G., Zhang, Y., Zhao, Y. and Zheng,
357 Z.: Position and orientation of the westerly jet determined Holocene rainfall
358 patterns in China, *Nat. Commun.*, [10](#), 2376, 2019.

- 359 Herzschuh, U., Kramer, A., Mischke, S. and Zhang, C.: Quantitative climate and
360 vegetation trends since the late glacial on the northeastern Tibetan Plateau deduced
361 from Koucha Lake pollen spectra. *Quaternary Research, Quat. Res.*, 71, 162–171,
362 2009.
- 363 Hijmans, R.J., Phillips, S., Leathwick, J. and Elith, J.: Dismo: Species Distribution
364 Modeling, version 1.0-12, available at: [http://CRAN.R-project.org/package/
365 dismo](http://CRAN.R-project.org/package/dismo), 2015.
- 366 Hill, M.O. and Gauch, H.G.: Detrended correspondence analysis: an improved
367 ordination technique, *Vegetatio*, 42, 41–58, 1980.
- 368 Juggins, S. and Birks, H.J.B.: Quantitative environmental reconstructions from
369 biological data, in: Birks, H.J.B., Lotter, A.F., Juggins, S. and Smol, J.P. (eds.),
370 *Tracking environmental change using lake sediments*, (vol. 5): Data handling
371 and numerical techniques, Springer, Dordrecht, 431–494, 2012.
- 372 Juggins, S.: Rioja: analysis of Quaternary Science Data version 0.7-3, available at:
373 <http://cran.r-project.org/web/packages/rioja/index.html>, 2012.
- 374 Li, J.F., Xie, G., Yang, J., Ferguson, D.F., Liu, X.D., Liu, H. and Wang, Y.F.: Asian
375 Summer Monsoon changes the pollen flow on the Tibetan Plateau, *Earth-Sci.
376 Rev.*, 202, 103114, 2020.
- 377 Liaw, A.: randomForest: Breiman and Cutler's Random Forests for Classification and
378 Regression, available at: [https://cran.r-project.org/web/packages/randomForest/
379 index.html](https://cran.r-project.org/web/packages/randomForest/index.html), 2018.
- 380 Ma, Q., Zhu, L., Wang, J., Ju, J., Lü, X., Wang, Y., Guo, Y., Yang, R., Kasper, T.,
381 Haberzettl, T. and Tang, L.: *Artemisia/Chenopodiaceae* ratio from surface lake
382 sediments on the central and western Tibetan Plateau and its application,
383 *Palaeogeogr. Palaeoclim. Palaeoecol.*, 479, 138–145, 2017.

- 384 Maher, L.J.: Statistics for microfossil concentration measurements employing
385 [sanmplessamples](#) spiked with marker grains, *Rev. Palaeobot. Palynol.*, 32, 153–
386 191, 1981.
- 387 Nychka, D., Furrer, R., Paige, J. and Sain, S.: *fields: Tools for spatial data*, version 9.6.1,
388 available at: <https://cran.r-project.org/web/packages/fields/>, 2019.
- 389 Oksanen, J., Blanchet, F.G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D.,
390 Minchin, P.R., O'Hara, R.B., Simpson, G.L., Solymos, P., Stevens, M.H.H.,
391 Szoecs, E. and Wagner, H.: *vegan: Community Ecology Package*, version 2.5-4,
392 available at: <https://cran.r-project.org/web/packages/vegan/index.html>, 2019.
- 393 Prentice, I.C.: Multidimensional scaling as a research tool in Quaternary palynology: a
394 review of theory and methods, *Rev. Palaeobot. Palynol.*, 31, 71–104, 1980.
- 395 R Core Team: *R, A language and environment for statistical computing*, R Foundation
396 for Statistical Computing, Vienna, 2019.
- 397 Shen, C., Liu, K.B., Tang, L. and Overpeck, J.T.: Quantitative relationships between
398 pollen rain and climate in the Tibetan Plateau. *Rev. Palaeobot. Palynol.*, 140, 61–
399 77, 2006.
- 400 [Sugita, S.: A model of pollen source area for an entire lake surface. *Quat. Res.*, 39, 369–](#)
401 [244, 1993.](#)
- 402 Tang, L., Mao, L., Shu, J., Li, C., Shen, C. and Zhou, Z.: *Atlas of Quaternary pollen*
403 *and spores in China*, Science Press, Beijing, 2017.
- 404 ter Braak, C.J.F. and Prentice, I.C.: A theory of gradient analysis, *Adv. Ecol. Res.*, 18,
405 271–317, 1988.
- 406 ter Braak, C.J.F. and Verdonschot, P.F.M.: Canonical correspondence analysis and
407 related multivariate methods in aquatic ecology, *Aquat. Sci.*, 57, 255–289, 1995.

408 Tian, F., Xu, Q., Li, Y., Cao, X., Wang, X. and Zhang, L.: Pollen assemblage
409 characteristics of lakes in the monsoon fringe area of China. Chinese Sci. Bull.,
410 53(21), 3354–3363, 2008.

411 Wang, B.: The Asian Monsoon, Springer, Chichester, 2006.

412 Wang, F.X., Qian, N.F., Zhang, Y.L. and Yang, H.Q.: Pollen Flora of China, Science
413 Press, Beijing, 1995.

414 Wu, Z.Y.: The vegetation of China. Science Press, Beijing, 1995 (in Chinese).

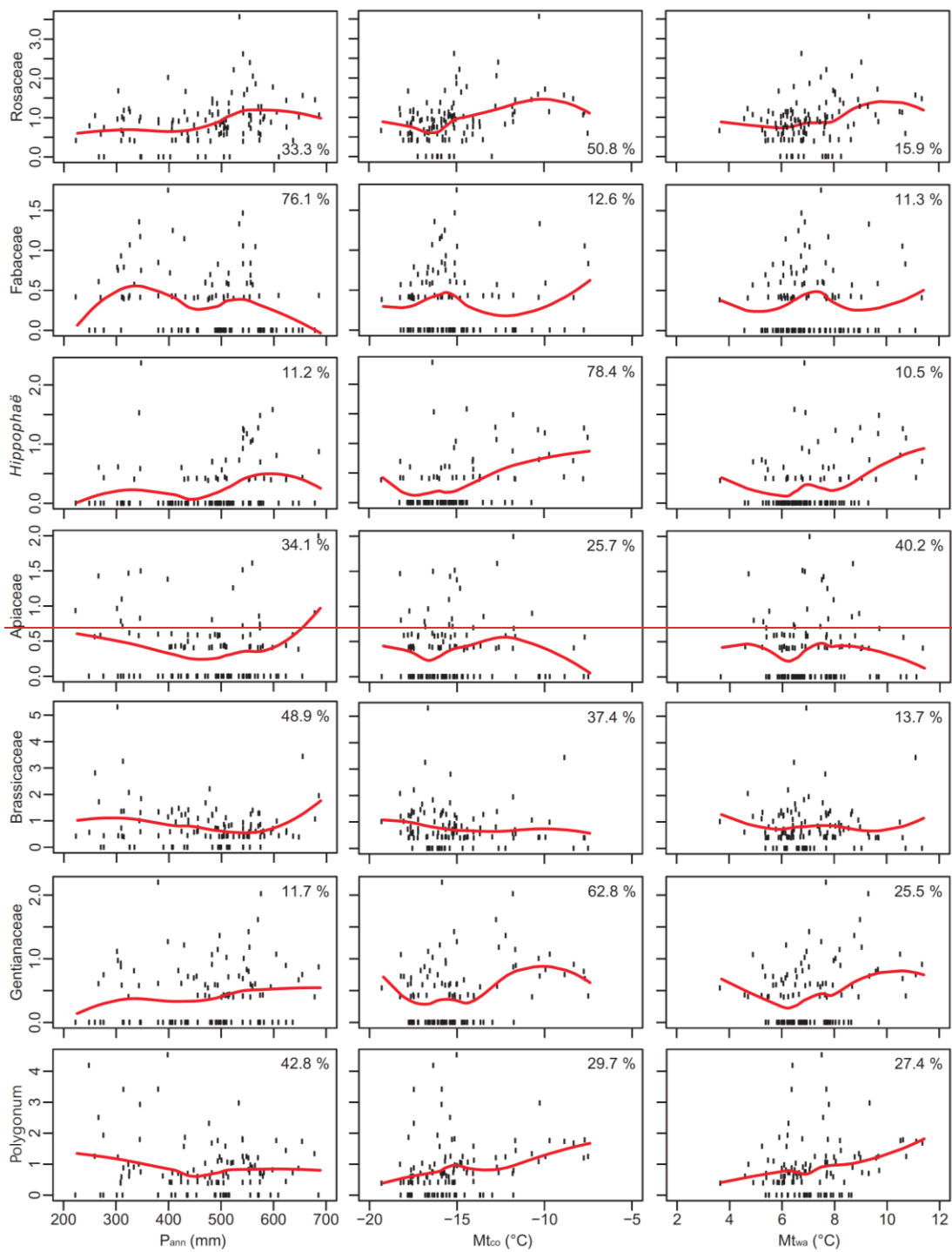
415

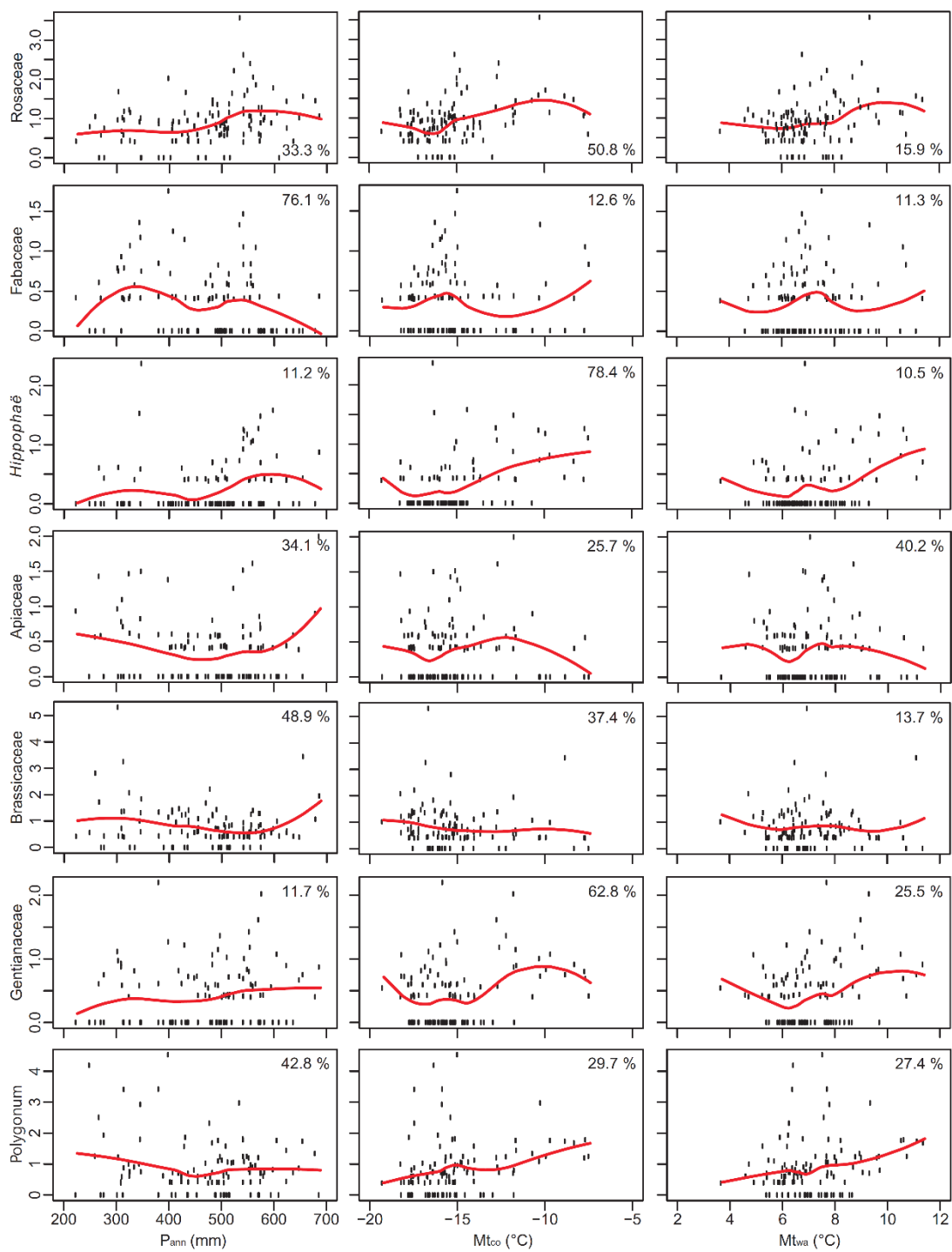
416

417

418 Appendix A

419 Boosted regression tree (BRT) modelled climate influences on pollen (seven common
420 or minor taxa) percentages. The pollen responses to three climatic ~~variables~~variables
421 (red curves) are fitted with a local polynomial regression (LOESS).





423

424

425 Appendix B

426 Importance (imp) of pollen taxa on the spatial distribution of P_{ann} ~~were~~ repeatedly
 427 assessed by the random forest algorithm (RF). Shown in bold are the pollen taxa
 428 selected for the P_{ann} reconstruction based on RF.

Taxa	imp-run1	imp-run2	imp-run3	imp-run4	imp-run5
<i>Abies</i>	-1.5723				
<i>Cedrus</i>	0.0000				
<i>Picea</i>	0.3104	3.4397	3.5811	2.1705	1.1599
<i>Pinus</i>	-1.6225				
<i>Alnus</i>	-0.3501				
<i>Betula</i>	5.8217	7.4399	7.4490	5.7763	5.9524
<i>Carpinus</i>	-1.2049				
<i>Castanea</i>	-1.4692				
<i>Corylus</i>	0.2806	-0.3715			
<i>Juglans</i>	0.0000				
Oleaceae	0.0000				
<i>Quercus</i>	-0.4776				
<i>Salix</i>	9.2463	9.6372	10.0018	9.4944	10.2897
<i>Ulmus</i>	-0.6041				
Chenopodiaceae	17.7282	18.0369	16.8653	16.3110	18.5089
<i>Ephedra</i>	2.8306	2.9972	4.4539	3.5096	4.0226
Ericaceae	0.0755	1.7893	-0.2415		
Euphorbiaceae	-0.9748				
Fabaceae	2.4847	2.5302	3.5031	3.2985	1.8323
<i>Hippophaë</i>	5.5569	3.5027	4.0142	3.1174	4.5627
Rhamnaceae	0.0000				
<i>Ilex</i>	0.0000				
<i>Nitraria</i>	-1.0010				
Rosaceae	3.0053	4.8099	2.9771	3.6032	4.3940
Tamaricaceae	-2.3780				
Apiaceae	-0.6466				
<i>Artemisia</i>	1.7355	-0.0902			
Asteraceae	2.3902	1.7955	1.1307	-1.0880	
Brassicaceae	1.7269	2.2776	1.4596	1.5560	1.5308
Caryophyllaceae	-0.0033				
Cyperaceae	9.9824	9.8975	11.1838	10.4553	10.3560
Balsaminaceae	0.0000				
Urticaceae	0.8534	-1.4774			
Gentianaceae	1.1305	-0.8603			
Lamiaceae	3.3097	2.6853	3.4047	2.2080	2.6588
Liliaceae	-0.5353				
Plantaginaceae	2.3294	1.3210	1.4498	0.8906	0.8763

Onagraceae	1.0010	-0.8613			
Papaveraceae	0.1148	1.0344	-1.7028		
Poaceae	13.8815	14.5295	14.7793	15.7914	16.2655
Polemoniaceae	-0.5507				
Polygonum	0.0523	2.4552	2.9776	1.9432	2.3618
<i>Rumex</i>	1.0010	0.0000			
Koenigia	5.4498	4.3961	3.3305	4.1574	4.9186
Primulaceae	-1.2283				
Ranunculaceae	6.4799	8.9763	7.6140	7.5498	5.5157
Saxifragaceae	0.9422	1.3283	1.8760	4.1134	2.3728
Scrophulariaceae	-1.0010				
Solanaceae	1.0010	-1.0008			
Thalictrum	2.9345	2.3850	2.6363	2.4267	3.3457
