

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*}, Fang Tian², Kai Li³, Jian Ni³, Xiaoshan Yu¹, Lina Liu¹, Nannan Wang¹

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

7 ² College of Resource Environment and Tourism, Capital Normal University, Beijing, 100048, China

8 ³ College of Chemistry and Life Sciences, Zhejiang Normal University, Jinhua, 321004, China

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

10

11

12 **Abstract**

13 A modern pollen dataset with an even distribution of sites is essential for pollen-based
14 past vegetation and climate estimations. As there were geographical gaps in previous
15 datasets covering the central and eastern Tibetan Plateau, lake surface-sediment
16 samples (n=117) were collected from the alpine meadow region on the Tibetan Plateau
17 between elevations of 3720 and 5170 m a.s.l. Pollen identification and counting were
18 based on standard approaches, and modern climate data were interpolated from a robust
19 modern meteorological dataset. A series of numerical analyses revealed that
20 precipitation is the main climatic determinant of pollen spatial distribution: Cyperaceae,
21 Ranunculaceae, Rosaceae, and *Salix* indicate wet climatic conditions, while Poaceae,
22 *Artemisia*, and Chenopodiaceae represent drought. Model performance of both
23 weighted-averaging partial least squares (WA-PLS) and the random forest (RF)
24 algorithm suggest that this modern pollen dataset has good predictive power in
25 estimating the past precipitation from pollen spectra from the eastern Tibetan Plateau.

26 In addition, a comprehensive modern pollen dataset can be established by combining
27 our modern pollen dataset with previous datasets, which will be essential for the
28 reconstruction of vegetation and climatic signals for fossil pollen spectra on the Tibetan
29 Plateau. Pollen datasets including both pollen counts and percentages for each sample
30 together with their site location and climatic data are available at the National Tibetan
31 Plateau Data Center (TPDC; DOI: 10.11888/Paleoenv.tpdc.271191).

32

33 **1 Introduction**

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

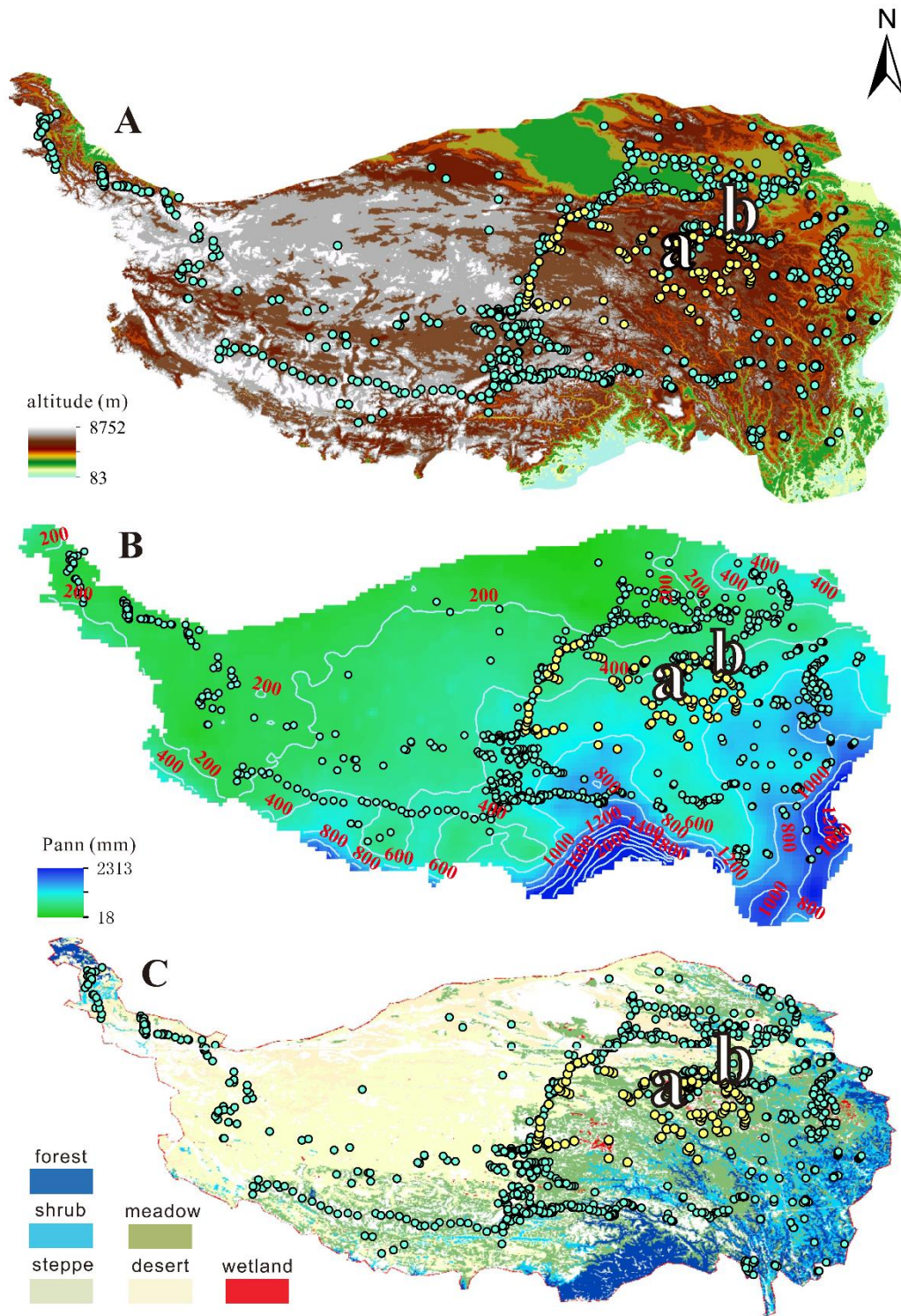
49 The available modern pollen datasets reveal that pollen assemblages on the Tibetan
50 Plateau are generally simple with Cyperaceae, *Artemisia*, Poaceae, and
51 Chenopodiaceae as the dominant taxa (e.g. Herzschuh et al., 2010; Cao et al., 2014),
52 with arboreal pollen taxa becoming more influential in the marginal areas (e.g. Ma et
53 al., 2017; Li et al., 2020). It is essential to identify the climatic indicators of the modern

54 pollen taxa (particular for the four dominant taxa) on the Tibetan Plateau, because the
55 climatic indicators derived from modern pollen datasets from the surrounding lowland
56 cannot be directly employed on the Tibetan Plateau. With our current modern pollen
57 dataset extracted from lake surface-sediments we aim to 1) fill a geographical gap and
58 thus establish a comprehensive modern pollen dataset covering the entire Tibetan
59 Plateau; 2) determine the climatic indicators for common pollen taxa from the alpine
60 meadow ecosystem; and 3) evaluate the predictive power of the modern dataset to
61 reconstruct past climate and assess the reliability of the random forest algorithm in
62 calibrating the pollen-climate relationship.

63

64 **2 Study area**

65 The elevation range of the lakes sampled for our pollen dataset is between 3720 and
66 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
67 75% quantile is 4550 m a.s.l.; Figure 1). Climate of this region is controlled by the
68 Asian Summer Monsoon in summer with warm and wet climatic conditions, and by
69 westerlies in winter with cold and dry conditions (Wang, 2006). The eastern and central
70 Tibetan Plateau containing these sampled lakes (with >4000 m a.s.l elevation) is
71 covered by alpine meadow with sporadic patches of subalpine shrub. The plant
72 communities of the alpine meadow are dominated by *Kobresia* species (Cyperaceae)
73 generally, with Ranunculaceae, Asteraceae, *Polygonum* (Polygonaceae), *Potentilla*
74 (Rosaceae), Fabaceae, and Caryophyllaceae as the common taxa. The subalpine shrub
75 is generally distributed on the northern slopes of mountains with *Salix oritrepha* and
76 *Potentilla fruticosa* as the main shrub components, while the herbaceous taxa
77 mentioned above are also common (Wu, 1995; Herzsuh et al., 2010; unpublished
78 vegetation survey).



79

80 **Figure 1** Spatial distribution of modern pollen samples (yellow dots: the 117 sampled
 81 lakes; bluish green dots: previous samples (surface-soils and lake surface-sediments)
 82 included in the dataset of Cao et al., 2014). A: Digital Elevation Model; B: isohyet map
 83 (mm); C: vegetation map. “a” and “b” indicate the locations of Koucha Lake and
 84 Xingxinghai Lake.

85 3 Materials and methods

86 3.1 Sample collecting and pollen processing

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

99 For pollen extraction, approximately 10 g (wet untreated sediment) per sample were
100 sub-sampled. Pollen samples were processed using standard acid-alkali-acid
101 procedures (including 10% HCl, 10% KOH, 40% HF and 9:1 mixture of acetic
102 anhydride and sulphuric acid successively; Fægri and Iversen, 1975) followed by 7-
103 µm-mesh sieving. A tablet with *Lycopodium* spores (27560 grains/tablet) was added to
104 each sample prior to pollen extraction as tracers (Maher, 1981). Pollen grains were
105 identified with the aid of modern pollen reference slides collected from the eastern and
106 central Tibetan Plateau (including 401 common species of alpine meadow; Cao et al.,
107 2020) and published atlases for pollen and spores (Wang et al., 1995; Tang et al., 2017).
108 More than 500 terrestrial pollen grains were counted for each sample, and more than
109 200 *Lycopodium* spores were counted for most of the samples (mean=270 grains;
110 median=480 grains), both of which ensure a reliable representation of the entire pollen
111 assemblage by the counted pollen data.

112 3.2 Data processing

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

128 To visualize the relationships between modern pollen assemblages and climatic
129 variables, ordination techniques were employed based on the square-root transformed
130 pollen data of 19 taxa (those present in at least 3 samples and with a $\geq 3\%$ maximum)
131 to stabilize variances and optimize the signal-to-noise ratio (Prentice, 1980). Detrended
132 correspondence analysis (DCA; Hill and Gauch, 1980) revealed that the length of the
133 first axis of the pollen data was 1.44 SD (standard deviation units), indicating a linear
134 response model is suitable for our pollen dataset (ter Braak and Verdonschot, 1995).
135 We performed redundancy analysis (RDA) to visualize the distribution of pollen
136 species and sampling sites along the climatic gradients, selecting the minimal adequate
137 model using forward selection and checking the variance inflation factors (VIF) at each
138 step. If VIF values were higher than 20, which indicates that some variables in the
139 model are co-linear, we stopped adding variables (ter Braak and Prentice, 1988). These

140 ordinations were performed using the *decorana* and *rda* functions in the *vegan* package
141 version 2.5-4 (Oksanen et al., 2019) for R.

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

147 The basic assumption of pollen-based past climate reconstruction assumes that pollen
148 taxa recorded in the modern calibration-set have similar ecological requirements as
149 those in the fossil spectra (Juggins and Birks, 2012); in other words, the modern
150 vegetation-climate relationship is assumed to be stable temporally through the target
151 period for reconstruction. To evaluate the potential of the pollen dataset for past climate
152 reconstruction, both the traditional method of weighted-averaging partial least squares
153 (WA-PLS) and a new approach using the random forest (RF) algorithm were run. WA-
154 PLS was performed using the *WAPLS* function in the *rioja* package version 0.7-3
155 (Juggins, 2012) for R using leave-one-out cross-validation, pollen percentages of the
156 19 selected pollen taxa were square-root transformed, and the number of WA-PLS
157 components used was selected using a randomization *t*-test (Juggins and Birks, 2012).
158 We performed the RF algorithm with the *randomForest* package (version 4.6-14; Liaw,
159 2018) in R. RF is an algorithm that integrates multiple decision trees, and the
160 importance of each explanatory variable is measured as the percentage increase in the
161 residual sum of squares after randomly shuffling the order of the variables to determine
162 which explanatory variable can be added to the model. In our study, the importance of
163 all pollen taxa on the spatial distribution of P_{ann} was estimated and the model
164 systematically optimized by a stepwise reduction in variables by deleting the least
165 important one. Our final RF model includes 19 pollen taxa (Appendix B), which all
166 make a positive contribution to the precipitation distribution. To assess the predictive
167 power of our pollen dataset, pollen spectra from Koucha Lake (covering the last 16 cal
168 ka BP (calibrated thousand years before 1950 CE); 34.0°N; 97.2°E, 4540 m a.s.l.;

169 Herzsuh et al., 2009) and Xingxinghai Lake (covering the last 7.5 cal ka BP; 34.8°N,
 170 98.1°E, 4228 m a.s.l.; Zhang et al., unpublished) were selected as the target fossil pollen
 171 datasets for quantitative reconstruction. A statistical significance test for all
 172 reconstructions was performed following the methods described in Telford and Birks
 173 (2011) using the *randomTF* function in the *palaeoSig* package version 1.1.2 for both
 174 WA-PLS and RF reconstruction methods separately (Telford, 2013).

175 4 Data description

176 Pollen assemblages of the dataset from alpine meadows are dominated by Cyperaceae
 177 (mean 68.4%, maximum 95.9%), with other herbaceous pollen taxa common including
 178 Poaceae (mean 10.3%, maximum 87.7%), Ranunculaceae (mean 4.8%, maximum
 179 33.6%), *Artemisia* (mean 3.7%, maximum 24.5%), and Asteraceae (mean 2.1%,
 180 maximum 33.6%). *Salix* (mean 0.4%, maximum 5.3%) is the major shrub taxon in these
 181 pollen assemblages, while arboreal taxa occur with low percentages generally (mean
 182 total arboreal percentage 0.9%, maximum 5.8%), mainly comprising *Pinus* (mean 0.3%,
 183 maximum 1.8%), *Betula* (mean 0.1%, maximum 0.9%), and *Alnus* (mean 0.1%,
 184 maximum 0.7%). Published vegetation data (e.g. Wu, 1995; Herzsuh et al., 2010)
 185 and our vegetation survey reveal that trees are absent from the alpine meadow
 186 communities within the study area, thus we believe the arboreal pollen with low
 187 abundances in the dataset will have been transported by wind from adjacent regions to
 188 the south and east. Generally, these pollen assemblages represent well the plant
 189 components in the alpine meadow communities, although they are influenced slightly
 190 by long-distance pollen transported by wind (Figure 2).

191 **Table 1** Summary statistics for parameters in the pollen dataset. Min.: minimum; Med.:
 192 median; Max.: maximum. Units for longitude and latitude are degrees, elevation is in
 193 m above sea level, Mt_{co} , Mt_{wa} and T_{ann} are °C, P_{ann} is mm, and pollen 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

Pollen taxa	Min.	Med.	Max.	Mean
<i>Abies</i>	0.00	0.00	0.38	0.01
<i>Cedrus</i>	0.00	0.00	0.19	0.00
<i>Picea</i>	0.00	0.00	2.52	0.10
<i>Pinus</i>	0.00	0.18	1.76	0.32
<i>Alnus</i>	0.00	0.00	0.67	0.11
<i>Betula</i>	0.00	0.00	0.94	0.11
<i>Carpinus</i>	0.00	0.00	0.63	0.06
<i>Castanea</i>	0.00	0.00	2.44	0.06
<i>Corylus</i>	0.00	0.00	1.88	0.07
<i>Juglans</i>	0.00	0.00	0.82	0.01
Oleaceae	0.00	0.00	0.16	0.00
<i>Quercus</i>	0.00	0.00	2.00	0.06
<i>Salix</i>	0.00	0.18	5.35	0.45
<i>Ulmus</i>	0.00	0.00	0.16	0.00
Chenopodiaceae	0.00	0.48	15.44	0.86
<i>Ephedra</i>	0.00	0.00	1.66	0.12
Ericaceae	0.00	0.00	0.19	0.01
Euphorbiaceae	0.00	0.00	0.19	0.00
Fabaceae	0.00	0.16	3.07	0.28
<i>Hippophaë</i>	0.00	0.00	5.62	0.27
Rhamnaceae	0.00	0.00	0.17	0.00
Apiaceae	0.00	0.16	3.98	0.32
<i>Artemisia</i>	0.19	2.43	24.51	3.68
Asteraceae	0.00	1.46	33.56	2.09
Brassicaceae	0.00	0.36	28.17	1.22
Caryophyllaceae	0.00	0.16	2.26	0.23
Cyperaceae	4.84	76.24	95.91	68.67
Balsaminaceae	0.00	0.00	0.14	0.00
Urticaceae	0.00	0.00	3.87	0.08
Gentianaceae	0.00	0.16	4.85	0.40
Lamiaceae	0.00	0.00	1.05	0.12
Liliaceae	0.00	0.00	0.50	0.04
Plantaginaceae	0.00	0.00	0.88	0.03
Onagraceae	0.00	0.00	0.34	0.00
Papaveraceae	0.00	0.00	0.82	0.03
Poaceae	0.39	4.90	87.74	10.28
Polemoniaceae	0.00	0.00	15.21	0.34
<i>Polygonum</i>	0.00	0.49	20.50	1.47
<i>Rumex</i>	0.00	0.00	1.64	0.03
<i>Koenigia</i>	0.00	0.00	2.96	0.39
Primulaceae	0.00	0.00	0.56	0.03
Ranunculaceae	0.00	3.47	33.62	4.88
Saxifragaceae	0.00	0.00	4.69	0.10
Scrophulariaceae	0.00	0.00	0.71	0.01
Solanaceae	0.00	0.00	0.69	0.01
<i>Thalictrum</i>	0.00	0.98	12.05	1.45

194

195 The region covered by these modern pollen samples has a P_{ann} gradient from 226 to 689
196 mm, and cold thermal conditions with low T_{ann} (-7.3 to 2.3 °C) and Mt_{co} (-19.2 to
197 -7.4 °C). A series of RDAs reveals that, relative to Mt_{co} and Mt_{wa} , P_{ann} explains more
198 pollen assemblage variation (10.8% as a sole predictor in RDA) in the dataset (Table
199 2). A biplot of the RDA shows that the direction of the P_{ann} vector has a smaller angle
200 with the positive direction of axis 1 (captures 43.2% of total inertia in the dataset) than
201 with the positive direction of axis 2 (10.3%), indicating that the major component of
202 axis 1 should be moisture. RDA axis 1, which is highly correlated with P_{ann} , divides the
203 pollen taxa into two groups generally: Cyperaceae, Ranunculaceae, Rosaceae, and *Salix*
204 indicating wet climatic conditions (located along the positive direction of P_{ann}), while
205 Poaceae, *Artemisia*, and Chenopodiaceae represent drought (located along the negative

206 direction of P_{ann} ; Figure 3). Axis 2 is highly correlated with the two temperature
207 variables; however these dominant pollen taxa have insignificant distributions along
208 the axis, hence temperature is the secondary climatic variable for the pollen dataset
209 relative to precipitation (Figure 3). Because of low occurrences and abundances for
210 some rare pollen taxa, BRT models are only performed for 14 dominant or common
211 pollen taxa. BRT modelling results suggest that P_{ann} is the main climatic determinant
212 for 9 out of 10 of the major pollen taxa with >0.6 prevalence, with Asteraceae an
213 exception having Mt_{co} as its main climatic determinant (68%; Table 3). BRT results
214 reveal that pollen abundances of Cyperaceae, Ranunculaceae, and *Salix* are positively
215 related to P_{ann} , while those of Poaceae, *Artemisia*, and Chenopodiaceae have a negative
216 relationship with P_{ann} , consistent with the RDA results (Figure 3 and 4; Appendix 1).

217

218 **5 Potential use of the modern pollen dataset**

219 Numerical analyses reveal that P_{ann} is the most important climatic determinant of pollen
220 distribution in the eastern Tibetan Plateau, hence, P_{ann} is selected as the target variable
221 in the calibration-set to assess the predictive power of this pollen dataset. Both
222 approaches (WA-PLS, RF) perform well with low RMSEP values (the root mean square
223 error of prediction) and high r^2 values (coefficient of determination between observed
224 and predicted climatic variables; Figure 5). However, the plots of observed vs. predicted
225 P_{ann} show a overestimate of P_{ann} for arid sites and an underestimate for wet sites (Figure
226 5). Hence, the inevitable “edge effects” should be treated with caution. Nevertheless,
227 reconstructions covering ca. 400–500 mm P_{ann} should be reliable because of the low
228 bias in the central part of the P_{ann} gradient (Figure 5).

229 Although the model performance of RF is not any better than that of WA-PLS, the
230 reconstruction produced by RF might be more reliable as suggested by the statistical
231 significance testing and comparison with modern observed P_{ann} for the two lakes
232 (Koucha Lake and Xingxinghai Lake). Statistical significance testing shows that the
233 proportion of variance in the fossil data explained by the WA-PLS reconstruction is

234 less than the 95% quantile of the variance explained by a reconstruction based on
 235 random environmental variables (999 trials) for the two lakes, while reconstructions
 236 produced by RF explain a higher proportion (Figure 6). In other words, reconstructions
 237 produced by RF might be controlled by the major pollen components, because the
 238 explained proportion of variance in the fossil pollen spectra is closer to that explained
 239 by the first PCA axis, while reconstructions by WA-PLS could be influenced more by
 240 the pollen taxa with low abundances (Figure 6). The hypothesis that WA-PLS is
 241 influenced more by low-abundance pollen taxa is supported by the high variation in
 242 reconstructed P_{ann} among the fossil pollen samples (Figure 7). Relative to
 243 reconstructions of WA-PLS, results of RF have lower temporal variation and fewer
 244 outliers, and the predicted P_{ann} by RF is closer to the observed P_{ann} for the two lakes
 245 (Koucha Lake, 500 mm; Xingxinghai Lake, 350 mm) than that by WA-PLS.

246

247 **Table 2** Summary statistics of redundancy analysis (RDA) of 19 pollen species and
 248 four climatic variables. VIF: variance inflation factor; P_{ann} : mean annual precipitation
 249 (mm); Mt_{co} : mean temperature of the coldest month ($^{\circ}\text{C}$); Mt_{wa} : mean temperature of
 250 the warmest month ($^{\circ}\text{C}$); T_{ann} : annual mean temperature ($^{\circ}\text{C}$).

Climatic variables	VIF	VIF	Climatic variables as	Marginal contribution based on	
	(without T_{ann})	(with T_{ann})	sole predictor	Explained variance	p -value
			Explained variance	Explained variance	
			(%)	(%)	
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	-	-	-

251

252

253

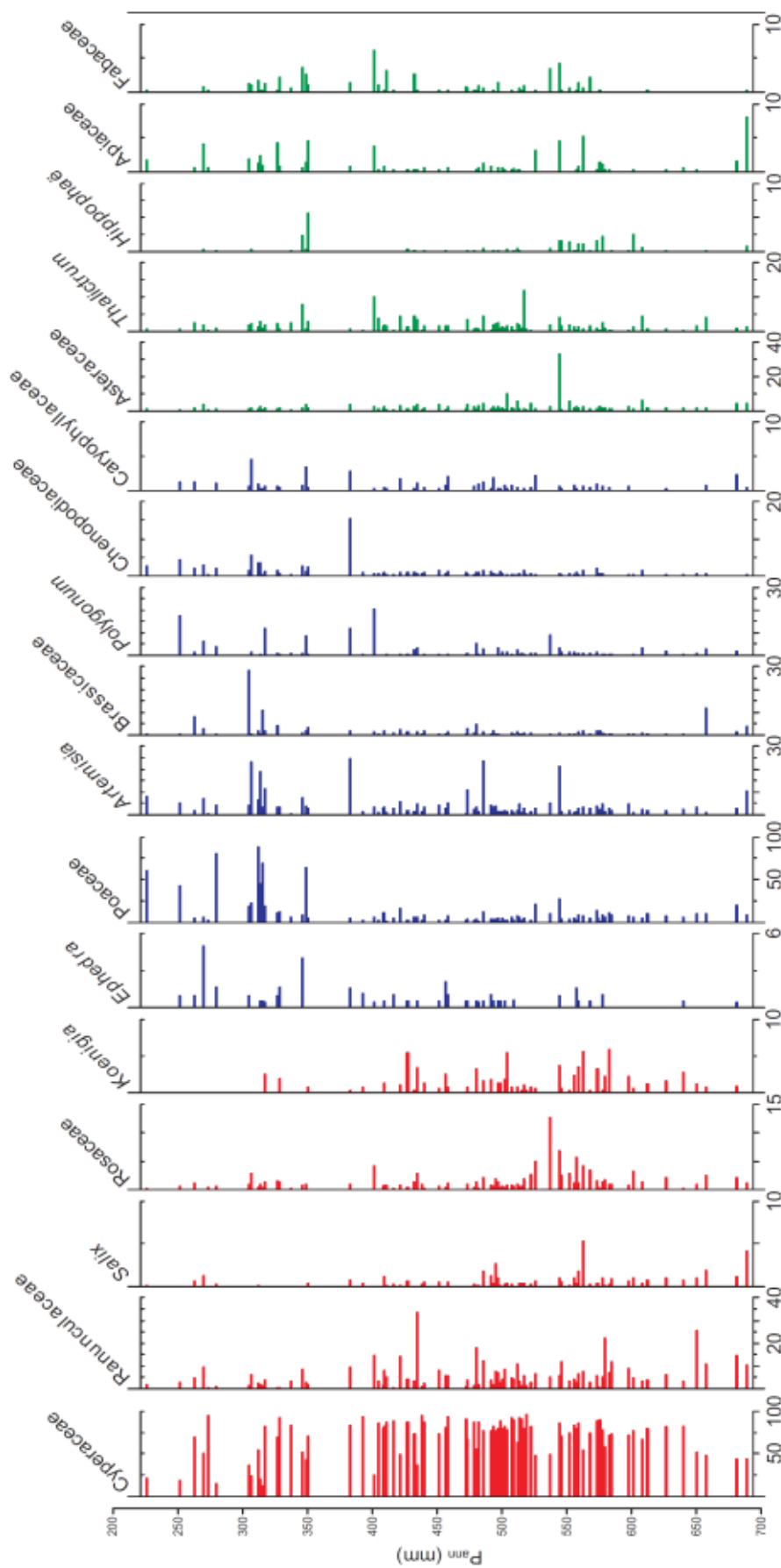
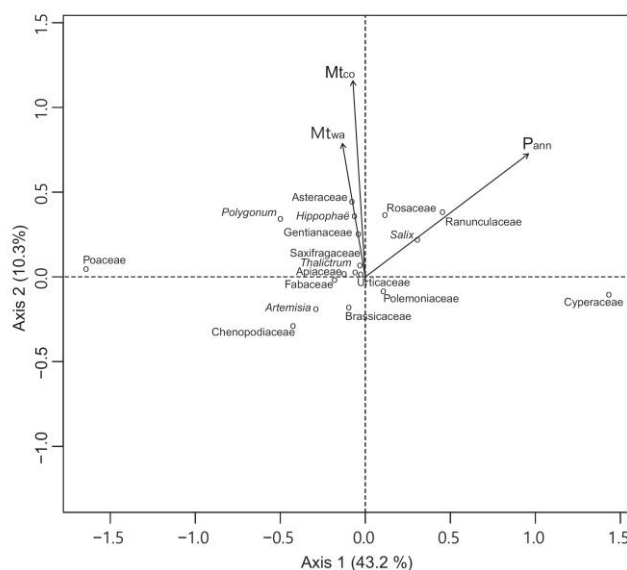


Figure 2 Pollen diagram showing the major taxa (percentage; %) of the 117 samples arranged by mean annual precipitation (P_{ann} ; mm). Pollen taxa with red bars are positively related to P_{ann} , those with blue bars are negatively related to P_{ann} , while the relationship is insignificant for those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

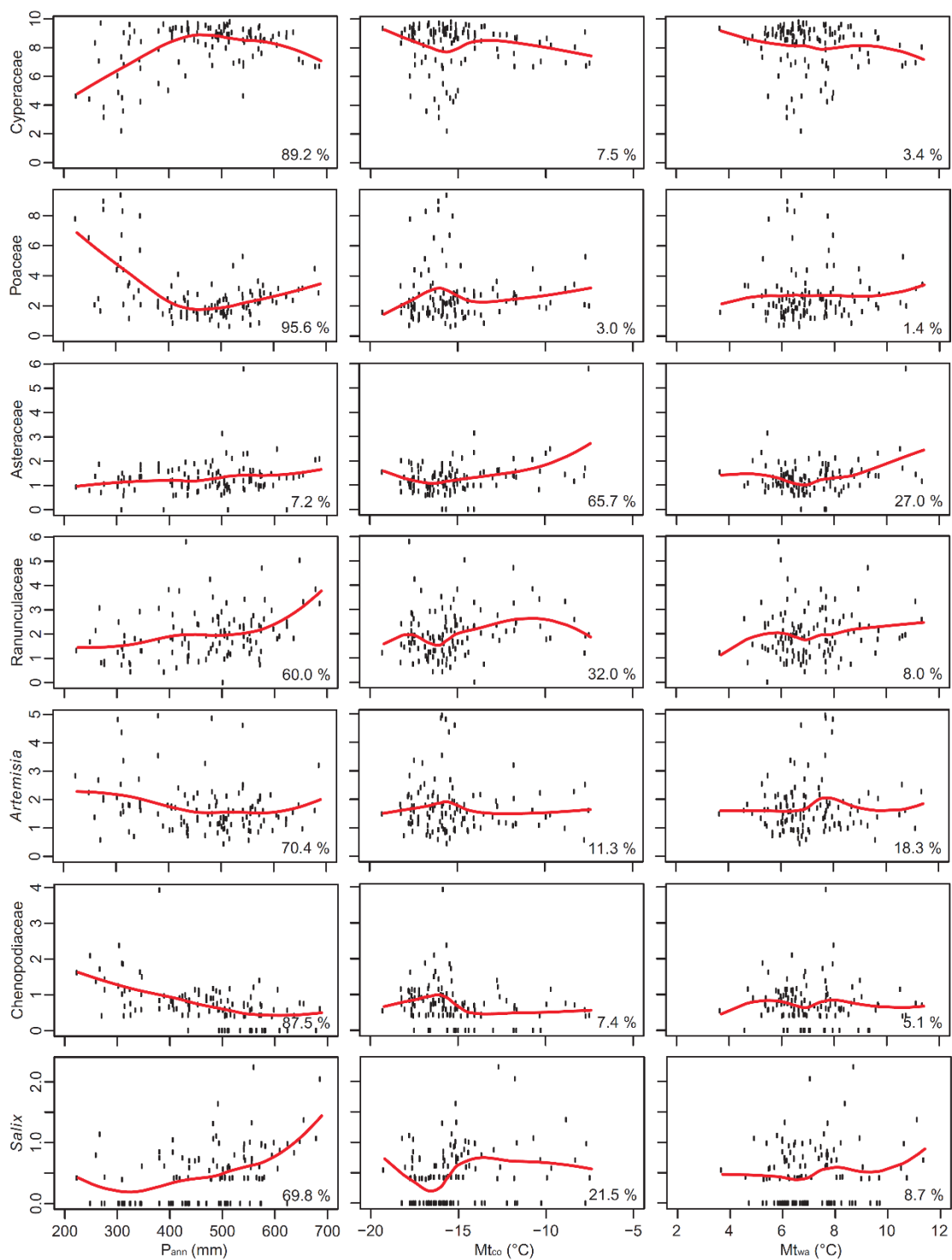


255

256 **Figure 3** Plot of the first two redundancy analysis (RDA) axes showing the
 257 relationships between 18 pollen taxa (circles) and 3 climatic variables (arrows). P_{ann} :
 258 mean annual precipitation (mm); Mt_{co} : mean temperature of the coldest month ($^{\circ}C$);
 259 Mt_{wa} : mean temperature of the warmest month ($^{\circ}C$).

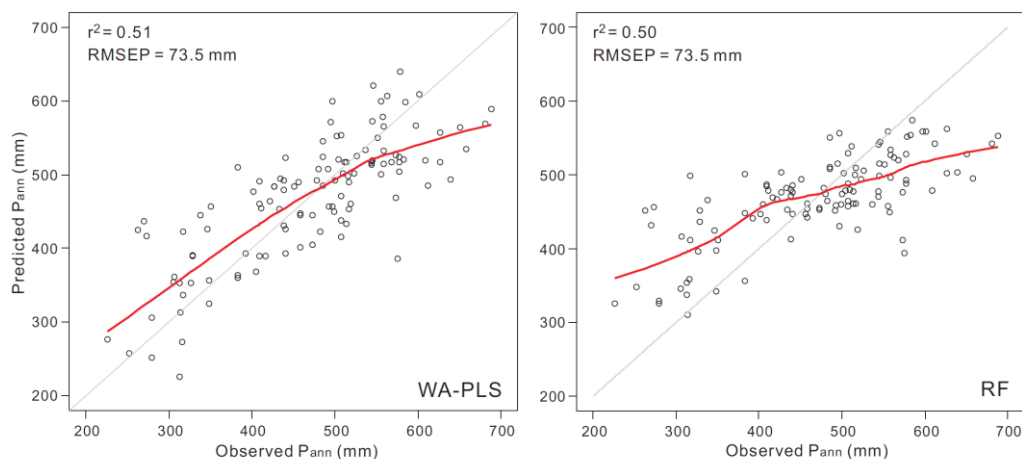
260 **Table 3** Relative influence of climatic variables to the spatial distributions of 14 pollen
 261 taxa based on boosted regression tree (BRT) models. For each variable, the relative
 262 influence is expressed as a percentage among the three variables. Pollen taxa are
 263 ordered by decreasing prevalence (the proportion of sites in which each taxon is
 264 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



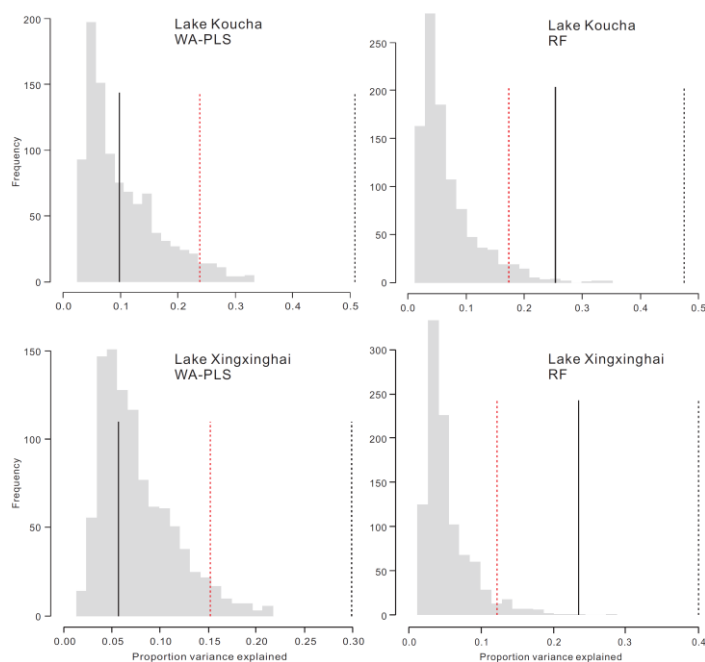
265

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



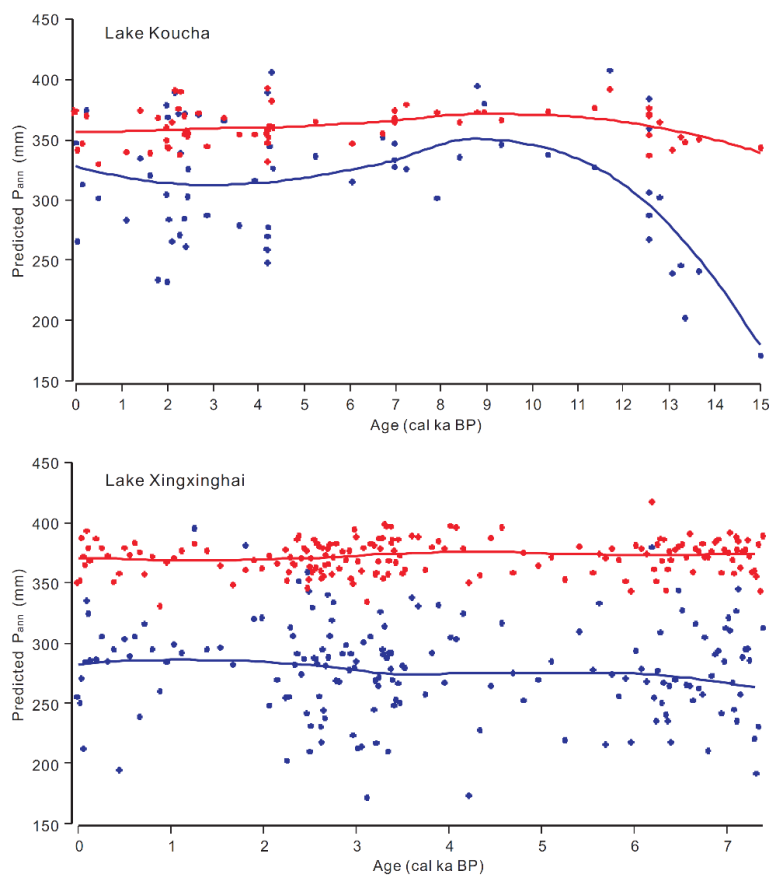
269

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



273

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



280

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

284 **6 Data availability**

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

288 **7 Summary**

289 We present a regional modern pollen dataset extracted from lake surface-sediments
 290 from the alpine meadow vegetation type on the Tibetan Plateau (eastern Tibetan Plateau,
 291 91.8° – 99.8° E and 31.6° – 35.5° N), including pollen counts and pollen percentages
 292 together with their positions and climatic data. Numerical analyses reveal that P_{ann} is
 293 the most important climatic determinant for pollen distribution in the dataset, and our

294 dataset behaves reliably and has good predictive power for past moisture reconstruction,
295 and the random forest algorithm is a potentially reliable approach in pollen-based past
296 environment reconstruction.

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

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

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

313 **References**

314 Birks, H.J.B., Heiri, O., Seppä, H. and Bjune, A.E.: Strengths and weaknesses of
315 quantitative climate reconstructions based on late-Quaternary biological proxies,
316 *Open Ecol. J.*, 3, 68–110, 2010.

317 Cao, X., Tian, F. and Ding, W.: Improving the quality of pollen-climate calibration-sets
318 is the primary step for ensuring reliable climate reconstructions, *Sci. Bull.*, 63,
319 1317–1318, 2018.

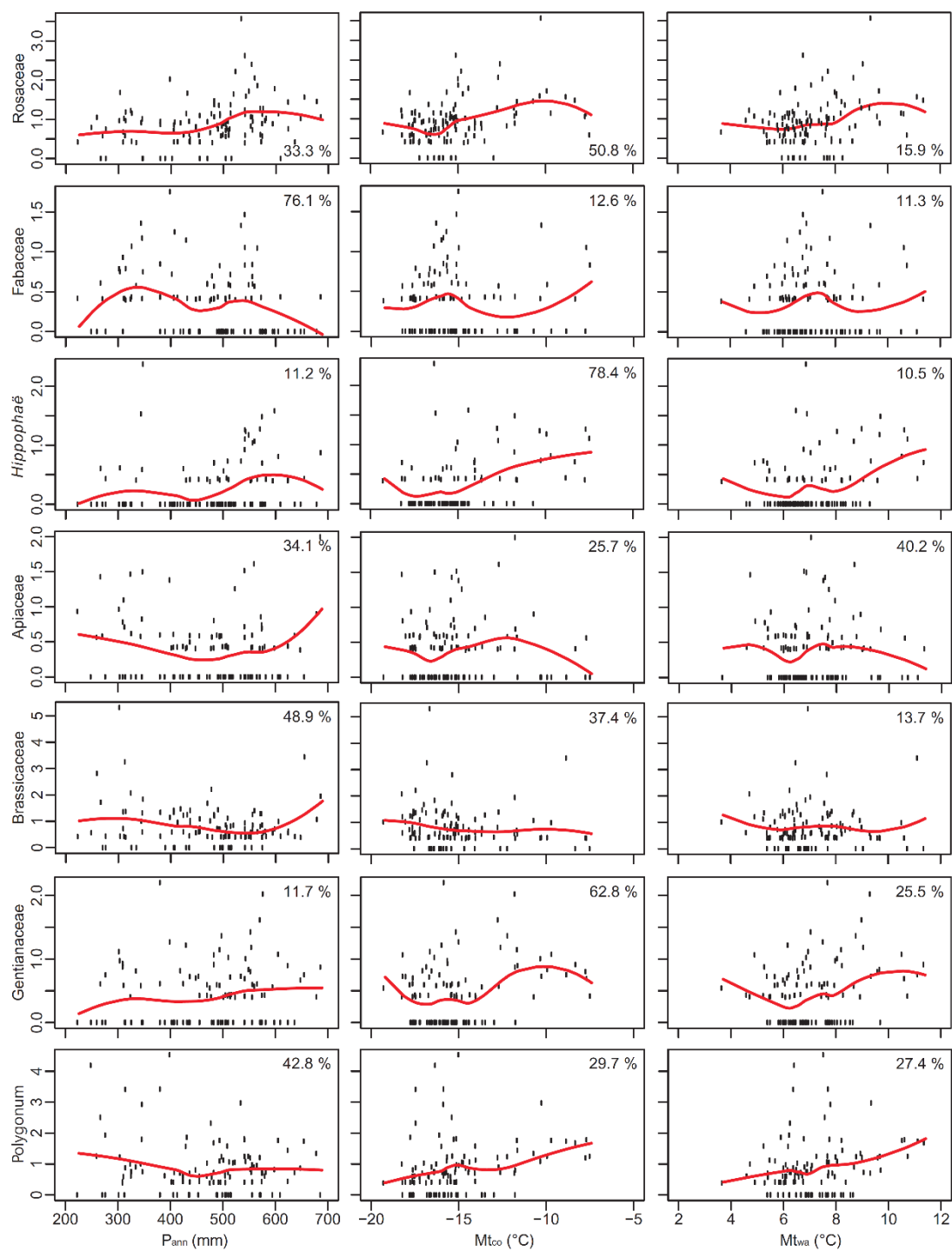
- 320 Cao, X., Tian, F., Li, K. and Ni, J.: Atlas of pollen and spores for common plants from
321 the east Tibetan Plateau. National Tibetan Plateau Data Center, DOI:
322 10.11888/Paleoenv.tpdc.270735, 2020.
- 323 Cao, X.Y., Herzsuh, U., Telford, R.J. and Ni, J.: A modern pollen-climate dataset
324 from China and Mongolia: assessing its potential for climate reconstruction, *Rev.*
325 *Palaeobot. Palynol.*, 211, 87–96, 2014.
- 326 Fægri, K. and Iversen, J.: Textbook of pollen analysis, Munksgaard, Copenhagen, 1975.
- 327 He, J., Yang, K., Tang, W., Lu, H., Qin, J., Chen, Y. and Li, X.: The first high-resolution
328 meteorological forcing dataset for land process studies over China, *Sci. Data*, 7,
329 25, DOI: 10.1038/s41597-020-0369-y, 2020.
- 330 Herzsuh, U., Birks, H.J.B., Mischke, S., Zhang, C. and Böhner, J.: A modern pollen-
331 climate calibration set based on lake sediments from the Tibetan Plateau and its
332 application to a Late Quaternary pollen record from the Qilian Mountains, *J.*
333 *Biogeogr.*, 37, 752–766, 2010.
- 334 Herzsuh, U., Cao, X., Laepple, T., Dallmeyer, A., Telford, R., Ni, J., Chen, F., Kong,
335 Z., Liu, G., Liu, K.-B., Liu, X., Stebich, M., Tang, L., Tian, F., Wang, Y.,
336 Wischnewski, J., Xu, Q., Yan, S., Yang, Z., Yu, G., Zhang, Y., Zhao, Y. and Zheng,
337 Z.: Position and orientation of the westerly jet determined Holocene rainfall
338 patterns in China, *Nat. Commun.*, 10, 2376, 2019.
- 339 Herzsuh, U., Kramer, A., Mischke, S. and Zhang, C.: Quantitative climate and
340 vegetation trends since the late glacial on the northeastern Tibetan Plateau deduced
341 from Koucha Lake pollen spectra. *Quat. Res.*, 71, 162–171, 2009.
- 342 Hijmans, R.J., Phillips, S., Leathwick, J. and Elith, J.: Dismo: Species Distribution
343 Modeling, version 1.0-12, available at: [http://CRAN.R-project.org/package/
344 dismo](http://CRAN.R-project.org/package/dismo), 2015.

- 345 Hill, M.O. and Gauch, H.G.: Detrended correspondence analysis: an improved
346 ordination technique, *Vegetatio*, 42, 41–58, 1980.
- 347 Juggins, S. and Birks, H.J.B.: Quantitative environmental reconstructions from
348 biological data, in: Birks, H.J.B., Lotter, A.F., Juggins, S. and Smol, J.P. (eds.),
349 Tracking environmental change using lake sediments (vol. 5): Data handling and
350 numerical techniques, Springer, Dordrecht, 431–494, 2012.
- 351 Juggins, S.: Rioja: analysis of Quaternary Science Data version 0.7-3, available at:
352 <http://cran.r-project.org/web/packages/rioja/index.html>, 2012.
- 353 Li, J.F., Xie, G., Yang, J., Ferguson, D.F., Liu, X.D., Liu, H. and Wang, Y.F.: Asian
354 Summer Monsoon changes the pollen flow on the Tibetan Plateau, *Earth-Sci. Rev.*,
355 202, 103114, 2020.
- 356 Liaw, A.: randomForest: Breiman and Cutler's Random Forests for Classification and
357 Regression, available at: [https://cran.r-project.org/web/packages/randomForest/](https://cran.r-project.org/web/packages/randomForest/index.html)
358 [index.html](https://cran.r-project.org/web/packages/randomForest/index.html), 2018.
- 359 Ma, Q., Zhu, L., Wang, J., Ju, J., Lü, X., Wang, Y., Guo, Y., Yang, R., Kasper, T.,
360 Haberzettl, T. and Tang, L.: *Artemisia/Chenopodiaceae* ratio from surface lake
361 sediments on the central and western Tibetan Plateau and its application,
362 *Palaeogeogr. Palaeoclim. Palaeoecol.*, 479, 138–145, 2017.
- 363 Maher, L.J.: Statistics for microfossil concentration measurements employing samples
364 spiked with marker grains, *Rev. Palaeobot. Palynol.*, 32, 153–191, 1981.
- 365 Nychka, D., Furrer, R., Paige, J. and Sain, S.: fields: Tools for spatial data, version 9.6.1,
366 available at: <https://cran.r-project.org/web/packages/fields/>, 2019.
- 367 Oksanen, J., Blanchet, F.G., Friendly, M., Kindt, R., Legendre, P., McGlenn, D.,
368 Minchin, P.R., O'Hara, R.B., Simpson, G.L., Solymos, P., Stevens, M.H.H.,
369 Szoecs, E. and Wagner, H.: vegan: Community Ecology Package, version 2.5-4,
370 available at: <https://cran.r-project.org/web/packages/vegan/index.html>, 2019.

- 371 Prentice, I.C.: Multidimensional scaling as a research tool in Quaternary palynology: a
372 review of theory and methods, *Rev. Palaeobot. Palynol.*, 31, 71–104, 1980.
- 373 R Core Team: R, A language and environment for statistical computing, R Foundation
374 for Statistical Computing, Vienna, 2019.
- 375 Shen, C., Liu, K.B., Tang, L. and Overpeck, J.T.: Quantitative relationships between
376 pollen rain and climate in the Tibetan Plateau. *Rev. Palaeobot. Palynol.*, 140, 61–
377 77, 2006.
- 378 Sugita, S.: A model of pollen source area for an entire lake surface. *Quat. Res.*, 39, 369–
379 244, 1993.
- 380 Tang, L., Mao, L., Shu, J., Li, C., Shen, C. and Zhou, Z.: Atlas of Quaternary pollen
381 and spores in China, Science Press, Beijing, 2017.
- 382 ter Braak, C.J.F. and Prentice, I.C.: A theory of gradient analysis, *Adv. Ecol. Res.*, 18,
383 271–317, 1988.
- 384 ter Braak, C.J.F. and Verdonschot, P.F.M.: Canonical correspondence analysis and
385 related multivariate methods in aquatic ecology, *Aquat. Sci.*, 57, 255–289, 1995.
- 386 Tian, F., Xu, Q., Li, Y., Cao, X., Wang, X. and Zhang, L.: Pollen assemblage
387 characteristics of lakes in the monsoon fringe area of China. *Chinese Sci. Bull.*,
388 53(21), 3354–3363, 2008.
- 389 Wang, B.: The Asian Monsoon, Springer, Chichester, 2006.
- 390 Wang, F.X., Qian, N.F., Zhang, Y.L. and Yang, H.Q.: Pollen Flora of China, Science
391 Press, Beijing, 1995.
- 392 Wu, Z.Y.: The vegetation of China. Science Press, Beijing, 1995 (in Chinese).
- 393
- 394

395 Appendix A

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



399

400

401 Appendix B

402 Importance (imp) of pollen taxa on the spatial distribution of P_{ann} was repeatedly
 403 assessed by the random forest algorithm (RF). Shown in bold are the pollen taxa
 404 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
