



HAL
open science

FROG: A global machine-learning temperature calibration for branched GDGTs in soils and peats

Pierre Véquaud, Alexandre Thibault, Sylvie Derenne, Christelle Anquetil, Sylvie Collin, Sergio Contreras, Andrew T Nottingham, Pierre Sabatier, Josef P Werne, Arnaud Huguet

► **To cite this version:**

Pierre Véquaud, Alexandre Thibault, Sylvie Derenne, Christelle Anquetil, Sylvie Collin, et al.. FROG: A global machine-learning temperature calibration for branched GDGTs in soils and peats. *Geochimica et Cosmochimica Acta*, 2022, 318, pp.468-494. 10.1016/j.gca.2021.12.007 . hal-03552122

HAL Id: hal-03552122

<https://hal.sorbonne-universite.fr/hal-03552122v1>

Submitted on 2 Feb 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

FROG: A GLOBAL MACHINE-LEARNING TEMPERATURE CALIBRATION FOR BRANCHED GDGTS IN SOILS AND PEATS

Pierre Véquaud^a, Alexandre Thibault^b, Sylvie Derenne^a, Christelle Anquetil^a, Sylvie Collin^a,
Sergio Contreras^c, Andrew T. Nottingham^{d,e}, Pierre Sabatier^f, Josef P. Werne^g, Arnaud
Huguet^{a*}

^aSorbonne Université, CNRS, EPHE, PSL, UMR METIS, Paris, 75005, France

^bAntea Group, Innovation Hub, 803 boulevard Duhamel du Monceau, Olivet, 45160, France

^cDepartamento de Química Ambiental, Facultad de Ciencias & Centro de Investigación en Biodiversidad y
Ambientes Sustentables (CIBAS), Universidad Católica de la Santísima Concepción, Casilla 297, Concepción,
Chile

^dSchool of Geosciences, University of Edinburgh, Crew Building, Kings Buildings, Edinburgh EH9 3FF United
Kingdom

^eSchool of Geography, University of Leeds, Leeds, United Kingdom

^fUniv. Savoie Mont Blanc, CNRS, EDYTEM, Le Bourget du Lac, 73776, France

^gDepartment of Geology and Environmental Science, University of Pittsburgh, Pittsburgh, PA 15260, USA

Abstract

Branched glycerol dialkyl glycerol tetraethers (brGDGTs) are a family of bacterial lipids which have emerged over time as robust temperature and pH paleoproxies in continental settings. Nevertheless, it was previously shown that other parameters than temperature and pH, such as soil moisture, thermal regime or vegetation can also influence the relative distribution of brGDGTs in soils and peats. This can explain a large part of the residual scatter in the global brGDGT calibrations with mean annual air temperature (MAAT) and pH in these settings. Despite improvements in brGDGT analytical methods and development of refined models, the root-mean-square error (RMSE) associated with global calibrations between brGDGT distribution and MAAT in soils and peats remains high (~ 5 °C). The aim of the present study was to develop a new global terrestrial brGDGT temperature calibration from a worldwide extended dataset (i.e. 775 soil and peat samples, i.e. 112 samples added to the previously available global calibration) using a machine learning algorithm. Statistical analyses highlighted five clusters with different effects of potential confounding factors in addition to MAAT on the relative abundances of brGDGTs. The results also revealed the limitations of using a single index and a simple linear regression model to capture the response of brGDGTs to temperature changes. A new improved calibration based on a random forest algorithm was

* Corresponding author. Tel: + 33-144-275-172; fax: +33-144-275-150.

E-mail address: arnaud.huguet@sorbonne-universite.fr (A. Huguet).

45 thus proposed, the so-called random **F**orest **R**egression for Pale**O**MAAT using brGDGTs
46 (**FROG**). This multi-factorial and non-parametric model allows to overcome the use of a single
47 index, and to be more representative of the environmental complexity by taking into account
48 the non-linear relationships between MAAT and the relative abundances of the individual
49 brGDGTs. The FROG model represents a refined brGDGT temperature calibration ($R^2 = 0.8$;
50 $RMSE = 4.01^\circ C$) for soils and peats, more robust and accurate than previous global soil
51 calibrations while being proposed on an extended dataset. This novel improved calibration was
52 further applied and validated on two paleo archives covering the last 110 kyr and the Pliocene,
53 respectively.

54

55 **Keywords:** branched GDGTs; global temperature calibration; soil; peat; machine
56 learning

57

58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91

1. Introduction

Investigating past climate variations is essential to understand and predict future environmental changes, especially in the context of global anthropogenic changes. To this aim, "indirect" indicators of past climates – so-called proxies – have been developed and used regularly since the last century, including those based on microbial lipids. Microorganisms are able to modify the lipid composition of their membranes to maintain a functional fluidity and permeability of the latter. The temperature and pH of the microorganism living environment are considered to be the predominant factors influencing the membrane lipid distribution (Lauber et al., 2009; Siles and Margesin, 2016; Hofmann et al., 2016; Shen et al., 2019).

Among microbial lipids, branched glycerol dialkyl glycerol tetraethers (brGDGTs) have been increasingly used as temperature and pH proxies in continental settings over the last 15 years. These membrane lipids are produced by still unidentified bacteria, although some of them may belong to the phylum *Acidobacteria* (Sinninghe Damsté et al., 2011, 2014, 2018). They are ubiquitous in terrestrial (Weijers et al., 2007; Peterse et al., 2012; De Jonge et al., 2014; Naafs et al., 2017a) and aquatic environments (Blaga et al., 2009; Peterse et al., 2009, 2015; Damsté et al., 2009; Tierney and Russell, 2009; Loomis et al., 2012; Weber et al., 2015). The analysis of brGDGTs, based on a large number of soils distributed worldwide showed that the relative distribution of these compounds is mainly related to mean annual air temperature (MAAT) and soil pH (Weijers et al., 2007; Peterse et al., 2012; De Jonge et al., 2014). The average number of pentane rings, reflected in the Cyclisation of Branched Tetraethers (CBT) index (Weijers et al., 2007; Peterse et al., 2012), has been correlated with soil pH, while the average number of methyl groups, referred to as the Methylation of Branched Tetraethers (MBT) index, has been initially correlated with mean annual mean air temperature (MAAT) and, to a lesser extent, soil pH (Weijers et al., 2007 ($r^2=0.77$; RMSE =4.8°C, $n=134$); Peterse et al., 2012 ($r^2=0.59$; RMSE=5.0°C; $n=176$)). More recently, new brGDGT isomers have been detected through improved analytical methods, with methyl groups being present in either 5th, 6th, 7th or 8th position (De Jonge et al., 2013, 2014; Hopmans et al., 2016; Ding et al., 2016). It was observed that 6-methyl isomers were strongly and predominantly dependent on soil pH and 5-methyl brGDGTs on temperature (De Jonge et al., 2014). This led to the development of a new MBT index excluding 6-methyl isomers - the MBT_{5Me} index - which correlates preferentially with MAAT ($r^2 = 0,64$; RMSE= 4.8°C; $n = 231$).

Hence, brGDGTs have emerged over time as robust temperature and pH paleoproxies in multiple types of settings – lakes (Blaga et al., 2009; Powers et al., 2010; Fawcett et al., 2011;

92 Harning et al., 2020), peatlands (Weijers et al., 2011b; Coffinet et al., 2018; Wu et al., 2020),
93 soil/paleosols (Ding et al., 2015 ; Lu et al., 2016 ; Feng et al., 2019 ; Wang et al., 2020) and
94 speleothems (Baker et al., 2019). Nevertheless, numerous studies showed that additional
95 parameters other than temperature and pH, such as soil moisture, thermal regime or vegetation
96 cover, may also influence the relative distribution of brGDGTs in peat and soils (Weijers et al.,
97 2011a ; Dirghangi et al., 2013; Huguet et al., 2010, 2013 ; Menges et al., 2014; Davtian et al.,
98 2016 ; Liang et al., 2019). This can explain a large part of the residual scatter in the global
99 brGDGT soil/peat calibrations with MAAT and pH (De Jonge et al., 2014; Naafs et al., 2017a,
100 b; Dearing Crampton-Flood et al., 2020).

101 To overcome these limitations, refinements in the global brGDGT calibrations were
102 proposed over the years. De Jonge et al (2014) developed a global temperature calibration (R^2
103 = 0.64; RMSE = 4.8°C; $n = 231$) based on the MBT'_{5ME} index, excluding the 6-methyl
104 brGDGTs isomers. Naafs et al. (2017a) showed that a stronger correlation between MBT'_{5ME}
105 and MAAT in soils could be obtained by excluding samples dominated by 6-methyl brGDGTs
106 (i.e. ratio of 5- vs. 6-methyl brGDGTs (IR_{6Me}) higher than 0.5; $n = 177$; $R^2 = 0.76$; RMSE =
107 4.1°C). More recently, Dearing Crampton-Flood et al. (2020) used Bayesian statistics instead
108 of more classically applied single linear regressions to investigate the relationship between
109 MBT'_{5Me} and MAAT in soils ($n = 353$; $R^2 = 0.64$; RMSE = 6 °C). Samples with IR_{6Me}>0.5
110 were included in this calibration, as excluding them did not significantly change the strength of
111 the correlation with MAAT. The robustness of the Bayesian approach relies on the fact that (i)
112 it considers a given index (e.g., the MBT'_{5Me} for brGDGTs) as the variable dependent on
113 environmental parameters, consistent with the fact that bacterial lipids are produced in response
114 to the variations of environmental parameters and that (ii) it avoids regression dilution
115 phenomena, in contrast with most of the models based on linear regressions.

116 Despite improvements in brGDGT analytical methods and development of refined
117 models, the RMSE associated with global calibrations between brGDGT distribution and
118 MAAT in soils and peat remains high (> 4 °C). Part of this uncertainty may be related to our
119 lack of understanding of the mechanism behind the relationship between MAAT and brGDGT
120 distribution. This relationship has initially been explained by a biophysiological mechanism,
121 i.e. the adjustment of the membrane lipid composition by the brGDGT-producing bacteria in
122 response to changes in environmental conditions (homeoviscous adaptation; Weijers et al.,
123 2007). Nevertheless, changes in bacterial community composition may also explain changes in
124 brGDGT distribution, as recently shown by the lipid characterization of Acidobacterial cultures
125 (Sinninghe Damsté et al., 2018) and the concomitant study of brGDGTs and bacterial

126 community composition in soils from well-documented experimental sites (De Jonge et al.,
127 2019, 2021).

128 Moreover, most of the previous global brGDGT calibrations in soils were based on a
129 correlation between MAAT and a single index (i.e. MBT_{5Me}; De Jonge et al., 2014; Naafs et
130 al., 2017a, b; Dearing Crampton-Flood et al., 2020), even though the relative distribution of
131 brGDGTs is likely to be concomitantly influenced by several environmental parameters. In
132 contrast, using relative abundances of bacterial lipids rather than a single index in models
133 appears more representative of the environmental complexity (Wang et al., 2020; Véquaud et
134 al., 2020; Dunkley Jones et al., 2020). In this way, multiple regression models were also used
135 to describe the relationships between brGDGT distribution and given environmental variables
136 (MAAT, pH) in soils (e.g. Peterse et al., 2012; De Jonge et al., 2014) or lakes (e.g. Pearson et
137 al., 2011; Russell et al., 2018). It was previously shown that the uncertainty in brGDGT
138 calibrations can be improved through the use of multiple regression methods vs. single predictor
139 methods (e.g. Loomis et al., 2012; Wang et al., 2020). Nevertheless, as other linear models, the
140 multiple regression ones cannot take into account non-linear influences, which may occur in
141 complex environmental settings. Such a limitation can be overcome using non-parametric
142 models such as machine-learning algorithms. Machine-learning models were very recently used
143 to develop global calibrations between the relative abundance of isoprenoid GDGTs and sea
144 surface temperature (SST) in marine settings (Dunkley Jones et al., 2020) and between the
145 relative abundance of bacterial 3-hydroxy fatty acids and MAAT in soils (Véquaud et al., 2020;
146 Wang et al., 2021). These models allow overcoming the use of a single index as they are based
147 on the whole suite of microbial lipids. They are built on a proportion of the total dataset
148 (randomly defined) and then tested on the rest of the dataset, considered as independent. Such
149 an approach improves the robustness of the model and avoids the phenomenon of regression
150 dilution. As they are non-parametric, they also capture non-linear environmental influences, in
151 line with the intrinsic complexity of the environmental settings.

152 In the present study, a machine-learning algorithm (random forest) was developed with
153 the aim of proposing a new global brGDGT calibration for MAAT reconstruction in soils and
154 peats with a reduced RMSE. It was based on an extended global dataset comprising 775 peat
155 and soil samples (with 112 samples added to the previous global brGDGT calibration by
156 Dearing Crampton-Flood et al., 2020). This dataset was statistically separated into clusters to
157 better understand the parameters affecting brGDGT distribution in soils at the global scale. The
158 clusters differed by the influence of environmental parameters – MAAT, mean annual
159 precipitation (MAP), soil pH and the number of frozen days during the year (FRS) on the

160 relative abundance of brGDGTs. This mechanistic approach highlighted the limitations of the
161 MBT'_{5Me}-MAAT relationship at the global scale and then led to the development of a refined
162 brGDGT temperature calibration (so-called FROG model) based on a random forest machine-
163 learning algorithm and the whole suite of individual brGDGTs. Alternative models were also
164 proposed to test the influence of confounding variables on the FROG calibration and potentially
165 further improve its accuracy.

166

167 **2. Materials and methods**

168 **2.1. Global soil dataset and environmental parameters**

169 The dataset of the present study is comprised of the globally distributed surface peat
170 and soil samples ($n=663$) used in previous brGDGT global calibrations (Weijers et al., 2007;
171 Peterse et al., 2012; De Jonge et al., 2014; Yang et al., 2015; Ding et al., 2015; Xiao et al., 2015;
172 Lei et al., 2016; Wang et al., 2016; Naafs et al., 2017b; Dearing Crampton-Flood et al., 2020).
173 This dataset was extended with 112 soil samples from 6 altitudinal transects located in France,
174 Italy, Tibet, Chile and Peru and for which brGDGT data were recently published (Huguet et al.,
175 2019; Véquaud et al., 2020, 2021). The details of the dataset ($n=775$) are provided in Table 1.

176 Actual MAAT, pH and MAP values measured from the nearest weather stations, when
177 available, were used to better determine the environmental reality, diversity and complexity.
178 Such values were available for most of the samples, i.e. 598 of the 775 samples. Nevertheless,
179 for the other samples, MAAT and pH values were extracted from the 0.5 gridded CRU TS v.
180 3.26 dataset (Harris et al., 2014), using the same approach as Dearing-Crampton Flood et al.
181 (2020). This approach would have been inappropriate for the 6 aforementioned altitudinal
182 transects, where large temperature variations derive from differences in elevation that can vary
183 across short distances, as noticed by Pérez-Angel et al. (2020).

184 To constrain the applicability of the MBT'_{5Me} as a temperature proxy in peat and soils,
185 Naafs et al. (2017a) used a thermal regime indicator, the Growing Degree Days (GDD). This
186 index is calculated by summing the daily temperatures above 0 °C over a year within a soil
187 (Choler 2018) and interpreted as a proxy of heat accumulation within the latter (McMaster and
188 Wilhelm, 1997; Choler, 2018). The GDD better reflects the growth temperatures encountered
189 by bacterial communities in soils and peats. Unfortunately, as the daily temperatures were not
190 available for the whole dataset of the present study, the GDD index could not be calculated.
191 Instead, another thermal regime indicator was used, the number of frozen days during the year
192 (FRS) for one location, proposed by Harris et al., (2014). The FRS was obtained for most of the
193 samples (i.e. those with site coordinates available, Table 1).

194 The MBT'_{5Me} index, reflecting the methylation level in 5-methyl isomers of GDGTs
195 and considered as related to MAAT, was calculated according to De Jonge et al. (2014; Eq. 1):

196

$$197 \quad MBT'_{5Me} = \frac{[Ia+Ib+Ic]}{[Ia+Ib+Ic] + [IIa+IIb+IIc] + [IIIa]} \quad (1)$$

198 The CBT' index was calculated as follows (De Jonge et al., 214; Eq. 2):

199
$$\text{CBT}' = \log \left(\frac{[Ic] + [IIa'] + [IIb'] + [IIc'] + [IIIa'] + [IIIb'] + [IIIc']}{[Ia] + [IIa + IIIa]} \right) \quad (2)$$

200 The IR_{6Me} reflects the relative abundance of 6- vs. 5-methyl brGDGTs, as proposed by Dang et
 201 al. (2016; Eq. 3):

202
$$\text{IR}_{6\text{Me}} = \log \left(\frac{[IIa'] + [IIb'] + [IIc'] + [IIIa'] + [IIIb'] + [IIIc']}{[IIa'] + [IIb' + IIc'] + [IIIa'] + [IIIb'] + [IIIc'] + [IIa] + [IIb] + [IIc] + [IIIa] + [IIIb] + [IIIc]} \right) \quad (3)$$

203
 204 The Roman numerals correspond to the different GDGT structures presented in De
 205 Jonge et al. (2014). The 6-methyl brGDGTs are denoted by an apostrophe after the Roman
 206 numerals for their corresponding 5-methyl isomers.

207

208 **2.2. Statistical analyses**

209 A Principal Component Analysis (PCA) was performed on the entire dataset with R
 210 software (version 4.0.3; R Core Team, 2020) to observe the distribution of the different samples
 211 based on their brGDGT relative abundances. A cluster classification of the samples based on
 212 the k-means method was proposed. In order to choose the optimal number of clusters, the ratio
 213 of Within-Cluster-Sum-of-Squares (WCSS) over the total sum of squares was calculated. The
 214 WCSS is the sum of squares of the distances of each data point in all clusters to their respective
 215 centroids. The optimal number of clusters corresponds to the minimum value of the ratio of the
 216 WCSS over the total sum of squares (a WCSS = 0 means one sample corresponds to one
 217 cluster). In order to choose the threshold for the optimal WCSS value, and so the optimal
 218 number of clusters, the elbow method was used. It consists in plotting the WCSS values against
 219 the number of clusters, then allowing to derive the optimal number of clusters.

220 Redundancy analysis (RDA) was first performed on the global dataset and then carried
 221 out on each cluster derived from the PCA analysis to evaluate and compare the influence of the
 222 environmental parameters on brGDGT distribution (i) at the global scale and (ii) in each cluster.
 223 RDA is a "constrained" analysis, used to directly visualize the variation in the lipid data as a
 224 function of the environmental variables. It allows not only assessing but also quantifying the
 225 influence of each explanatory variable (i.e. environmental variables) on the distribution of
 226 bacterial lipids. RDA yields the influence of each variable, with regard to the statistical
 227 variance, on the pool of bacterial lipids, and allowed a quantification in percent of the influence
 228 of each parameter (i.e. conditional effect). Conditional effects summarize the effects of each
 229 variable taking into account the effect of variables with the greatest influence (Braak and
 230 Smilauer, 2002). RDA analyses were performed on centered and standardized data using the

231 CANOCO v. 5.04 software (Braak and Smilauer, 2002). The relationships between each
232 variable and the dimensions of the RDA were investigated using the corresponding r-values
233 and the percentages of variance.

234 In order to refine the threshold of the community index (CI) proposed by De Jonge et
235 al. (2019), all linear regressions between MAAT and MBT'_{5Me} on the global dataset were tested
236 by successive iteration of CI values from 0 to 1 (0.001 step) using the R software, version 4.0.3
237 (R Core Team, 2014).

238

239 **2.3. Machine learning: Random forest model**

240

241 The random forest algorithm was used to develop a global calibration between
242 brGDGT relative abundances and MAAT. The random forest algorithm is a supervised learning
243 method notably used for regressions (e.g. Ho, 1995; Breiman, 2001; Denisko and Hoffman,
244 2018). This model works by building a multitude of decision trees from a training dataset and
245 producing the mean prediction of the individual trees. Decision tree learning is one of the
246 predictive modeling approaches used to move from observations to conclusions about the target
247 value of an item.

248 In order to calculate the model based on the random forest algorithm, the global dataset
249 was divided into two subsets: a training dataset and a test dataset. The training dataset
250 corresponds to the samples used to fit the model. The test dataset corresponds to the samples
251 used to provide an unbiased evaluation of the model previously fit on the training dataset. The
252 training phase required for the random forest regression was performed on 75% of the sample
253 set (which allow to neglect the overfit of the model), with 500 trees and an iteration of ten-fold
254 cross-validations per model. The cross validation allows the optimization of the
255 hyperparameters (number of variables in each node and minimal node size) of the models. Data
256 selection was performed randomly on the dataset, but with a stratification modality according
257 to the MAAT to limit the impact of extreme values. Then, the robustness and precision of the
258 different models, developed from the random forest algorithm, were tested on the remaining 25
259 % of samples, considered as an independent dataset. Random forest models were performed
260 with R software, version 4.0.3 (R Core Team, 2014) using the packages tidymodels (version
261 0.1.02)- ranger (version 0.12.1).

262 A R package with a web-application is available on a GITHUB repository ([paleoFROG](#))
263 for the reconstruction of brGDGT-derived MAAT using the FROG models proposed in the
264 present study.

265 The performances of the random forest model were compared with those of the
266 Bayesian models, BayMBT and BayMBT₀, proposed by Dearing Crampton-Flood et al. (2020).
267 The latter were performed with the MATLAB code available from the GITHUB repository of
268 Jessica Tierney (<https://github.com/jesstierney>; Dearing Crampton-Flood et al., 2020), using
269 MATLAB, version 9.8. The *prior* mean was set to the MAAT mean for all soil samples (10°C),
270 with a *prior* standard deviation of 30°C.

271

272 **3. Results**

273 **3.1. Principal component analysis and clustering on the global dataset**

274 In order to explore the global dataset and understand which samples could explain the
275 scattering on the global MBT'_{5Me}-MAAT calibration, we performed a statistical clustering of
276 the extended peat and soil dataset without any *a priori* assumptions on the basis of their
277 brGDGT fractional abundances. With this aim, a Principal Component Analysis (PCA) was
278 performed on the entire brGDGT dataset (Fig. 1). The first 3 axes of the PCA carry most of the
279 variance (70.3%; Figs. 1a, b, c). Consequently, the description of the analysis will be restricted
280 to these axes. A cluster classification of the samples based on the k-means method was
281 performed, yielding 5 clusters (Fig. 1d), based on the Within Cluster Sum of squares and the
282 elbow method. The distribution of the samples between the different clusters is heterogeneous
283 (between 76 and 230 samples), with various proportions of soil and peat samples (Table 2). The
284 clusters are well-differentiated, with different means and amplitudes for MAAT, FRS and pH
285 (Fig. 2), and also based on their geographical locations (Fig. 3). Clusters B and D contain a
286 larger proportion of peat samples representative of acidic environments, which can explain the
287 lower pH values by ca. 1 to 2 units compared to the other clusters (Gorham, 1991; Killups,
288 2005; Dedysh et al., 2006; Comont et al., 2006) (Table 2, Fig. 2). Cluster A shows samples
289 mainly distributed over tropical and subtropical latitudes (Fig. 3) associated with high MAAT
290 (22.4 ± 6.0 °C) and rather high MAP (1069 ± 385 mm/yr; Table 2, Fig. 2). Cluster B samples
291 are distributed over temperate to subtropical latitudes, with precipitation amounts (1237 ± 643
292 mm/yr) as high as for samples from cluster A, but lower MAAT (16.0 ± 7.5 °C; Fig. 3; Table
293 2). Clusters A and B are characterized by comparable and higher MAATs than those from the
294 other clusters, and conversely lower FRS (Fig. 2, Table 2). Samples of cluster C are mostly
295 distributed in China and correspond to loess samples (Fig. 3). Within this cluster, MAP ($453 \pm$
296 643 mm/yr) is the lowest of all clusters and MAAT (6.7 ± 5.1 °C) is on average lower than in
297 clusters A and B, associated with a higher FRS (Table 2; Fig. 2). The samples from clusters D

298 and E show similar geographical distributions, mostly in the northern hemisphere, at temperate
299 latitudes, and even polar latitudes (Fig. 3). This results in lower MAP and MAAT especially
300 for cluster D (784 ± 457 mm/yr and 3.9 ± 5.9 °C, respectively), and a high FRS, similar to cluster
301 C (Table 2; Fig. 2). Thus, the statistical differentiation of the brGDGT dataset into different
302 clusters is reflected through various descriptive environmental parameters.

303

304

305 **3.2. BrGDGT distribution in the global dataset and associated clusters**

306

307 The fractional abundances of the individual brGDGTs were determined in the global
308 dataset and in the five clusters statistically derived from the latter (Fig. 4). In the global dataset,
309 the acyclic brGDGTs *Ia*, *Ila* and *Ila'* were predominant. Distinct brGDGT distributions were
310 observed in each cluster. In cluster A, the tetra-methylated brGDGTs *Ia* and *Ib* as well as penta-
311 methylated brGDGT *Ila* are the most abundant. Acyclic brGDGT *Ia* is largely predominant (ca.
312 75% of total brGDGT relative abundance) in cluster B. In cluster C, 6-methyl acyclic isomers
313 of the penta- and hexa-methylated brGDGTs (*Ila'* and *IIla'*) and brGDGT *Ia* represent
314 altogether ca. 65% of the total brGDGT relative abundance. The brGDGT distribution of cluster
315 D is dominated by acyclic compounds *Ia* and *Ila*. In cluster E, the 6-methyl brGDGTs are
316 slightly more abundant than the 5-methyl isomers, with acyclic brGDGTs *Ia*, *Ila*, *Ila'*, *IIla*,
317 *IIla'* and monocyclic brGDGTs *Ib*, *Ilb*, *Ilb'* representing each between ca. 10 and 20% of total
318 brGDGT relative abundance. The obvious differences in brGDGT distribution between the 5
319 clusters are also reflected in the indices derived from these compounds. Thus, the MBT'_{5Me} is
320 higher in clusters A and B (mean 0.88 ± 0.09 and 0.82 ± 0.13 , respectively) than in clusters C
321 (0.56 ± 0.14) as well as D and E (0.47 ± 0.09 and 0.49 ± 0.12 , respectively, Fig. 5a). Regarding
322 the CBT', it is much lower in clusters B and D (mean -1.29 ± 0.55 and -0.95 ± 0.63 ,
323 respectively) than in cluster C (0.40 ± 0.24), E (0.13 ± 0.25 ; Fig. 5b) and A (-0.16 ± 0.33),.
324 Similarly, the relative abundance of 6-methyl vs. 5-methyl brGDGTs (IR_{6Me} ratio; Eq. 3) is
325 much lower in clusters B and D (mean ~ 0.2) than in the other three clusters (mean comprised
326 between 0.62 and 0.80; Fig. 5c).

327 **3.3. Relationships between MBT'_{5Me} and MAAT**

328 A strong and significant correlation between MAAT and MBT'_{5Me} is observed when
329 considering the total soil dataset (Supp. Fig. 1; Eq. 4):

330

331 $MAAT (^{\circ}C) = 35.98 \times MBT'_{5Me} - 12.74$ ($n=775$; $R^2= 0.65$, $RMSE= 5.2^{\circ}C$) (4)

332

333 Nevertheless, this global calibration shows a considerable scatter. The linear
334 regressions between the MBT'_{5Me} and MAAT were further explored for each cluster derived
335 from the PCA analysis (Fig. 1). Clusters A and B show strong significant linear relationships
336 ($R^2 = 0.61$ and 0.77 , respectively; $p < 0.0001$) between MAAT and MBT'_{5Me} (Fig. 6a, b;)
337 associated with improved RMSE ($3.8^{\circ}C$ and $3.6^{\circ}C$, respectively) compared to the global
338 calibration (Supp. Fig. 1; Eq. 4). In contrast, for the other clusters (Figs. 4c, d, e), especially D
339 and C, significant ($p < 0.0001$) but weak relationships between MBT'_{5Me} and MAAT are
340 observed (Fig. 6).

341 To further investigate the influence of the proposed environmental variables (MAAT,
342 MAP, pH, FRS) on the brGDGT relative abundance in (i) the global dataset and (ii) the different
343 clusters, RDA was performed (Fig. 7). Regarding clusters A and B, the first two axes explain
344 63.9% and 74.6% of the total variance of the dataset, with an explained fitted variation of 97.3%
345 and 99.9%, respectively (i.e. explained fitted variation; relationship between the fractional
346 abundances of brGDGTs and the selected environmental variables, calculated as the sum of all
347 constrained eigenvalues) (Fig. 7a, b; Table 3). The first axis of the RDA for clusters A and B is
348 well correlated with MAAT ($r=0.80$; $r=0.90$ respectively), FRS ($r=-0.80$; $r=-0.81$), MAP
349 ($r=0.77$; $r=0.78$) and pH ($r=-0.75$; $r=-0.60$) (Fig. 5a, b; Table 3). Axis 2 of cluster A is mainly
350 correlated with the FRS ($r=0.56$) and to a lesser extent pH ($r=-0.34$) and MAAT ($r=-0.26$), while
351 axis 2 of cluster B is predominantly negatively correlated with pH ($r=-0.76$) and to a lesser
352 extent with MAAT ($r=-0.41$; Fig. 7a, b; Table 3). The quantification of the combined influence
353 of the different environmental variables on the brGDGT distribution shows a predominant
354 effect of the FRS (39.6%) and to a lesser extent pH (18.6%) for cluster A and MAAT (57.3%)
355 and to a lesser extent pH (16.7%) for cluster B (Table 3). The predominant influence of the
356 thermal regime (MAAT and FRS) on the relative distribution of brGDGTs in these two clusters,
357 despite large variation in pH range, especially in cluster B (Fig. 2), explains why the linear
358 regressions between MBT'_{5Me} and MAAT are stronger (Fig. 6) than that observed for the global
359 dataset (Supp. Fig. 1; Eq. 4).

360 The first two axes of the RDAs for clusters C and D explain 34.9% and 53.4% of the
361 total inertia of the dataset, with an explained fitted variation of 95.7% and 99.4%, respectively
362 (Fig. 7c, d; Table 3). For these two clusters, axis 1 is strongly negatively correlated with pH
363 ($r=-0.91$ and -0.96 , respectively) and positively correlated with MAP values ($r=0.59$ and 0.48 ,
364 respectively). Axis 2 is controlled by the thermal regime, being mainly correlated with FRS

365 (r=0.74 and -0.68 for clusters C and D, respectively) and MAAT (r=-0.78 and 0.39,
366 respectively). The quantification of the combined influence of the environmental variables on
367 brGDGT distribution in clusters C and D shows a predominant effect of soil pH (25.6% and
368 47.4%, respectively) and, only to a much lesser extent, MAP and MAAT (<6%) (Table 3). This
369 is consistent with the weak correlation between MBT'_{5Me} and MAAT for cluster C, and absence
370 of correlation for cluster D as well as the high scattering of the corresponding values (Fig. 6c,
371 d).

372 The first two axes of the RDA for cluster E explain 28.4% of the total inertia of the
373 dataset and the explained fitted variation is 94.40% (Fig. 7e; Table 3). Axis 1 is mainly
374 correlated with MAAT (r=-0.76), FRS (r=0.87), pH (r=0.71) and to a lesser extent, MAP (r=-
375 0.41; Table 3). Axis 2, on the other hand, is mainly influenced by MAAT (r=0.62), and to a
376 lesser extent pH (r=0.41; Table 3). When examining the combined effect of the different
377 environmental variables in this cluster, it appears that brGDGT distribution is mainly and
378 significantly controlled by FRS (15.2%) and to a lesser extent by MAAT (6.6%) and pH (7.7%;
379 Table 3). The major influence of the thermal regime (FRS, MAAT) on brGDGT distribution in
380 cluster E is consistent with the relationship ($R^2=0.44$) observed between MBT'_{5Me} and MAAT
381 (Fig. 6e). Nevertheless, in contrast with cluster A, the additional influence of pH may explain
382 the moderate determination coefficient of this correlation.

383 Regarding the global dataset, the first two axes explain 69.5% of the total variance of
384 the dataset (Fig. 7f; Table 3) and the selected environmental variables explain 99.0% of the
385 variance of the brGDGT relative abundances (Table 3). Axis 1 is strongly controlled by pH (r=-
386 0.92) and MAP (r=0.72) and to a lesser extent by FRS (r=-0.54) and MAAT (r=0.48). Axis 2 is
387 strongly correlated with MAAT (r=0.83), and to a lesser extent FRS (r=-0.58) followed by MAP
388 (r=0.31) and pH (r=0.34). The quantification of the combined influence of the environmental
389 variables on brGDGT distribution in the global dataset shows a predominant effect of soil pH
390 (51.8%) and to a lesser extent MAAT (15.7%), with only a minor influence of FRS and MAP
391 (<2%; Table 3).

392
393

394 **4. Discussion**

395 **4.1. Constraints on the MBT'_{5Me}-MAAT relationship in soils**

396 *4.1.1. Global level*

397 The MBT'_{5Me} was shown to be linearly and strongly correlated with MAAT in the
398 present extended soil dataset (Supp. Fig. 1; Eq. 2), as previously observed at the global level
399 (e.g., De Jonge et al., 2014; Dearing Crampton-Flood et al., 2020). Nevertheless, in line with
400 these previous studies, the RMSE remains high (5.2 °C). This scatter may have multiple
401 sources, such as the fact that the brGDGT calibrations are achieved against MAAT, whereas
402 brGDGT-producing bacteria live in soils. Soil temperature is not necessarily equivalent to
403 MAAT and also depends on the vegetation cover (e.g. Wang et al., 2020), which may explain
404 part of the scatter. Moreover, another source of uncertainty may be related to the fact that
405 climatic data derived from the nearest weather stations (gridded datasets) are often used to
406 develop brGDGT calibrations, while they may not appropriately reflect the local air
407 temperatures nor the soils ones, as recently reported by Pérez-Angel et al. (2020). Last, part of
408 the remaining uncertainty in the brGDGT-MAAT calibrations may be due to the influence of
409 other environmental parameters than MAAT on brGDGT distribution, as discussed below.

410 The RDA analysis performed on our global dataset showed that soil pH was the main
411 environmental control on brGDGT distribution besides MAAT (Fig. 7), as also previously
412 reported (e.g. Weijers et al., 2007; Peterse et al., 2012; De Jonge et al., 2014; Naafs et al.,
413 2017a). As expected, 6-methyl brGDGTs were all located in the left quadrant along axis 2,
414 mainly controlled by pH, in line with the positive correlations previously observed between
415 these compounds and pH (De Jonge et al., 2014; Dang et al., 2016). These isomers were
416 purposefully excluded from the calculation of the MBT so that it is no more related to pH and
417 only to MAAT (De Jonge et al., 2014). Nevertheless, the relative abundances of two of the main
418 brGDGTs involved in the MBT'_{5Me} (*Ia* and *Iia*) were shown to be significantly correlated with
419 pH ($R = 0.52$ and 0.32 , respectively; $p < 0.001$; Sup. Table 1) in the present dataset, as also
420 previously observed (De Jonge et al., 2014). Such correlations with pH were even higher than
421 those observed with MAAT ($R^2 = 0.28$ and 0.24 for compounds *Ia* and *Iia*, respectively; $p <$
422 0.001 ; Sup. Table 1), which may explain part of the remaining uncertainty in the MBT'_{5Me}
423 relationship. Very recently, De Jonge et al. (2021) highlighted the importance of taking into
424 account the effect of soil pH on MBT'_{5Me} values and associated temperature reconstructions, as
425 soil pH was shown to be the main factor responsible for concomitant changes in brGDGT
426 distribution and bacterial community composition in mid- and high-latitude experimental sites
427 and hypothesized that such conclusions were also valid at the global scale.

428 Soil moisture has also been suggested to have an effect on the relative abundance of
429 brGDGT distribution (e.g. Loomis et al., 2010; Dirghanghi et al., 2013; Menges et al., 2014;
430 Dang et al., 2016; Naafs et al., 2017a), with weak or no linear relationships between
431 MBT/MBT'_{5Me} and MAAT in arid soils (MAP < 500 mm/yr). The relative soil moisture is
432 related to pH variations, with arid soils mainly being alkaline (Naafs et al., 2017a). This may
433 have a role on the diversity of bacterial communities (Lauber et al., 2009; Shen et al., 2019).
434 Alternatively and/or complementarily, it was suggested that brGDGT producers may change
435 their membrane composition in response to soil moisture changes (Loomis et al., 2010; Dang
436 et al., 2016). The relative soil moisture may also impact the capacity of a soil to retain heat
437 (Idso et al., 1975; Davidson et al., 1998; Balleza et al., 2014; Dang et al., 2016), indirectly
438 influencing the methylation degree of brGDGTs. Dang et al (2016) and Naafs et al. (2017a)
439 especially showed that the MBT'_{5Me} was only significantly correlated with MAAT when the
440 ratio of 5- vs. 6-methyl brGDGTs (IR_{6Me} ; Eq. 3) was lower than 0.5. Therefore, to potentially
441 improve the accuracy of the global MBT'_{5Me}-MAAT calibration (Supp. Fig. 1), the total dataset
442 was divided into two subgroups based on a threshold value of 0.5 for the IR_{6Me} ratio as proposed
443 by Dang et al. (2016) and Naafs et al. (2017a). Two subgroups with similar number of samples
444 ($n=389$ for $IR_{6Me}>0.5$; $n=384$ for $IR_{6Me}<0.5$) were thus obtained. The linear regressions between
445 MBT'_{5Me} and MAAT in the two subgroups were statistically similar (Supp. Fig. 2), even though
446 a slightly higher determination coefficient and lower RMSE were observed when $IR_{6Me}<0.5$.
447 Moreover, the regressions obtained for the two subgroups did not show obvious improvements
448 with the one derived from the total dataset. Therefore, a separation of the present dataset based
449 on the IR_{6Me} values does not appear necessary, as previously observed by Dearing Crampton-
450 Flood (2020) for their dataset, showing that the relative abundance of 6- vs. 5-methyl brGDGTs
451 only has a limited influence on the MBT'_{5Me}-MAAT relationship at the global scale.

452 More recently, De Jonge et al. (2019) suggested that the response of brGDGTs to
453 temperature changes is strongly dependent on the nature of bacterial communities present in
454 soils. These authors initially showed that the distribution of brGDGTs in a set of geothermally
455 warmed soils from Iceland changed when the average annual soil temperature was above 14°C.
456 This sudden change in brGDGT distribution coincided with an abrupt shift in the bacterial
457 community composition. A relative increase in brGDGT *Ia* vs. homologues *IIa* and *IIIa* was
458 observed in the soils with annual soil temperature > 14 °C (warm soil cluster). This was
459 reflected in a change in the community index (CI) proposed by De Jonge et al. (2019; Eq. 5):

$$460 \quad CI = [Ia]/([Ia] + [IIa] + [IIIa]) \quad (5)$$

461

462 The CI index is similar to the MBT'_{5Me} (Eq. 1), except that it excludes the compounds
463 containing cyclopentyl moieties, i.e. those suspected to be pH-sensitive compounds. De Jonge
464 et al. (2019) proposed a CI threshold of 0.64 to separate the geothermal soils with an annual
465 soil temperature higher than 14°C (CI >0.64) and those with a lower temperature (CI
466 <0.64). This observation was extended to the global scale, revealing two distinct clusters over
467 the entire peat and soil dataset compiled by De Jonge et al. (2019) – one considered as a “cold”
468 subgroup ($n=251$ soils and peats; MAAT between -8.3°C and 18.2°C) and another one as a
469 “warm” subgroup ($n=195$ soils and peats; MAAT between 0.4°C and 27.1°C) – with different
470 responses to temperature and pH. The slopes and determination coefficients of the MBT'_{5Me}-
471 MAAT relationship were significantly different in the two subgroups, which may explain part
472 of the uncertainty in the MBT'_{5Me}-MAAT correlation at the global scale. The extended dataset
473 proposed in the present study ($n=775$) was thus divided into two subgroups based on the CI
474 threshold of 0.64 proposed by De Jonge et al. (2019). A strong linear relationship between
475 MBT'_{5Me} and MAAT was observed for the warm cluster ($R^2=0.71$), while it was much weaker
476 for the cold subgroup ($R^2=0.20$; Fig. 8), as previously observed by De Jonge et al. (2019) on
477 their global dataset. The discrimination of samples in clusters based on their different brGDGT
478 signature and response to environmental changes allows better understanding the MBT'_{5Me}-
479 MAAT relationship.

480 The CI threshold (0.64) defined by De Jonge et al. (2019) was based on a smaller
481 number of soils ($n=446$) than available in the present study. To refine this value based on the
482 present extended sample set ($n=775$), all linear regressions between MAAT and MBT'_{5Me} were
483 tested by successive iteration for CI values from 0 to 1 (0.001 step; Supp. Fig. 3a). The refined
484 threshold is 0.69 (corresponding to the best adjusted R^2), still in agreement with the clusters
485 presented in this study (Fig. 5d) and those of De Jonge et al. (2019). It leads to only limited
486 changes in the slopes and intercepts of the MBT'_{5Me}-MAAT relationships in the warm and cold
487 clusters (Sup. Fig.3b) in comparison with those obtained for the previously defined threshold
488 of 0.64 (Fig. 8). The CI might be considered first and foremost in any use of the MBT'_{5Me} index
489 to reconstruct paleo-MAATs, as proposed by De Jonge et al. (2019).

490

491 4.1.2. *Sample clustering*

492 Complementarily to the empirical approach of De Jonge et al. (2019) described above,
493 we used statistical tools to (i) classify the samples of the present peat and soil dataset based on
494 their brGDGT distribution and (ii) investigate and compare the influence of the environmental
495 factors on the brGDGT distribution in each cluster (Fig. 1).

496 Clusters A and B, which encompassed soils from temperate and (sub)tropical areas,
497 were characterized by similarly high MBT'_{5Me} values (mean > 0.8; Fig. 5a) and CI > 0.64 (Fig.
498 5d) and can be related to “warm” groups as defined by De Jonge et al. (2019). These two clusters
499 differed by the more acidic nature of samples from cluster B than from cluster A, reflected in
500 the much lower CBT' values in cluster B, consistent with the positive relationship usually
501 observed between CBT' and pH (De Jonge et al., 2014). In line with the increase in the
502 fractional abundance of 6-methyl brGDGTs with pH previously observed in soils (De Jonge et
503 al., 2014), the samples from cluster A were also characterized by higher IR_{6Me} ratio than those
504 from cluster B. Despite these differences related to pH, the brGDGT distributions of these
505 clusters were mainly impacted by the thermal regime (Table 3). Thus, moderate ($R^2 > 0.25$) to
506 strong correlations ($R^2 > 0.50$, $p < 0.001$) between acyclic 5-methyl brGDGTs (*Ia*, *IIa* and *IIIa*)
507 and MAAT were obtained in clusters A and B, respectively, as previously observed at the global
508 level (De Jonge et al., 2014; Naafs et al., 2017a). This was reflected in the strong positive
509 correlations observed between MBT'_{5Me} and MAAT in clusters A and B (Fig. 6a,b), with non-
510 significant differences between slopes, intercepts and RMSE, highlighting the overall similar
511 response of brGDGT source microorganisms to temperature changes in soil and peat samples
512 from the two “warm” subgroups.

513 In contrast with clusters A and B, the three other clusters were characterized by CI <
514 0.64 (Fig. 5d), corresponding to “cold” groups (De Jonge et al., 2019). They encompassed soil
515 and peat samples from cold to (sub)temperate zones (Fig. 3), with similar range of MBT'_{5Me}
516 values, which were much lower than those of clusters A and B, consistent with the increase of
517 the methylation degree of brGDGTs at lower temperatures (Weijers et al., 2007). Samples from
518 cluster D were more acidic than those from cluster C and E, leading to distinct brGDGT
519 distributions, and especially lower average CBT' values in cluster D. The differences in
520 brGDGT distribution between the three “cold” clusters were also reflected in the IR_{6Me} ratio,
521 the highest values of the latter in cluster C being consistent with the higher relative abundance
522 of 6-methyl vs. 5-methyl brGDGTs generally observed in arid/alkaline soils (e.g. De Jonge et
523 al., 2014; Naafs et al., 2017a), as those from cluster C (Table 2). In addition, the three clusters
524 largely differed in their dependence to environmental parameters. Thus, brGDGT distributions

525 in clusters C and D were predominantly influenced by pH and to a lesser extent by the thermal
526 regime (Table 3). This led to weak (or no) correlations ($R^2 < 0.1$) between individual brGDGTs
527 and MAAT (or FRS) (Supp. Tables 4, 5), explaining in turn the weak relationships between the
528 MBT'_{5Me} and MAAT in these two clusters (Fig. 6c, d). In contrast with clusters C and D, the
529 influence of the thermal regime (MAAT/FRS) on brGDGT distribution in samples from cluster
530 E was higher than the one of pH, with weak to moderate correlations (R^2 0.2-0.45) only between
531 the relative abundance of tetramethylated brGDGTs *Ib* and *Ic* and hexamethylated brGDGT
532 IIIa and MAAT (Sup. Table 6; Fig. 7e), as also observed at the global level (De Jonge et al.,
533 2014). Nevertheless, no correlations between brGDGTs *Ia/IIIa* and MAAT were observed in
534 cluster E. This contrasts with observations made at the global level in the previous (De Jonge
535 et al., 2014; Naafs et al., 2017a) or present soil datasets (Sup. Table 1), where these compounds
536 were considered as temperature-sensitive. This explains the more moderate correlation between
537 the MBT'_{5Me} and MAAT in cluster E (R^2 0.44) than in clusters A and B or the total dataset (R^2
538 > 0.6 ; Fig. 6). It should be noted that in the cold soil cluster of De Jonge et al. (2019), the MAAT
539 was similarly only correlated with the relative abundances of brGDGTs *Ia* and *IIIc*, leading to
540 a weak correlation between MBT'_{5Me} and MAAT (R^2 0.28). Nevertheless, in the global cluster
541 cold defined by De Jonge et al. (2019), brGDGT *Ia* was significantly negatively correlated with
542 pH (R^2 0.69; $p < 0.001$), in contrast with cluster E (R^2 0.16; $p = 0.0003$). Such a difference may
543 be related to the smaller size ($n = 77$) and different samples constituting cluster E vs. the global
544 cold cluster ($n = 251$) of De Jonge et al. (2019).

545 Overall, the brGDGT distribution was differently impacted by environmental variables
546 in each of the clusters of the present study (Fig. 7; Table 3), the effect of the thermal regime
547 being predominant only in the two warm clusters (A and B) and to a lesser extent cold cluster
548 E. We also observed a different dependency of the brGDGT distribution to temperature in the
549 two warm clusters vs. cold cluster E, leading to distinct correlations between MBT'_{5Me} and
550 MAAT for the two types of clusters, consistent with previous observations by de Jonge et al.
551 (2019). Despite the smaller size of the cold cluster E ($n = 77$) vs. the two warm clusters A ($n = 76$)
552 and B ($n = 174$) altogether, the relationship between MBT'_{5Me} and MAAT for the global dataset
553 ($n = 775$) was observed to be driven by the one of cluster E, as revealed by the similar slopes
554 and intercepts (Fig. 6). This implies that the MBT'_{5Me} proxy is much more influenced by the
555 temperature changes encountered in cold cluster E than in the other warm clusters. Such a
556 difference in sensitivity between warm and cold groups has also been previously reported by
557 De Jonge et al. (2019), who suggested that different bacterial communities, with different
558 brGDGT fingerprints, may be associated with the warm and cold groups. Similarly, this shift

559 in bacterial communities could at least partly explain the different MBT'_{5Me}-MAAT
560 relationships in clusters A/B vs. E. Additionally, the thermal regime may influence the activity
561 of brGDGT-producing microorganisms and associated biosynthesis, as several previous studies
562 suggested that brGDGTs may be preferentially produced during warm seasons (Weijers et al.,
563 2011b; Huguet et al., 2013; Deng et al., 2016). As the thermal regime is much higher in the cold
564 clusters (C, D and E), they should be the most impacted by the seasonal production of
565 brGDGTs, thus weakening the relationships between MAAT and brGDGT distribution. Raberg
566 et al. (2021) very recently proposed to replace MAAT by a warm season index (mean
567 temperature of months above freezing) to take into account the thermal regime effect (especially
568 large at high-latitude), thus improving brGDGT calibrations in lake sediments. Similarly,
569 Dearing Crampton-Flood et al. (2020) showed that the MBT'_{5Me} in soils and peats was better
570 correlated with the average temperature of months above 0 °C (BayMBT₀ model, R² 0.70) than
571 with MAAT (BayMBT, R² 0.64), with a reduced RMSE (3.8 vs. 6.0 °C for the BayMBT₀ and
572 BayMBT, respectively).

573 In any case, the fact that (i) a strong relationship between MBT'_{5Me} and MAAT is only
574 observed for the warmer clusters (Fig. 6), (ii) the MBT'_{5Me}-MAAT relationship at the global
575 level is driven by the moderate correlation of cold subgroup E containing only ca. 10% of the
576 samples from the total dataset and (iii) weak correlations between MBT'_{5Me} and MAAT are
577 observed for ca. 55% of the dataset (clusters C and D) highlights the limitations of using a single
578 index and a simple linear regression model to capture the response of brGDGTs to MAAT
579 changes. Thus, very recently, Pérez-Angel et al. (2020) showed that the fractional abundance
580 of brGDGT *Ia* was non-linearly related to MAAT in a dataset comprised of tropical soils
581 ($n=175$). Other models taking into account such non-linear trends should be complementarily
582 developed to better reflect this complex response to MAAT changes. Machine-learning models
583 can be used to this aim and will be tested in the following.

584

585 **4.2. A novel global calibration between MAAT and brGDGT distribution using the**
586 **random forest algorithm**

587 *4.2.1. Model development*

588 A machine-learning model – the random forest algorithm – was tested to potentially
589 derive stronger and more accurate global MAAT calibration from brGDGT distributions than
590 single linear regressions based on the MBT'_{5Me} index.

591 The random forest algorithm was first "trained" during the learning phase to estimate
592 MAAT from brGDGT relative abundances. During this phase, the model produces decision
593 trees that will automatically discriminate the compounds whose relative abundances are not
594 influenced by temperature, and thus selects those to be used to estimate MAAT, while taking
595 in account the interdependency of environmental parameters. The developed model was called
596 “random **F**orest **R**egression for Pale**O**MAAT using brGDGTs” (FROG).

597 As the 5-methyl brGDGTs are considered to be mainly correlated with MAAT in
598 contrast with 6-methyl isomers (De Jonge et al., 2014), the random forest model was initially
599 applied to the fractional abundances of the 5-methyl brGDGTs only (FROG_{5Me}). After training
600 on the sampling dataset (75% of the global dataset; $n=583$), a strong ($R^2 = 0.80$) global
601 calibration between 5-methyl brGDGT relative abundance and MAAT was obtained (Fig. 9a),
602 with a RMSE of 4.12 °C. The random forest algorithm was then separately trained and tested
603 using the relative abundances of all brGDGTs including the 6-methyl isomers (FROG). This
604 "alternative" global calibration ($R^2=0.81$ and RMSE =4.09 °C; Fig. 9b) appeared similar to the
605 FROG_{5Me} model (Fig. 9a). The RMSE of both random forest calibrations (FROG_{5Me} and
606 FROG) are lower than those of the simple linear regression between MBT'_{5Me} and MAAT (Fig.
607 9c, d; Table 4), with no clear trends in the residuals of the models. As observed by Naafs et al.
608 (2017a) and Dearing Crampton-Flood et al. (2020) in their global MBT'_{5Me}-MAAT calibration,
609 a major part of the uncertainty in the FROG model is likely due to the wide dispersion of
610 predicted temperatures for MAAT < 10°C (Fig. 9c, d), corresponding to samples from high
611 latitudes and/or high elevations. The characteristics of the FROG and FROG_{5Me} model for the
612 whole dataset and test dataset (R^2 , RMSE, variance in residuals) were provided in Table 4 and
613 are comparable, which shows that the FROG models do not overfit the data or, if so, to an
614 extremely moderate extent.

615 Even though the random forest algorithm agnostically selected the set of brGDGTs
616 which best describes the variability of MAAT, it presents the advantage of not being a blackbox,
617 as the mechanism behind it can be described. First, the independent variations in the individual

618 brGDGT variations with estimated MAAT for the FROG model can be compared, with non-
619 linear trends for all the compounds (Sup. Fig. 4), showing the interest of using a non-linear
620 model to describe the relationship between brGDGT distribution and MAAT. For example, the
621 relative abundance of tetramethylated brGDGT *Ia* was observed to non-linearly increase with
622 predicted MAAT, in contrast with compounds *Ia* and *IIIa* (Sup. Fig. 4). This is consistent with
623 the trends obtained from the linear regression models both on the present extended dataset (Sup.
624 Table 1) and on previous soil datasets (e.g. De Jonge et al., 2014; Naafs et al., 2017a), reflecting
625 the decrease in the methylation degree and thus increase in the MBT'_{5Me} with MAAT (Supp.
626 Fig. 1; Eq. 2).

627 Second, the weight of the different variables used to define the random forest model
628 with MAAT could be quantified using the permutation importance method (Fig 9e, f; Breiman,
629 2001). This method consists in considering each compound separately and determining how
630 much model performance decreases if this compound is removed from the model. The
631 brGDGTs *Ia*, *Ia*, *IIIa* were the homologues predominantly used by the FROG_{5Me} model (Fig.
632 9e), consistent with the major influence of these compounds in the MBT'_{5Me} (De Jonge et al.,
633 2014) and the existing linear relationship between this index and MAAT (Sup. Fig. 1). Other
634 brGDGT homologues, considered as being less temperature sensitive than the aforementioned
635 acyclic brGDGTs (e.g. De Jonge et al., 2014), were also taken into account by the FROG_{5Me}
636 model, especially compounds *Ib* and *Iib*, which contain one cyclopentyl moiety. These
637 compounds were similarly selected by a stepwise forward selection method to develop brGDGT
638 calibration with MAAT, e.g. for East African lacustrine sediments (Russell et al., 2018) or
639 Chinese soils (Wang et al., 2020). Therefore, different statistical models suggest that
640 homologues *Ib* and *Iib* can improve the brGDGT-MAAT relationship in various environments,
641 as it might take into account the covariations between MAAT and other environmental
642 parameters. In contrast, brGDGTs *Ic* and *Iic*, which are low-abundant in the peat and soil
643 dataset, have only a minor weight in the FROG_{5Me} model (Fig. 9e).

644 Regarding the random forest model established with all brGDGTs (FROG model), the
645 predominant homologues were also the acyclic brGDGTs *Ia*, *Ia*, and *IIIa*, and to a lesser extent,
646 compound *Ib* as well as compounds *IIIa'*, *Ia'*, i.e. 6-methyl isomers, pointing to their
647 contribution to MAAT reconstruction. 6-methyl brGDGTs were usually considered as being
648 predominantly influenced by pH variations (De Jonge et al., 2014b). Nevertheless, they were
649 also included in some recent local soil calibrations between brGDGT fractional abundances and
650 MAAT in China (Wang et al., 2020) and Mongolia (Dugerdil et al., 2021) and in a regional
651 calibration based on tropical soils ($n=175$; Pérez-Angel et al., 2020). These compounds were

652 shown to be part of those sensitive to temperature and improved the correlation with MAAT
653 (R^2 and RMSE) in comparison with linear MBT'_{5Me}-MAAT regression. It should also be noted
654 that both linear and non-linear regression models were tested for the development of the
655 aforementioned pan-tropical soil brGDGT calibration with MAAT (Pérez-Angel et al., 2020).
656 The non-linear multiple regression showed slight improvements over the linear multiple
657 regression based on brGDGT fractional abundances, which was interpreted as showing the
658 importance of taking into account non-linear influences when establishing brGDGT
659 calibrations in warm environments. In agreement with these recent results, the FROG model
660 integrates both types of brGDGT isomers, including 6-methyl ones, thus reflecting the
661 environmental reality, as the whole distribution of brGDGTs, and not only 5-methyl ones, can
662 be concomitantly and non-linearly influenced by environmental parameters including MAAT
663 (Sup. Fig. 4).

664 The respective weights of the different brGDGTs were observed to differ between the
665 FROG and FROG_{5Me} calibrations (Fig. 9). Nevertheless, the comparison of the statistical
666 characteristics of these models do not allow favoring one vs. another, as both of them showed
667 similar determination coefficients, RMSE and estimation ranges (Table 4). Therefore, both
668 models will be tested on modern soil samples and sedimentary archives and compared in section
669 4.4.

670 *4.2.2. Comparison of the FROG model with previous global soil calibrations*

671 The FROG/FROG_{5Me} model were compared with previously published global soil/peat
672 linear calibrations based on Ordinary Least Square (De Jonge et al., 2014), Deming (Naafs et
673 al., 2017a) and Bayesian (Dearing Crampton-Flood et al., 2020) regressions. The FROG and
674 FROG_{5Me} models are characterized by a higher determination coefficient and lower RMSE (Fig.
675 9; $R^2=0.81$ vs. 0.8 and RMSE=4.09 °C vs. 4.12 °C, respectively) than available calibrations (De
676 Jonge et al., 2014: $R^2=0.61$, RMSE=4.8°C; $n=231$; Naafs et al., 2017a: $R^2=0.71$, RMSE=
677 4.1°C, $n=177$; Dearing Crampton-Flood et al., 2020: $R^2=0.64$, RMSE=6.0 °C, $n=343$; Table 4),
678 while being obtained from a larger dataset ($n=775$).

679 In order to make a direct statistical comparison of the different models, the calibrations
680 by De Jonge et al. (2014), Naafs et al. (2017a) and Dearing Crampton-Flood et al. (2020;
681 BayMBT model) were applied to the same dataset as the one used for the random forest model
682 (Fig. 10; Table 4). This confirms that the FROG model performs better in terms of robustness
683 and accuracy than the previously published global soil and peat calibrations (De Jonge et al.,
684 2014; Naafs et al., 2017a; Dearing-Crampton Flood et al., 2020). The FROG model and the
685 BayMBT both increases the upper limit of MAAT estimation by > 3 °C and 5 °C, respectively,

686 in comparison with the calibrations by De Jonge et al. (2014) and Naafs et al. (2017a; Table 4).
687 In contrast, the lower limit of MAAT estimates for the random forest model is higher than that
688 of the other linear calibrations by De Jonge et al. (2014), Naafs et al. (2017a) and especially
689 Dearing Crampton-Flood et al. (2020). This difference is related to the fact that the extended
690 peat and soil dataset contains a lower number of samples collected under cold climates (12%
691 with MAAT < 0°C). The low representation of such type of samples is a limitation for the
692 training phase of the random forest model. In contrast, by definition, the linear models are able
693 to reach higher or lower limits, even without the presence of « extreme » samples in the dataset.
694 Nevertheless, machine-learning algorithms are flexible and the lower limit of temperature
695 prediction of the FROG model could be decreased by analyzing brGDGTs in a larger number
696 of soil/peat samples from cold settings and adding them to the model.

697 The higher robustness and accuracy of the FROG model compared to those based on
698 the MBT_{5Me} (Table 4) could be explained by its non-parametric nature and the fact that it takes
699 into account non-linear influences on the brGDGT distribution (Supp. Fig. 4), unlike the linear
700 models (De Jonge et al., 2014; Naafs et al., 2017a; Dearing Crampton-Flood et al., 2020). In
701 addition, one of the prerequisites of the Bayesian model, as recently proposed by Dearing
702 Crampton-Flood et al. (2020), is to determine a *prior*, before analysis. Nevertheless, as shown
703 in this study, the influence of environmental parameters on brGDGT distribution in soils/peat
704 is sample-dependent (e.g. Fig. 7; Table 3). Therefore, it seems difficult to determine a *prior*
705 adapted to all the samples in such a large dataset. Bayesian models could be more efficient
706 when applied to local/regional calibrations, where the *prior* can be determined more precisely.
707 In any case, this study shows that efficient statistical approaches are useful to improve brGDGT
708 calibrations.

709

710 **4.3. Development of alternative MAAT calibrations**

711 Although the FROG MAAT calibration presents a strong determination coefficient,
712 the associated RMSE is still ca. 4 °C (Table 4). To test the influence of the thermal regime and
713 MAP on the FROG model and potentially improve its accuracy, alternative submodels based
714 on different subsets of the extended dataset were developed.

715

716 *4.3.1. Alternative model considering the influence of thermal regime*

717 The FRS was shown to have a significant impact on the brGDGT distribution of the
718 extended soil dataset based on RDA analyses, especially on the 5-methyl isomers (Fig. 7; Table

719 3). The FRS can be considered as an indirect indicator of the thermal regime which may
720 influence the growth of brGDGT source microorganisms and the production of these lipids.
721 Nevertheless, the impact of the thermal regime on brGDGT distribution and concentration was
722 the object of contrasting observations. The latter was not observed to be affected by seasonal
723 temperature variations in mid-latitude soils (Weijers et al., 2011b; Lei et al., 2016). In contrast,
724 several studies suggested that brGDGTs may be preferentially produced during summer in mid-
725 to high latitude peats (Weijers et al., 2011a; Huguet et al., 2013) or soils (e.g. Deng et al., 2016),
726 reflecting an effect of the seasonality on the brGDGT production, with an enhancement of the
727 latter during the seasonal optimum, for example at the warm season when the soil is not frozen.
728 An increase in the microbial biomass production was reported in unfrozen soils (Schimel and
729 Clein, 1996; Nedwell, 1999; Schimel et al., 2007).

730 To take into account the potential influence of the thermal regime on brGDGT
731 distribution, an alternative Bayesian calibration was proposed by Dearing Crampton-Flood et
732 al. (2020). This model (BayMBT₀) estimates the average temperature of all months that have
733 an average temperature above 0 °C using the MBT'_{5Me} index. The BayMBT₀ model (R²= 0.70;
734 RMSE = 3.8°C; lower limit = 0.9°C, upper limit = 27.1°C; Dearing Crampton-Flood et al.,
735 2020) improves the strength and accuracy of the BayMBT model but does not allow
736 reconstructing negative temperatures. To evaluate the influence of the thermal regime on the
737 FROG model, the same approach as that proposed by Dearing Crampton-Flood et al. (2020)
738 was applied.

739 In the present extended dataset, monthly temperatures are available for 661 out of the
740 775 soil samples. This excludes some of the samples, collected in the French Alps, Peruvian
741 Andes, Mts Pollino, Shegyla and Italy (Véquaud et al., 2020). The alternative calibration based
742 on the average temperature of all months that have a temperature above 0 °C (FROG₀) is as
743 strong as the global random forest calibration (FROG; Fig. 11), with an R²= 0.84, and is more
744 accurate (RMSE = 2.5 °C) than the latter (Table 4, Supp. Fig. 6), even though it is based on a
745 slightly reduced dataset (*n*=661). The FROG₀ calibration also performs better than the
746 BayMBT₀ model (R² = 0.56; RMSE = 4.1°C), with a higher R², lower RMSE but with a slightly
747 lower range of MAAT estimation (Table 4). Overall, the present study, through RDA analyses
748 (Fig. 7) and the FROG₀ model, highlights the effect of the thermal regime on brGDGT
749 distribution. In particular, it can be assumed that the activity of brGDGT source organisms will
750 be reduced within frozen soils. In addition, snowfall can isolate the soil from the atmospheric
751 compartment, decoupling atmospheric temperatures from soil temperatures during the cold
752 season. This confirms the interest of using alternative models including this effect to improve

753 the accuracy of MAAT reconstruction, as also suggested by Dearing Crampton-Flood et al.
754 (2020).

755

756 *4.3.2. Alternative model considering the influence of MAP*

757 As previously discussed, the brGDGT distribution may be largely impacted by soil
758 moisture (e.g., Dirghangi et al., 2013; Naafs et al., 2017a). In the present study, MAP was used
759 as a proxy of the soil water content (SWC), which varies as a first approximation according to
760 the precipitation regime, although other factors can play a role, such as relief (topography),
761 evapotranspiration, grain size or vegetation cover (Crave and Gascuel-Odoux, 1997; Gómez-
762 Plaza et al., 2001).

763 The influence of MAP was tested with the FROG₅₀₀ random forest model. This
764 alternative to the FROG model excludes all the samples with MAP < 500 mm/year,
765 corresponding to soils previously defined as dry, and generally described as alkaline soils from
766 arid regions, poor in organic matter (Peterse et al., 2012; Naafs et al., 2017a; Dearing Crampton-
767 Flood et al., 2020). The FROG₅₀₀ model contains 442 samples and only shows a slight
768 improvement in MAAT reconstruction in comparison with the FROG calibration, with a
769 slightly higher determination coefficient ($R^2 = 0.85$) and lower RMSE (3.5°C; Table 4). The
770 slight decrease in the RMSE of the FROG₅₀₀ vs. the FROG model (by ca. 0.5 °C) is likely
771 related to the large reduction of the dataset size by more than 300 samples. Therefore, the use
772 of the FROG₅₀₀ model, as the BayMBT₅₀₀ proposed by Dearing Crampton-Flood et al. (2020),
773 does not seem preferable, as its performance is only slightly better than the original FROG
774 calibration, even though it contains a much lower number of samples.

775

776 **4.4. Paleo application of the FROG global calibration**

777 The statistical characteristics including determination coefficients and RMSE are not
778 sufficient enough to discriminate between calibrations for paleotemperature reconstructions, as
779 recently noticed by Dugerdil et al. (2021). Therefore, the performance and validity of the FROG
780 models were tested and compared with the temperature record from Pliocene sediments from
781 the North Sea basin (Dearing Crampton-Flood et al., 2018, 2020) and from a Chinese loess-
782 paleosol sequence covering the last 110 kyr (Gao et al., 2012; Lu et al., 2016; Wang et al.,
783 2020). These archives were the object of previous paleostudies, providing a context for the
784 interpretation of the MAAT data from the FROG models.

785

786 *4.4.1. Application of the FROG models to sediments from the Pliocene*

787 The performance and reliability of the FROG models were tested using a paleorecord
788 based on Pliocene sediments from the North Sea basin (Dearing Crampton-Flood et al., 2018,
789 2020). BrGDGTs were previously analyzed in this archive, with subsequent reconstruction of
790 past temperature variations during the Pliocene using the MBT'_{5Me} (Dearing Crampton-Flood
791 et al., 2018). More recently, the validity of the BayMBT₀ model proposed by Dearing
792 Crampton-Flood et al. (2020) was tested on this archive and compared with the MAAT records
793 derived from the soil brGDGT calibrations by De Jonge et al. (2014) and Naafs et al. (2017a).

794 The aforementioned calibrations, as well as the FROG/FROG_{5Me} models (Fig.11a),
795 showed an overall decrease in MAAT during the Pliocene, consistent with the global cooling
796 that is documented in the literature (Lisiecki and Raymo, 2005; Dearing Crampton-Flood et al.,
797 2018, 2020). The MAAT records derived from the different models showed similar qualitative
798 trends over the reconstructed period, even though the range of variation appeared slightly
799 smaller (3-4°C) for the FROG model than for the other calibrations (Fig. 11a). The FROG and
800 FROG_{5Me} calibrations provided similar results, implying that the inclusion of the 6-methyl
801 isomers in the models did not significantly change the paleotemperature reconstructions in the
802 present case.

803 The absolute temperatures reconstructed by the different brGDGT models (Fig. 11a)
804 were generally lower than those expected in NW Europe (13-14 °C) over the Pliocene and
805 derived from pollen assemblages and model outputs (Dearing Crampton-Flood et al., 2020 and
806 references therein). The BayMBT record showed the lowest MAAT estimates, with very large
807 oscillations and negative temperatures during the late Pliocene, in contrast with the other global
808 calibrations (Fig. 11a). Nevertheless, it is unlikely that temperatures went below 0°C during the
809 Pliocene based on the presence of pollen assemblage of warm-adapted species found in Dutch
810 sediments of this period, as specified by Dearing Crampton-Flood et al. (2020). That is why
811 Dearing Crampton-Flood et al. (2020) favored the use of the BayMBT₀ model, based on the
812 mean temperature of all months above 0 °C, for MAAT reconstruction, as it is more accurate
813 (RMSE = 3.8°C) than the BayMBT model (RMSE = 6°C). The BayMBT₀ model indeed
814 provided higher absolute MAAT estimates than the BayMBT and much more consistent with
815 temperatures estimated from other proxies (i.e. 13-14°C; Fig. 11b).

816 Although the FROG model did not result in such negative reconstructed MAAT, the
817 FROG₀ model was also applied to the Pliocene archive to compare with the BayMBT₀ record
818 (Fig. 11b). Both models showed a decreasing trend in temperature over the Pliocene, although
819 the BayMBT₀ model displayed larger oscillations (between 3.5 °C and 16.2°C; mean 10.7°C)

820 than the FROG₀ one (between 10 °C and 13.1°C; mean 11.7°C). As the MAAT estimates from
821 the FROG₀ model are associated with a smaller error (2.5 °C; Table 4) than the other global
822 brGDGT calibrations including the BayMBT₀ one (3.8 °C; Dearing Crampton-Flood et al.,
823 2020), an improvement in the accuracy of the paleoreconstruction over this period can be
824 anticipated when using the FROG₀ model. This model, unlike the FROG model, estimates the
825 mean temperature of all months above 0°C, which can be considered as more reflecting warm
826 season temperature. Unlike the FROG model, the FROG₀ model can neglect the thermal and
827 nival influences on the samples, and so on the source organisms of brGDGTs. We suggest that
828 the FROG₀ model should be used in addition to the FROG model rather than alone, in order to
829 obtain complementary information on annual and seasonal temperatures dynamics in
830 paleoclimate studies.

831

832 *4.4.2. Application of the FROG model to a Chinese loess-paleosol sequence*

833 The FROG calibration was also applied to the Lantian loess-paleosol sequence (LPS),
834 located in the southern Chinese Loess Plateau (Fig. 1 in Gao et al., 2012) and covering the last
835 110 kyr (Gao et al., 2012). BrGDGT data from this 8.5 m sequence and associated
836 paleotemperature reconstructions were previously published (Gao et al., 2012; Lu et al., 2016;
837 Wang et al., 2020). Recently, Wang et al. (2020) analysed brGDGTs in 149 modern soils
838 covering a large climate gradient in China and calibrated brGDGT distribution against both
839 mean annual soil temperature (MAST) and MAAT. They applied these local MAST and MAAT
840 brGDGT calibrations as well as the global MAAT calibration (MAT_{mr}) by De Jonge et al.
841 (2014) to the LPS sequence over the last 60 kyr and showed that the local MAST calibration
842 provided more reasonable variations in the past continental temperatures than the local (Wang
843 et al., 2020) and global (Peterse et al., 2012; De Jonge et al., 2014) MAAT calibrations. This
844 was notably related to past changes in vegetation coverage which may affect the relationship
845 between MAST and MAAT.

846 Regarding the modern period, Wang et al. (2020) collected six surface soil samples
847 adjacent to the Lantian LPS to serve as a reference for the present time. The reconstructed
848 MAAT based on the FROG (11.9 ± 0.8 °C) and FROG_{5Me} (11.9 ± 0.9 °C) for these soils are in
849 agreement with the recorded MAAT at this site (12.6 °C; Wang et al., 2020), suggesting that
850 these calibrations can at least be applied for modern MAAT reconstruction in the region. At
851 this site, MAAT₀ can be considered as close to MAAT, as a winter (December; January;
852 February) mean air temperature of ~1 °C was reported (Gao et al., 2012). The FROG₀ model
853 provided temperature estimates (14.3 ± 0.3 °C) consistent with the expected one. In contrast

854 with the different FROG calibrations, the temperature estimates derived from the BayMBT₀
855 model (16.8 ± 2.0 °C) were much higher than expected, with a large variability between the 6
856 soil samples. As for the local MAST calibration by Wang et al. (2020), it provided temperatures
857 (12.8 ± 1.4 °C) consistent with those measured in soils nearby as previously reported (Wang et
858 al., 2020) (Fig. 12a).

859 So as to assess the reliability of the calibration, the local MAST calibration by Wang
860 et al. (2020) was applied to the whole LPS covering the last 110 kyr and compared with the
861 MAAT records derived from the different FROG models (FROG, FROG_{5Me} and FROG₀) as
862 well as the BayMBT₀ (Dearing-Crampton Flood et al., 2020). The BayMBT₀ model was chosen
863 rather than the BayMBT one, as it provides most accurate temperature reconstructions, as
864 reported by Dearing Crampton-Flood et al. (2020) and specified above. Nevertheless, it should
865 be noted that the BayMBT₀ model, as the FROG₀ ones, allow the reconstruction of the mean
866 temperature of all months > 0 °C (MAAT₀) instead of the MAAT.

867 When applied to the Lantian LPS, the different calibrations showed the same
868 qualitative trends (Fig. 12b). Thus, the temperature oscillated between 110 and 60 kyr (with
869 different amplitudes depending on the calibration) and then showed a continuous cooling trend
870 between 60 kyr and 30 kyr, the lowest values being reached between ca. 22 and 27 kyr BP,
871 corresponding to the local Last Glacial Maximum (LGM), as previously observed by Gao et al.
872 (2012) and Lu et al. (2016). Then, the temperature increased rapidly and peaked at the Early
873 Holocene before decreasing toward the present-day values. Such general trends are similar to
874 those previously reported for the Lantian LPS (Gao et al., 2012; Lu et al., 2016; Wang et al.,
875 2020). Local insolation (35 °N) was previously shown to be the dominant factor impacting the
876 temperature records of this site, as similarly observed in other LPSs of the southern plateau
877 (Peterse et al., 2011, 2014; Jia et al., 2013). Thus, the temperature maxima at ca. 12-15 kyr, 63
878 kyr and 82 kyr BP were consistent with insolation maxima, while the low temperature observed
879 at ca. 90 kyr, 75 kyr, 25 kyr and the late Holocene coincided with insolation minima (Wang et
880 al., 2008, Fig. 1; Gao et al., 2012, Fig. 3; Lu et al., 2016). The brGDGT-derived temperature
881 records and the $\delta^{18}\text{O}$ record from Chinese speleothems (Wang et al., 2008), related to monsoon
882 intensity, displayed roughly similar patterns (Fig. 12b), showing the relationship existing
883 between the temperature and precipitation intensity, even though the temperature record was
884 generally observed to precede the $\delta^{18}\text{O}$ record, as already noticed in the Lantian LPS (Gao et
885 al., 2012; Wang et al., 2020) and other sites of the southern plateau (Peterse et al., 2011, 2014;
886 Jia et al., 2013). The FROG, FROG_{5Me} and FROG₀ calibrations showed especially well-defined
887 peaks over the period between 75 kyr and 110 kyr BP, corresponding to the glacial and

888 interglacial substages of Marine Isotope Stage 5. These extrema, also apparent when using the
889 local MAST calibration by Wang et al. (2020) and, with a more reduced amplitude, the
890 BayMBT₀ model, were similarly observed in the Chinese speleothem $\delta^{18}\text{O}$ record (Wang et al.,
891 2008; Fig. 12b), although the exact timing differed probably in relation to age modelling
892 uncertainties (Wang et al., 2008).

893 The results derived from the brGDGT calibrations were in general agreement with the
894 previously published records and climate models. Thus, the temperature estimates for the LGM
895 derived from the FROG/FROG_{5Me}/FROG₀ and local MAST (Wang et al., 2020) calibrations
896 were, respectively, ca. 2 °C and 4 °C lower than those of the present-day surficial sediments
897 (Fig. 12a), consistent with the difference of ca. 2-4 °C derived from East Asian climate models
898 (Ju et al., 2007). As for the BayMBT₀ calibration, the difference in temperature estimates
899 between LGM and present-day (ca. 3 °C, respectively) should be interpreted with care, as the
900 temperature derived from the top-core sediment (ca. 17 °C) is abnormally high compared to the
901 recorded MAAT (ca. 12 °C), as discussed above for soils surrounding the Lantian LPS, hence
902 BayMBT₀ seems to be less relevant than FROG models.

903 After the local LGM (at ca. 21-24 kyr BP), the temperature was shown to increase by
904 ca. 8 to 11°C (depending on the calibration) between the LGM and the peak at the Early
905 Holocene (Fig. 12b), consistent with the increase of ca. 10 °C observed during this deglaciation
906 based on lacustrine records from central eastern Europe (Sanchi et al., 2014) and western North
907 America (Feakins et al., 2019). The FROG model showed a higher warming trend (by ca. 3 °C)
908 than the FROG_{5Me} one. Such a difference is likely related to the inclusion of the 6-methyl
909 isomers in the FROG model. Loess deposits are developed under arid conditions, which favors
910 the domination of 6-methyl vs. 5-methyl brGDGTs (De Jonge et al., 2014; Naafs et al., 2017a),
911 as it is the case in the Lantian LPS (average IR_{6Me} over the whole sequence: 0.58 ± 0.07). The
912 high abundance of 6-methyl brGDGTs in arid/alkaline soils (i.e. IR_{6Me} > 0.5) was shown to
913 make complex the applicability of the MBT'_{5Me} (Naafs et al., 2017a). In the Lantian LPS, this
914 led to the development of specific local calibrations, such as the one proposed by Wang et al.
915 (2020) to reconstruct MAST, based on stepwise regression method and including several 6-
916 methyl brGDGTs (*Ila'*, *Ilb'*, *IIla'*) to improve temperature reconstruction. Similarly, it may not
917 be excluded that the presence of 6-methyl brGDGTs in the FROG model could help in
918 improving paleotemperature reconstructions in comparison with the FROG_{5Me} calibration
919 containing only 5-methyl brGDGTs. Nevertheless, except in the early Holocene, the difference
920 in temperature estimates between the FROG and FROG_{5Me} models was generally <1.5 °C,
921 which is much lower than the RMSE of these calibrations (ca. 4 °C; Table 4). As both of them

922 provided similar qualitative trends, this does not favor one calibration vs. another for
923 paleotemperature reconstruction of this archive.

924 To conclude, the FROG, FROG_{5Me} and FROG₀ calibrations were able to accurately
925 reconstruct present-day temperatures at the Lantian LPS, in contrast with the BayMBT₀ model.
926 The reliability of the FROG model was further demonstrated when applied to the whole
927 sedimentary record. It showed documented climatic variations, with a reduced error (4°C)
928 compared to previous global soil calibrations (Peterse et al., 2012; De Jonge et al., 2014), and
929 consistent with the trends derived from a local MAST calibration (Wang et al., 2020).

930

931 5. Conclusions

932 Several global brGDGT calibrations for MAAT reconstruction in soils and peats have
933 been proposed over the last years. Nevertheless, the uncertainty in brGDGT-based temperature
934 estimates is still substantial, largely due to the influence of the various environmental variables
935 in addition to MAAT on brGDGT distribution. A statistical clustering and analysis of the
936 globally distributed brGDGT dataset allowed hierarchizing the parameters affecting brGDGT
937 distribution in soils and thus the MBT'_{5Me}-MAAT relationship at the global scale. pH was
938 shown to be the main environmental control on brGDGT distribution, followed by MAAT, over
939 the whole dataset. The five statistical clusters were well-differentiated based on environmental
940 parameters (MAAT, FRS, MAP, pH) and geographical locations and were characterized by
941 distinct brGDGT distributions. A strong relationship between MBT'_{5Me} and MAAT was only
942 observed for the warmer clusters while the MBT'_{5Me}-MAAT relationship at the global level
943 was shown to be driven by the moderate correlation corresponding to a cold subgroup
944 containing only ca. 10% of the samples from the total dataset. This highlighted the limitations
945 of using a single index and a simple linear regression model to capture the response of brGDGTs
946 to temperature changes.

947 A new improved MAAT calibration based on random forest algorithm was then
948 proposed, the so-called Random Forest Regression for PaleOMAAT using brGDGTs (**FROG**).
949 The FROG model, which is multi-factorial and non-parametric, appears to be more robust and
950 accurate than previous global calibrations while being proposed on an extended soil and peat
951 dataset. This is related to the fact that it takes into account the non-linear influences on the
952 relationships between MAAT and the relative abundances of individual brGDGTs. Finally, the
953 FROG model was applied to two existing paleorecords and compared with available
954 calibrations, showing its suitability for paleoreconstructions. Application of this new model

955 should improve the accuracy of brGDGT-based MAAT reconstructions in soils and peats,
956 especially in environments where the MBT_{5Me} shows some limitations because of potential
957 confounding factors. As the random forest algorithm is adaptative and flexible, the FROG
958 model, freely available to the community through an R package, could be easily further
959 improved by the implementation of additional samples in the dataset. The machine-learning
960 approach proposed in this study for calibrating the brGDGT-MAAT relationship could be
961 applied to other settings, such as lacustrine ones.

962

963 **Research data.** FROG models presented in this study are freely available using a R package
964 with a web-application on a GITHUB repository ([paleoFROG](#)). The soil dataset used in this
965 study will be added on Pangaea.

966

967 **Acknowledgments.** We thank Sorbonne Université for a PhD scholarship to Pierre Véquaud
968 and the Labex MATISSE (Sorbonne Université) for financial support. The EC2CO programme
969 (CNRS/INSU – BIOHEFECT/MICROBIEN) is thanked for funding of the SHAPE project.
970 Arnaud Huguet and Sergio Contreras are grateful for funding of the ECOS SUD/ECOS ANID
971 #C19U01/190011 project. Andrew T. Nottingham was supported by the UK Natural
972 Environment Research Council (NERC), grant NE/T012226.

973

974 **References**

975 Baker A., Blith A.J., Jex C.N., Mcdonald J.A., Woltering M., Khan S.J., 2019. Glycerol dialkyl
976 glycerol tetraethers (GDGT) distributions from soil to cave: Refining the speleothem
977 paleothermometer. *Organic Geochemistry* **136**, 103890.

978

979 Balleza D., Garcia-Arribas A. B., Sot J., Ruiz-Mirazo K. and Goñi F. M. (2014) Ether- versus
980 Ester-Linked Phospholipid Bilayers Containing either Linear or Branched Apolar
981 Chains. *Biophysical Journal* **107**, 1364–1374.

982 Blaga C. I., Reichart G.-J., Heiri O. and Sinninghe Damsté J. S. (2009) Tetraether membrane
983 lipid distributions in water-column particulate matter and sediments: a study of 47
984 European lakes along a north–south transect. *J Paleolimnol* **41**, 523–540.

985 Braak C. J. F. ter and Smilauer P. (2002) *CANOCO Reference Manual and CanoDraw for*
986 *Windows User's Guide: Software for Canonical Community Ordination (version 4.5).*,
987 www.canoco.com, Ithaca NY, USA.

988 Breiman L. (2001) Random Forests. *Machine Learning* **45**, 5–32.

- 989 Choler P. (2018) Winter soil temperature dependence of alpine plant distribution: Implications
 990 for anticipating vegetation changes under a warming climate. *Perspectives in Plant*
 991 *Ecology, Evolution and Systematics* **30**, 6–15.
- 992 Coffinet S., Huguet A., Bergonzini L., Pedentchouk N., Williamson D., Anquetil C., Gałka M.,
 993 Kołaczek P., Karpińska-Kołaczek M., Majule A., Laggoun-Défarge F., Wagner T. and
 994 Derenne S. (2018) Impact of climate change on the ecology of the Kyambangunguru
 995 crater marsh in southwestern Tanzania during the Late Holocene. *Quaternary Science*
 996 *Reviews* **196**, 100–117.
- 997 Comont L., Laggoun-Défarge F. and Disnar J.-R. (2006) Evolution of organic matter indicators
 998 in response to major environmental changes: The case of a formerly cut-over peat bog
 999 (Le Russey, Jura Mountains, France). *Organic Geochemistry* **37**, 1736–1751.
- 1000 Crave A. and Gascuel- Odoux C. (1997) The Influence of Topography on Time and Space
 1001 Distribution of Soil Surface Water Content. *Hydrological Processes* **11**, 203–210.
- 1002 Dang XinYue, Xue J., Yang H. and Xie S. (2016) Environmental impacts on the distribution of
 1003 microbial tetraether lipids in Chinese lakes with contrasting pH: Implications for
 1004 lacustrine paleoenvironmental reconstructions. *Sci. China Earth Sci.* **59**, 939–950.
- 1005 Dang Xinyue, Yang H., Naafs B. D. A., Pancost R. D. and Xie S. (2016) Evidence of moisture
 1006 control on the methylation of branched glycerol dialkyl glycerol tetraethers in semi-arid
 1007 and arid soils. *Geochimica et Cosmochimica Acta* **189**, 24–36.
- 1008 Davidson E. A., Belk E. and Boone R. D. (1998) Soil water content and temperature as
 1009 independent or confounded factors controlling soil respiration in a temperate mixed
 1010 hardwood forest. *Global Change Biology* **4**, 217–227.
- 1011 Davtian N., Ménot G., Bard E., Poulénard J. and Podwojewski P. (2016) Consideration of soil
 1012 types for the calibration of molecular proxies for soil pH and temperature using global
 1013 soil datasets and Vietnamese soil profiles. *Organic Geochemistry* **101**, 140–153.
- 1014 De Jonge C., Hopmans E. C., Stadnitskaia A., Rijpstra W. I. C., Hofland R., Tegelaar E. and
 1015 Sinninghe Damsté J. S. (2013) Identification of novel penta- and hexamethylated
 1016 branched glycerol dialkyl glycerol tetraethers in peat using HPLC–MS2, GC–MS and
 1017 GC–SMB-MS. *Organic Geochemistry* **54**, 78–82.
- 1018 De Jonge C., Hopmans E. C., Zell C. I., Kim J.-H., Schouten S. and Sinninghe Damsté J. S.
 1019 (2014) Occurrence and abundance of 6-methyl branched glycerol dialkyl glycerol
 1020 tetraethers in soils: Implications for palaeoclimate reconstruction. *Geochimica et*
 1021 *Cosmochimica Acta* **141**, 97–112.
- 1022 De Jonge C., Radujković D., Sigurdsson B. D., Weedon J. T., Janssens I. and Peterse F. (2019)
 1023 Lipid biomarker temperature proxy responds to abrupt shift in the bacterial community
 1024 composition in geothermally heated soils. *Organic Geochemistry* **137**, 103897.
- 1025 De Jonge C., Kuramae E.E., Radujković D., Weedon J.T., Janssens I.A., Peterse F. (2021). The
 1026 influence of soil chemistry on branched tetraether lipids in mid- and high latitude soils:
 1027 Implications for brGDGT-based paleothermometry. *Geochimica et Cosmochimica Acta*
 1028 **310**, 95-112.
 1029

- 1030 Dearing Crampton-Flood E., Peterse F., Munsterman D. and Sinninghe Damsté J. S. (2018)
 1031 Using tetraether lipids archived in North Sea Basin sediments to extract North Western
 1032 European Pliocene continental air temperatures. *Earth and Planetary Science Letters*
 1033 **490**, 193–205.
- 1034 Dearing Crampton-Flood E., Tierney J. E., Peterse F., Kirkels F. M. S. A. and Sinninghe
 1035 Damsté J. S. (2020) BayMBT: A Bayesian calibration model for branched glycerol
 1036 dialkyl glycerol tetraethers in soils and peats. *Geochimica et Cosmochimica Acta* **268**,
 1037 142–159.
- 1038 Dedysh S. N., Pankratov T. A., Belova S. E., Kulichevskaya I. S. and Liesack W. (2006)
 1039 Phylogenetic Analysis and In Situ Identification of Bacteria Community Composition
 1040 in an Acidic Sphagnum Peat Bog. *Appl. Environ. Microbiol.* **72**, 2110–2117.
- 1041 Denisko D. and Hoffman M. M. (2018) Classification and interaction in random forests. *Proc*
 1042 *Natl Acad Sci USA* **115**, 1690–1692.
- 1043 Ding S., Schwab V. F., Ueberschaar N., Roth V.-N., Lange M., Xu Y., Gleixner G. and Pohnert
 1044 G. (2016) Identification of novel 7-methyl and cyclopentanyl branched glycerol dialkyl
 1045 glycerol tetraethers in lake sediments. *Organic Geochemistry* **102**, 52–58.
- 1046 Ding S., Xu Y., Wang Y., He Y., Hou J., Chen L. and He J.-S. (2015) Distribution of branched
 1047 glycerol dialkyl glycerol tetraethers in surface soils of the Qinghai–Tibetan Plateau:
 1048 implications of brGDGTs-based proxies in cold and dry regions. *Biogeosciences* **12**,
 1049 3141–3151.
- 1050 Dirghangi S. S., Pagani M., Hren M. T. and Tipple B. J. (2013) Distribution of glycerol dialkyl
 1051 glycerol tetraethers in soils from two environmental transects in the USA. *Organic*
 1052 *Geochemistry* **59**, 49–60.
- 1053 Dugerdil L., Joanin S., Peyron O., Jouffroy-Bapicot I., Vanni re B., Boldgib B., Unkelbach J.,
 1054 Behling H., M not G. (2021). Climate reconstructions based on GDGT and pollen
 1055 surface datasets from Mongolia and Baikal area: calibrations and applicability to
 1056 extremely cold–dry environments over the Late Holocene. *Climate of the Past* **17**, 1199–
 1057 1226.
- 1058
 1059 Dunkley Jones T., Eley Y. L., Thomson W., Greene S. E., Mandel I., Edgar K. and Bendle J.
 1060 A. (2020) OPTiMAL: a new machine learning approach for GDGT-based
 1061 palaeothermometry. *Climate of the Past* **16**, 2599–2617.
- 1062 Fawcett P. J., Werne J. P., Anderson R. S., Heikoop J. M., Brown E. T., Berke M. A., Smith S.
 1063 J., Goff F., Donohoo-Hurley L., Cisneros-Dozal L. M., Schouten S., Sinninghe Damst 
 1064 J. S., Huang Y., Toney J., Fessenden J., WoldeGabriel G., Atudorei V., Geissman J. W.
 1065 and Allen C. D. (2011) Extended megadroughts in the southwestern United States
 1066 during Pleistocene interglacials. *Nature* **470**, 518–521.
- 1067 Feakins S.J., Wu M.S., Ponton C., Tierney J.E. (2019). Biomarkers reveal abrupt switches in
 1068 hydroclimate during the last glacial in southern California. *Earth and Planetary Science*
 1069 *Letters* **515**, 164-172.

- 1070 Gao L., Nie J., Clemens S., Liu W., Sun J., Zech R., Huang Y. (2012). The importance of solar
1071 insolation on the temperature variations for the past 110 kyr. *Palaeogeography,*
1072 *Palaeoclimatology, Palaeoecology* **317-318**, 128-133.
- 1073
- 1074 Gómez-Plaza A., Martínez-Mena M., Albaladejo J. and Castillo V. M. (2001) Factors
1075 regulating spatial distribution of soil water content in small semiarid catchments.
1076 *Journal of Hydrology* **253**, 211–226.
- 1077 Gorham E. (1991) Northern Peatlands: Role in the Carbon Cycle and Probable Responses to
1078 Climatic Warming. *Ecological Applications* **1**, 182–195.
- 1079 Harning D.J., Curtin L., Geirsdóttir Á., D'Andrea W.J., Miller G.H., Sepúlveda J. (2020). Lipid
1080 Biomarkers Quantify Holocene Summer Temperature and Ice Cap Sensitivity in
1081 Icelandic Lakes. *Geophysical Research Letters* **47**, e2019GL085728.
- 1082
- 1083 Harris I., Jones P. D., Osborn T. J. and Lister D. H. (2014) Updated high-resolution grids of
1084 monthly climatic observations - the CRU TS3.10 Dataset: UPDATED HIGH-
1085 RESOLUTION GRIDS OF MONTHLY CLIMATIC OBSERVATIONS. *Int. J.*
1086 *Climatol.* **34**, 623–642.
- 1087 Hofmann K., Lamprecht A., Pauli H. and Illmer P. (2016) Distribution of Prokaryotic
1088 Abundance and Microbial Nutrient Cycling Across a High-Alpine Altitudinal Gradient
1089 in the Austrian Central Alps is Affected by Vegetation, Temperature, and Soil Nutrients.
1090 *Microb Ecol* **72**, 704–716.
- 1091 Hopmans E. C., Schouten S. and Sinninghe Damsté J. S. (2016) The effect of improved
1092 chromatography on GDGT-based palaeoproxies. *Organic Geochemistry* **93**, 1–6.
- 1093 Huguet A., Coffinet S., Roussel A., Gayraud F., Anquetil C., Bergonzini L., Bonanomi G.,
1094 Williamson D., Majule A. and Derenne S. (2019) Evaluation of 3-hydroxy fatty acids
1095 as a pH and temperature proxy in soils from temperate and tropical altitudinal gradients.
1096 *Organic Geochemistry* **129**, 1–13.
- 1097 Huguet A., Fosse C., Laggoun-Défarge F., Delarue F. and Derenne S. (2013) Effects of a short-
1098 term experimental microclimate warming on the abundance and distribution of branched
1099 GDGTs in a French peatland. *Geochimica et Cosmochimica Acta* **105**, 294–315.
- 1100 Huguet A., Fosse C., Laggoun-Défarge F., Toussaint M.-L. and Derenne S. (2010) Occurrence
1101 and distribution of glycerol dialkyl glycerol tetraethers in a French peat bog. *Organic*
1102 *Geochemistry* **41**, 559–572.
- 1103 Huguet A., Francez A.-J., Jusselme M. D., Fosse C. and Derenne S. (2014) A climatic chamber
1104 experiment to test the short term effect of increasing temperature on branched GDGT
1105 distribution in Sphagnum peat. *Organic Geochemistry* **73**, 109–112.
- 1106 Idso S. B., Schmugge T. J., Jackson R. D. and Reginato R. J. (1975) The utility of surface
1107 temperature measurements for the remote sensing of surface soil water status. *Journal*
1108 *of Geophysical Research (1896-1977)* **80**, 3044–3049.

- 1109 Jia G.D., Rao Z.G., Zhang J., Li Z.Y., Chen F.H. (2013). Tetraether biomarker records from a
1110 loess-paleosol sequence in the western Chinese Loess Plateau. *Frontiers in*
1111 *Microbiology* **4** <https://doi.org/10.3389/fmicb.2013.0019>.
1112
- 1113 Ju L., Wang H., Jiang D. (2007). Simulation of the Last Glacial Maximum climate over East
1114 Asia with a regional climate model nested in a general circulation model.
1115 *Palaeogeography, Palaeoclimatology, Palaeoecology* **248**, 376-390.
1116
- 1117 Killops S. (2005) Introduction to Organic Geochemistry, 2nd edn (paperback) ed. V. Killops.
1118 *Geofluids* **5**, 236–237.
- 1119 Lauber C. L., Hamady M., Knight R. and Fierer N. (2009) Pyrosequencing-Based Assessment
1120 of Soil pH as a Predictor of Soil Bacterial Community Structure at the Continental Scale.
1121 *Appl. Environ. Microbiol.* **75**, 5111–5120.
- 1122 Lei Y., Yang H., Dang X., Zhao S. and Xie S. (2016) Absence of a significant bias towards
1123 summer temperature in branched tetraether-based paleothermometer at two soil sites
1124 with contrasting temperature seasonality. *Organic Geochemistry* **94**, 83–94.
- 1125 Liang J., Russell J. M., Xie H., Lupien R. L., Si G., Wang J., Hou J. and Zhang G. (2019)
1126 Vegetation effects on temperature calibrations of branched glycerol dialkyl glycerol
1127 tetraether (brGDGTs) in soils. *Organic Geochemistry* **127**, 1–11.
- 1128 Lisiecki L. E. and Raymo M. E. (2005) A Pliocene-Pleistocene stack of 57 globally distributed
1129 benthic $\delta^{18}\text{O}$ records. *Paleoceanography* **20**.
- 1130 Loomis S.E., Russell J.M., Sinninghe Damsté J.S. (2010). Distributions of branched GDGTs in
1131 soils and lake sediments from western Uganda: Implications for a lacustrine
1132 paleothermometer. *Organic Geochemistry* **42**, 739-751.
1133
- 1134 Loomis S. E., Russell J. M., Ladd B., Street-Perrott F. A. and Sinninghe Damsté J. S. (2012)
1135 Calibration and application of the branched GDGT temperature proxy on East African
1136 lake sediments. *Earth and Planetary Science Letters* **357–358**, 277–288.
- 1137 Lu H., Liu W., Wang H. and Wang Z. (2016) Variation in 6-methyl branched glycerol dialkyl
1138 glycerol tetraethers in Lantian loess–paleosol sequence and effect on paleotemperature
1139 reconstruction. *Organic Geochemistry* **100**, 10–17.
- 1140 McMaster G. and Wilhelm W. (1997) Growing degree-days: one equation, two interpretations.
1141 *Publications from USDA-ARS / UNL Faculty*.
- 1142 Menges J., Huguet C., Alcañiz J. M., Fietz S., Sachse D. and Rosell-Melé A. (2014) Influence
1143 of water availability in the distributions of branched glycerol dialkyl glycerol tetraether
1144 in soils of the Iberian Peninsula. *Biogeosciences* **11**, 2571–2581.
- 1145 Naafs B. D. A., Gallego-Sala A. V., Inglis G. N. and Pancost R. D. (2017a) Refining the global
1146 branched glycerol dialkyl glycerol tetraether (brGDGT) soil temperature calibration.
1147 *Organic Geochemistry* **106**, 48–56.
- 1148 Naafs B. D. A., Inglis G. N., Zheng Y., Amesbury M. J., Biester H., Bindler R., Blewett J.,
1149 Burrows M. A., del Castillo Torres D., Chambers F. M., Cohen A. D., Evershed R. P.,

- 1150 Feakins S. J., Gałka M., Gallego-Sala A., Gandois L., Gray D. M., Hatcher P. G.,
 1151 Honorio Coronado E. N., Hughes P. D. M., Huguet A., Könönen M., Laggoun-Défarge
 1152 F., Lähteenoja O., Lamentowicz M., Marchant R., McClymont E., Pontevedra-Pombal
 1153 X., Ponton C., Pourmand A., Rizzuti A. M., Rochefort L., Schellekens J., De
 1154 Vleeschouwer F. and Pancost R. D. (2017b) Introducing global peat-specific
 1155 temperature and pH calibrations based on brGDGT bacterial lipids. *Geochimica et*
 1156 *Cosmochimica Acta* **208**, 285–301.
- 1157 Nedwell D. B. (1999) Effect of low temperature on microbial growth: lowered affinity for
 1158 substrates limits growth at low temperature. *FEMS Microbiol Ecol* **30**, 101–111.
- 1159 Pearson E.J., Juggins S., Talbot H.M., Weckström J., Rosén P., Ryves D.B., Roberts S.J.,
 1160 Schmidt R. (2011). A lacustrine GDGT-temperature calibration from the Scandinavian
 1161 Arctic to Antarctic: Renewed potential for the application of GDGT-paleothermometry
 1162 in lakes. *Geochimica et Cosmochimica Acta* **75**, 6225-6238.
 1163
- 1164 Pérez-Angel L.C., Sepúlveda J., Molnar P., Montes C., Rajagopalan B., Snell K., Gonzalez-
 1165 Arango C., Dildar N. (2020). Soil and air temperature calibrations using branched
 1166 GDGTs for the Tropical Andes of Colombia: Toward a Pan-tropical calibration.
 1167 *Geochemistry, Geophysics, Geosystems* **21**, e2020GC008941
 1168
- 1169 Peterse F., Kim J.-H., Schouten S., Kristensen D. K., Koç N. and Sinninghe Damsté J. S. (2009)
 1170 Constraints on the application of the MBT/CBT palaeothermometer at high latitude
 1171 environments (Svalbard, Norway). *Organic Geochemistry* **40**, 692–699.
- 1172 Peterse F., Prins M.A., Beets C.J., Troelstra S.R., Zheng H., Gu Z., Schouten S., Sinninghe
 1173 Damsté J.S. (2011). Decoupled warming and monsoon precipitation in East Asia over
 1174 the last deglaciation. *Earth and Planetary Science Letters* **301**, 256-264.
 1175
- 1176 Peterse F., van der Meer J., Schouten S., Weijers J. W. H., Fierer N., Jackson R. B., Kim J.-H.
 1177 and Sinninghe Damsté J. S. (2012) Revised calibration of the MBT–CBT
 1178 paleotemperature proxy based on branched tetraether membrane lipids in surface soils.
 1179 *Geochimica et Cosmochimica Acta* **96**, 215–229.
- 1180 Peterse F., Martínez-García A., Zhou B., Beets C.J., Prins M.A., Zheng H., Eglinton T.I. (2014).
 1181 Molecular records of continental air temperature and monsoon precipitation variability
 1182 in East Asia spanning the past 130,000 years. *Quaternary Science Reviews* **83**, 76-82.
 1183
- 1184 Peterse F., Moy C. M. and Eglinton T. I. (2015) A laboratory experiment on the behaviour of
 1185 soil-derived core and intact polar GDGTs in aquatic environments. *Biogeosciences* **12**,
 1186 933–943.
- 1187 Powers L., Werne J. P., Vanderwoude A. J., Sinninghe Damsté J. S., Hopmans E. C. and
 1188 Schouten S. (2010) Applicability and calibration of the TEX86 paleothermometer in
 1189 lakes. *Organic Geochemistry* **41**, 404–413.
- 1190 Raberg J.H., Harning D.J., Crump S.E., de Wet G., Blumm A., Kopf S., Geirsdóttir A., Miller
 1191 G.H., Sepúlveda J. (2021). Revised fractional abundances and warm-season
 1192 temperatures substantially improve brGDGT calibrations in lake sediments.
 1193 *Biogeosciences* **18**, 3579-3603.

- 1194
- 1195 Russell J.M., Hopmans E.C., Loomis S.E., Liang J., Sinninghe Damsté J.S. (2018).
1196 Distributions of 5- and 6-methyl branched glycerol dialkyl glycerol tetraethers
1197 (brGDGTs) in East African lake sediments: Effects of temperature, pH, and new
1198 lacustrine paleotemperature calibrations. *Organic Geochemistry* **17**, 56-69.
1199
- 1200 Sanchi L., Ménot G., Bard E. (2014). Insights into continental temperatures in the northwestern
1201 Black Sea area during the Last Glacial period using branched tetraether lipids.
1202 *Quaternary Science Reviews* **84**, 98-108.
- 1203 Schimel J., Balsler T. C. and Wallenstein M. (2007) Microbial Stress-Response Physiology and
1204 Its Implications for Ecosystem Function. *Ecology* **88**, 1386–1394.
- 1205 Schimel J. P. and Clein J. S. (1996) Microbial response to freeze-thaw cycles in tundra and
1206 taiga soils. *Soil Biology and Biochemistry* **28**, 1061–1066.
- 1207 Shen C., Shi Y., Fan K., He J.-S., Adams J. M., Ge Y. and Chu H. (2019) Soil pH dominates
1208 elevational diversity pattern for bacteria in high elevation alkaline soils on the Tibetan
1209 Plateau. *FEMS Microbiol Ecol* **95**.
- 1210 Siles J. A. and Margesin R. (2016) Abundance and Diversity of Bacterial, Archaeal, and Fungal
1211 Communities Along an Altitudinal Gradient in Alpine Forest Soils: What Are the
1212 Driving Factors? *Microb Ecol* **72**, 207–220.
- 1213 Sinninghé Damsté J. S., Ossebaar J., Abbas B., Schouten S. and Verschuren D. (2009) Fluxes
1214 and distribution of tetraether lipids in an equatorial African lake: Constraints on the
1215 application of the TEX86 palaeothermometer and BIT index in lacustrine settings.
1216 *Geochimica et Cosmochimica Acta* **73**, 4232–4249.
- 1217 Sinninghé Damsté J. S., Rijpstra W. I. C., Foesel B. U., Huber K. J., Overmann J., Nakagawa
1218 S., Kim J. J., Dunfield P. F., Dedysh S. N. and Villanueva L. (2018) An overview of the
1219 occurrence of ether- and ester-linked iso-diabolic acid membrane lipids in microbial
1220 cultures of the Acidobacteria: Implications for brGDGT paleoproxies for temperature
1221 and pH. *Org. Geochem.* **124**, 63–76.
- 1222 Sinninghé Damsté J. S., Rijpstra W. I. C., Hopmans E. C., Foesel B. U., Wüst P. K., Overmann
1223 J., Tank M., Bryant D. A., Dunfield P. F., Houghton K. and Stott M. B. (2014) Ether-
1224 and Ester-Bound iso-Diabolic Acid and Other Lipids in Members of Acidobacteria
1225 Subdivision 4. *Appl. Environ. Microbiol.* **80**, 5207–5218.
- 1226 Sinninghé Damsté J. S., Rijpstra W. I. C., Hopmans E. C., Weijers J. W. H., Foesel B. U.,
1227 Overmann J. and Dedysh S. N. (2011) 13,16-Dimethyl Octacosanedioic Acid (iso-
1228 Diabolic Acid), a Common Membrane-Spanning Lipid of Acidobacteria Subdivisions 1
1229 and 3. *Appl. Environ. Microbiol.* **77**, 4147–4154.
- 1230 Tierney J. E. and Russell J. M. (2009) Distributions of branched GDGTs in a tropical lake
1231 system: Implications for lacustrine application of the MBT/CBT paleoproxy. *Organic*
1232 *Geochemistry* **40**, 1032–1036.
- 1233 Véquaud P., Derenne S., Anquetil C., Collin S., Poulénard J., Sabatier P. and Huguet A. (2021)
1234 Influence of environmental parameters on the distribution of bacterial lipids in soils

- 1235 from the French Alps: Implications for paleo-reconstructions. *Organic Geochemistry*
1236 **153**, 104194.
- 1237 Véquaud P., Derenne S., Thibault A., Anquetil C., Bonanomi G., Collin S., Contreras S.,
1238 Nottingham A., Sabatier P., Salinas N., Scott W. P., Werne J. P. and Huguet A. (2020)
1239 Development of global temperature and pH calibrations based on bacterial 3-hydroxy
1240 fatty acids in soils. *Biogeosciences Discussions*, 1–40.
- 1241 Wang Y.J., Cheng H., Edwards R.L., Kong X.G., Shao X.H., Chen S.T., Wu J.Y., Jiang X.Y.,
1242 Wang X.F., An Z.S. (2008). Millennial- and orbital-scale changes in the East Asian
1243 monsoon over the past 224,000 years. *Nature* **451**, 1090-1093.
1244
- 1245 Wang H., An Z., Lu H., Zhao Z. and Liu W. (2020) Calibrating bacterial tetraether distributions
1246 towards in situ soil temperature and application to a loess-paleosol sequence.
1247 *Quaternary Science Reviews* **231**, 106172.
- 1248 Wang H., Liu W. and Lu H. (2016) Appraisal of branched glycerol dialkyl glycerol tetraether-
1249 based indices for North China. *Organic Geochemistry* **98**, 118–130.
- 1250 Weber Y., De Jonge C., Rijpstra W. I. C., Hopmans E. C., Stadnitskaia A., Schubert C. J.,
1251 Lehmann M. F., Sinninghe Damsté J. S. and Niemann H. (2015) Identification and
1252 carbon isotope composition of a novel branched GDGT isomer in lake sediments:
1253 Evidence for lacustrine branched GDGT production. *Geochimica et Cosmochimica Acta*
1254 **154**, 118–129.
- 1255 Weijers J. W. H., Schouten S., van den Donker J. C., Hopmans E. C. and Sinninghe Damsté J.
1256 S. (2007) Environmental controls on bacterial tetraether membrane lipid distribution in
1257 soils. *Geochimica et Cosmochimica Acta* **71**, 703–713.
- 1258 Weijers J. W. H., Steinmann P., Hopmans E. C., Schouten S. and Sinninghe Damsté J. S.
1259 (2011a) Bacterial tetraether membrane lipids in peat and coal: Testing the MBT–CBT
1260 temperature proxy for climate reconstruction. *Organic Geochemistry* **42**, 477–486.
1261
- 1262 Weijers J.W., Bernhardt B., Peterse F., Werne J.P., Dungait J.A., Schouten S., Sinninghe
1263 Damsté J.S. (2011b). Absence of seasonal patterns in MBT-CBT indices in
1264 mid-latitude soils. *Geochimica et Cosmochimica Acta* **75**, 3179-3190.
1265
- 1266 Xiao W., Xu Y., Ding S., Wang Y., Zhang X., Yang H., Wang G. and Hou J. (2015) Global
1267 calibration of a novel, branched GDGT-based soil pH proxy. *Organic Geochemistry* **89–**
1268 **90**, 56–60.
- 1269 Yang H., Lü X., Ding W., Lei Y., Dang X. and Xie S. (2015) The 6-methyl branched tetraethers
1270 significantly affect the performance of the methylation index (MBT') in soils from an
1271 altitudinal transect at Mount Shennongjia. *Organic Geochemistry* **82**, 42–53.
- 1272
- 1273
- 1274

Figure and table captions

Figure 1. Principal Component Analysis performed on the global dataset ($n=767$) with the k-means clustering. (a) PC2 vs PC1, (b) PC3 vs PC1, (c) Variance explained (%) for each component, (d) Optimal number of clusters according to the elbow method, based on the observation of the Within Cluster Sum of Squares (WCSS).

Figure 2. Boxplot showing the distribution of the 4 environmental variables considered: (a) pH, (b) MAAT ($^{\circ}\text{C}$), (c) MAP (mm/year), (d) FRS. Interquartile range (IQR) = $Q_3 - Q_1$ where Q_3 is the 75th percentile and Q_1 is the 25th percentile. Outliers are defined with a 1.5 coefficient on the IQR. Letters on the panel show the differences between each cluster according to Kruskal-Wallis and Dunn post-hoc tests.

Figure 3. Spatial distribution of samples in the global dataset. The colors correspond to the different clusters.

Figure 4. Fractional abundances of the individual brGDGTs determined in the (a) Cluster A, (b) Cluster B, (c) Cluster C, (d) Cluster D, (e) Cluster E and (f) in the global dataset.

Figure 5. Boxplot showing the distribution of the (a) $\text{MBT}'_{5\text{Me}}$, (b) CBT' , (c) $\text{IR}_{6\text{Me}}$, (d) Community Index (CI; Eq. 5; defined by De Jonge et al. (2019) for the 5 clusters defined after PCA analysis (Fig. 2). The CI thresholds of 0.64 and 0.69 separating “warm” and “cold” groups as proposed by De Jonge et al. (2019) and in the present study, respectively, are represented in panel (d). Letters on the panels show the differences between each cluster according to Kruskal-Wallis and Dunn post-hoc tests.

Figure 6. Linear regressions between $\text{MBT}'_{5\text{Me}}$ and MAAT ($^{\circ}\text{C}$) for (a) cluster A, (b) cluster B, (c) cluster C, (d) cluster D, (e) cluster E, (f) the whole dataset. Grey dots correspond to soils and green dots to peats samples.

Figure 7. Redundancy analysis between brGDGT distribution and environmental variables for (a) cluster A, (b) cluster B, (c) cluster C, (d) cluster D and (e) cluster E.

Figure 8. Comparison of the linear regressions between $\text{MBT}'_{5\text{Me}}$ and MAAT ($^{\circ}\text{C}$) in the two subgroups derived from the extended dataset: “cold” cluster ($\text{CI} < 0.64$) and “warm” cluster ($\text{CI} > 0.64$). The community index (CI; Eq. 3) was defined by De Jonge et al. (2019). The dashed line corresponds to the linear relationship between MAAT ($^{\circ}\text{C}$) and $\text{MBT}'_{5\text{Me}}$ in the global dataset used in this study ($n=775$).

Figure 9. MAAT predicted by the random forest models using the relative abundances of (a) 5-methyl brGDGTs (FROG_{5Me}) and (b) all brGDGTs (FROG model). Residuals of the (c) FROG_{5Me} model and (d) FROG model with all brGDGTs plotted against predicted MAAT. Importance of the individual brGDGTs in (e) the FROG_{5Me} and (f) FROG model with all brGDGTs, according to the permutation importance method (Breiman, 2001). These results were obtained from the test dataset. Grey dots correspond to soils and green dots to peats.

Figure 10. Comparison of the global $\text{MBT}'_{5\text{Me}}$ calibration proposed by (a) De Jonge et al. (2014), (b) Naafs et al. (2017a) and (c) Dearing Crampton-Flood et al. (2020; BayMBT model) using the same dataset as the FROG model ($n = 192$). Grey dots correspond to soils and green dots to peats.

Figure 11. Reconstructed MAAT for the Pliocene marine sediment sequence from the Hank core located in the Netherlands (Dearing Crampton-Flood et al., 2018, 2020) (a) derived from the calibrations by De Jonge et al. (2014), Naafs et al. (2017a), Dearing Crampton-Flood et al. (2020; BayMBT) and from the FROG/FROG_{5Me} models and (b) derived from the calibrations by Dearing Crampton-Flood et al. (2020; BayMBT₀) and from the FROG₀ model. Grey zones are the 95% intervals for the FROG and FROG₀ models

Figure 12. (a) Estimated temperatures from the 6 surface soils collected adjacent to the Lantian LPS, serving as a reference for the present time (Wang et al., 2020), derived from the BayMBT₀ model (Dearing Crampton-Flood et al. 2020) and FROG/FROG_{5Me}/FROG₀ models (this study). (b) Comparison of the MAAT estimates from the Lantian LPS sequence covering the last 110 kyr, derived from local MAST calibration (Wang et al., 2020) the BayMBT₀ model (Dearing Crampton-Flood et al. 2020) and FROG/FROG_{5Me}/FROG₀ models (this study), with a $\delta^{18}\text{O}$ record from a Chinese speleothem (in green; Wang et al., 2008).

Table 1. Location, number, and references for soils and peat samples used to establish the new global brGDGT calibration proposed in this study. Available parameters for the different sampling sites are shown: MAAT (Mean Annual Air Temperature (°C), pH, MAP (Mean Annual Precipitation (mm/yr), FRS (Number of frost days per year). Asterisks represent the new samples added to the global dataset.

Table 2. Quantitative description of the 5 clusters obtained after k-means on PCA on the brGDGT distribution of the total dataset.

Table 3. RDA correlation coefficients of the selected environmental variables along axes 1 and 2 for each cluster, and quantification of the influence of the different environmental variables on brGDGT relative abundances. Statistical significance ($p < 0.05$) is shown with an asterisk.

Table 4. Characteristics of the different brGDGT models compared in this study to estimate MAAT in terrestrial settings: R^2 , RMSE (or RMSEP; i.e. Root Mean Square Error of Prediction, for the results on the test dataset), variance of the residuals and the upper and lower limits of estimation. The "training" samples (75%) were used to develop the different machine learning models, which were then tested on the remaining sample set. Characteristics are presented for the test dataset and all the dataset (under brackets, in italics). Previous calibrations from De Jonge et al. (2014), Naafs et al. (2017a) and Dearing Crampton-Flood et al. (2020) are indicated with asterisks.

Supplementary figure 1 Linear relationship between MBT'_{5Me} and MAAT (°C) in the global dataset used in this study ($n=775$).

Supplementary figure 2 Comparison of the linear regressions between MBT'_{5Me} and MAAT (°C) in the global dataset for samples with $\text{IR}_{6\text{Me}} > 0.5$ and $\text{IR}_{6\text{Me}} < 0.5$.

Supplementary figure 3 (a) Estimation of the community index thresholds, for each linear regression between MAAT and MBT'_{5Me}; (b) Comparison of the linear regressions between MBT'_{5Me} and MAAT (°C) in the two subgroups derived from the extended dataset: "cold" cluster ($\text{CI} < 0.69$) and "warm" cluster ($\text{CI} > 0.69$). The dashed line corresponds to the linear relationship between MAAT (°C) and MBT'_{5Me} in the global dataset used in this study ($n=775$)

Supplementary figure 4. Partial plots of the individual brGDGT variations in the FROG model proposed to estimate MAAT.

Supplementary figure 5. Partial plots of the individual brGDGT variations in the FROG_{5Me} model proposed to estimate MAAT.

Supplementary figure 6. (a) MAAT predicted by the FROG₀ model using the relative abundances of all brGDGTs. (b) Residuals of the FROG₀ model. (c) Importance of the individual brGDGTs in the FROG₀ model, according to the permutation importance method (Breiman, 2001). These results were obtained from the test dataset.

Supplementary Table 1. (a) Correlation matrix between the fractional abundances of the brGDGTs, and the MAAT, MAP, pH, FRS, CBT', MBT'_{5Me}, the community index (CI) and the IR_{6Me} in the global dataset presented in this study (n=775). The values given are the R². (b) *p*-values of the correlations shown in the supplementary table 1a.

Supplementary Table 2. (a) Correlation matrix between the fractional abundances of the brGDGTs, and the MAAT, MAP, pH, FRS, CBT', MBT'_{5Me}, the community index (CI) and the IR_{6Me} in the Cluster A. The values given are the R². (b) *p*-values of the correlations shown in supplementary table 2a.

Supplementary Table 3. (a) Correlation matrix between the fractional abundances of the brGDGTs, and the MAAT, MAP, pH, FRS, CBT', MBT'_{5Me}, the community index (CI) and the IR_{6Me} in the Cluster B. The values given are the R². (b) *p*-values of the correlations shown in supplementary table 3a.

Supplementary Table 4. (a) Correlation matrix between the fractional abundances of the brGDGTs, and the MAAT, MAP, pH, FRS, CBT', MBT'_{5Me}, the community index (CI) and the IR_{6Me} in the Cluster C. The values given are the R². (b) *p*-values of the correlations shown in the supplementary table 4a

Supplementary Table 5. (a) Correlation matrix between the fractional abundances of the brGDGTs, and the MAAT, MAP, pH, FRS, CBT', MBT'_{5Me}, the community index (CI) and the IR_{6Me} in the Cluster D. The values given are the R². (b) *p*-values of the correlations shown in the supplementary table 5a.

Supplementary Table 6. (a) Correlation matrix between the fractional abundances of the brGDGTs, and the MAAT, MAP, pH, FRS, CBT', MBT'_{5Me}, the community index (CI) and the IR_{6Me} in the Cluster E. The values given are the R². (b) *p*-values of the correlations shown in the supplementary table 6a.

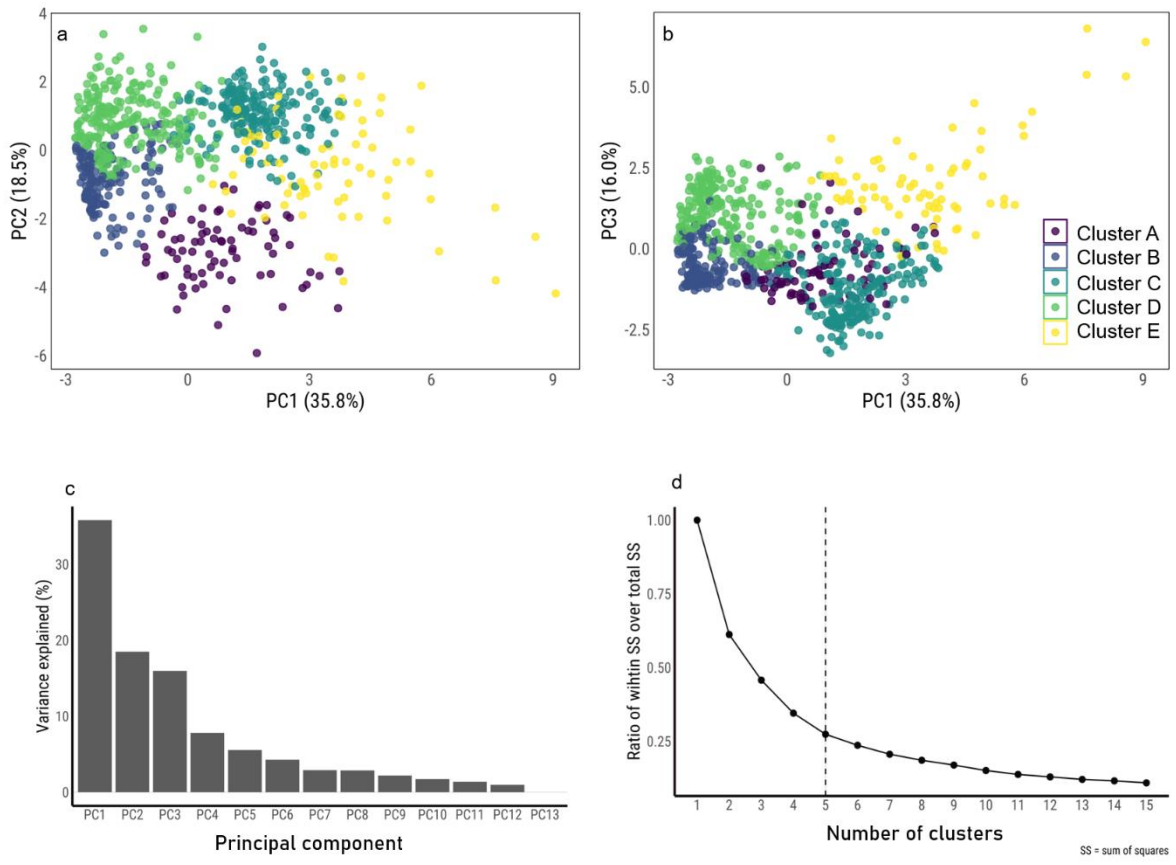


Figure 1

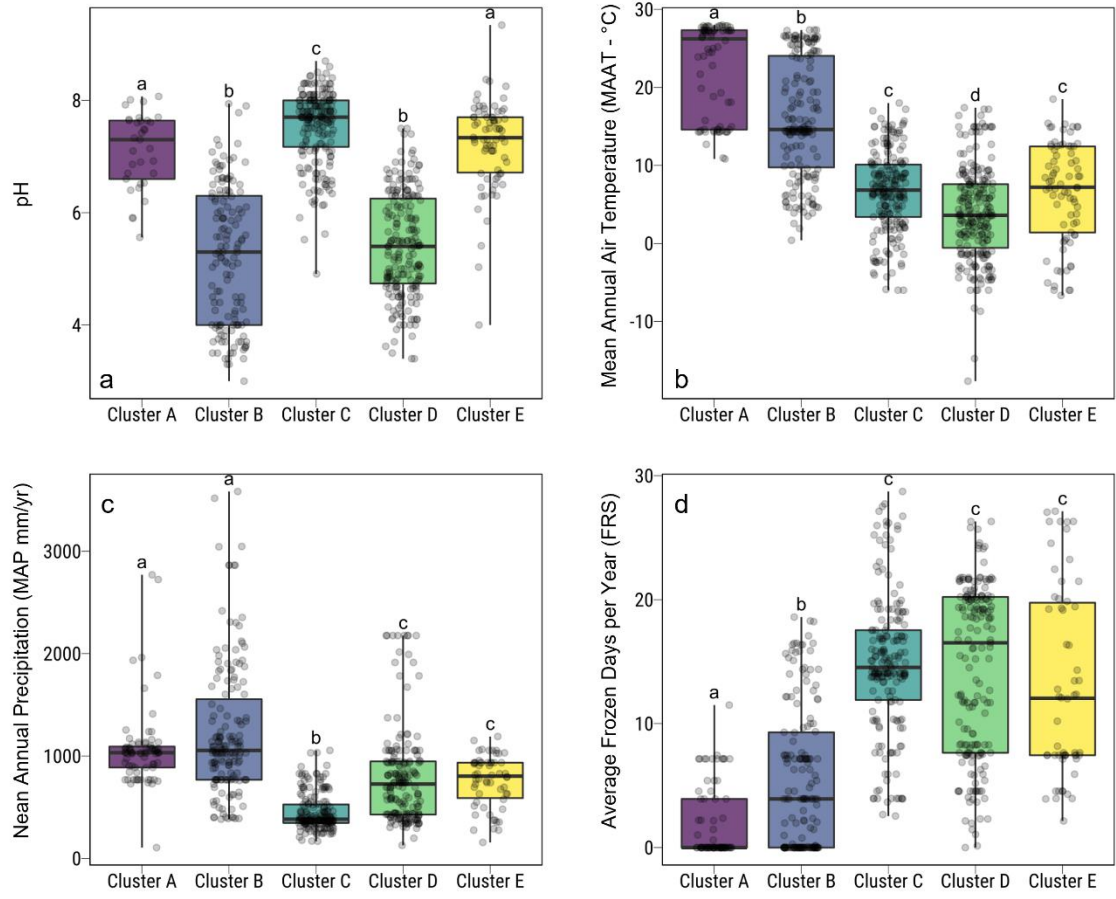


Figure 2

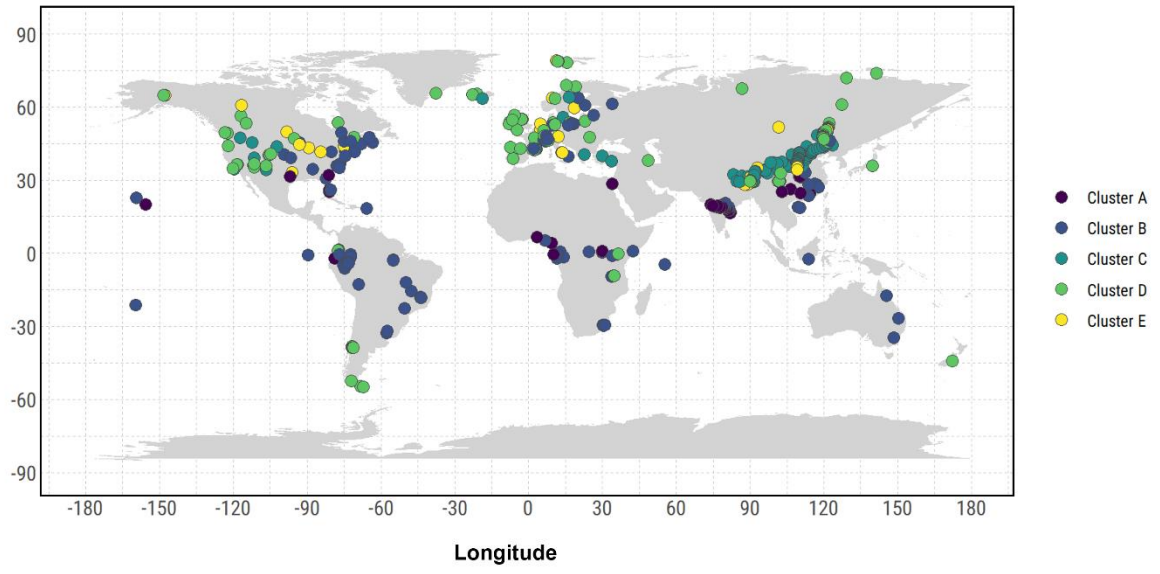


Figure 3

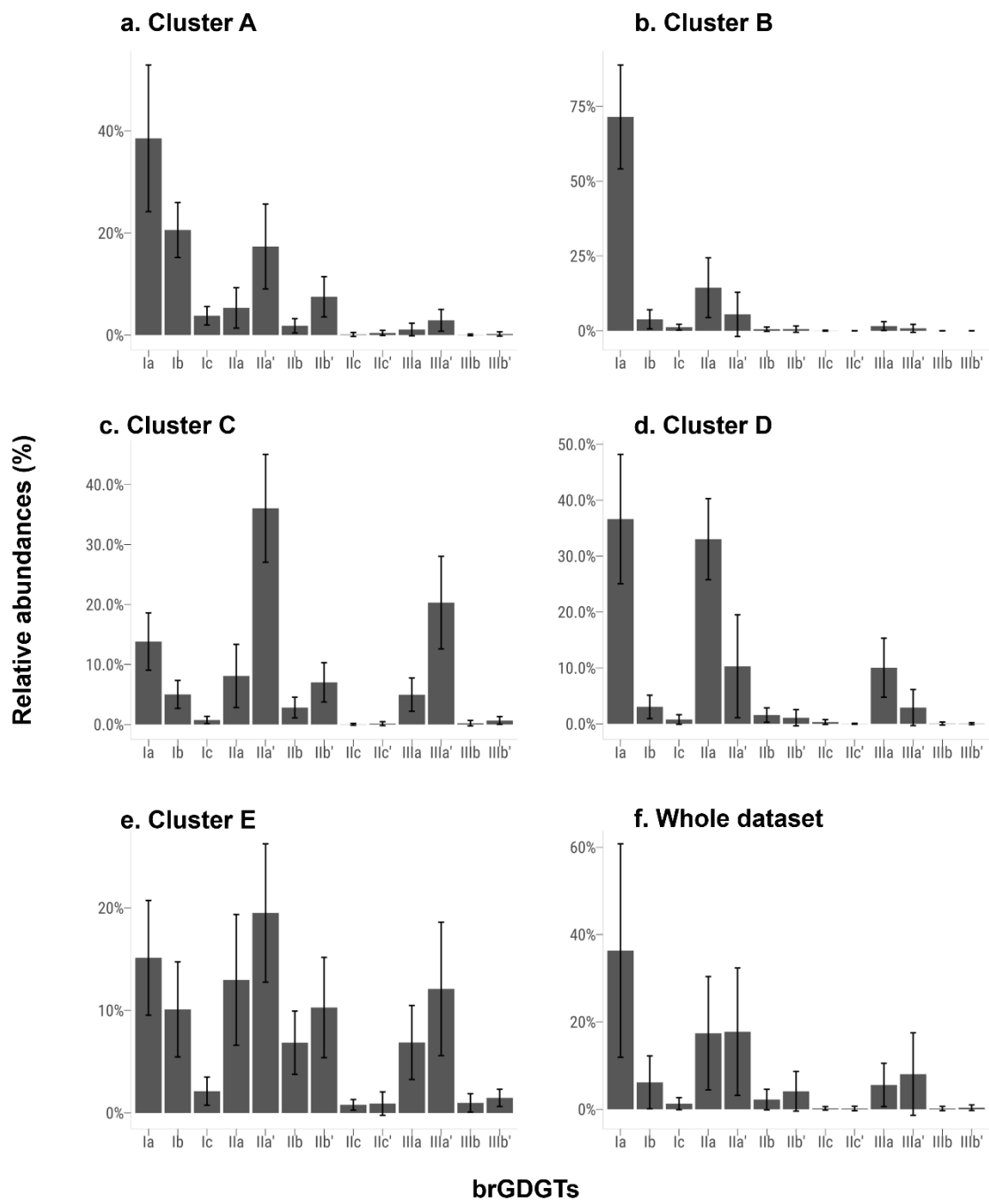


Figure 4

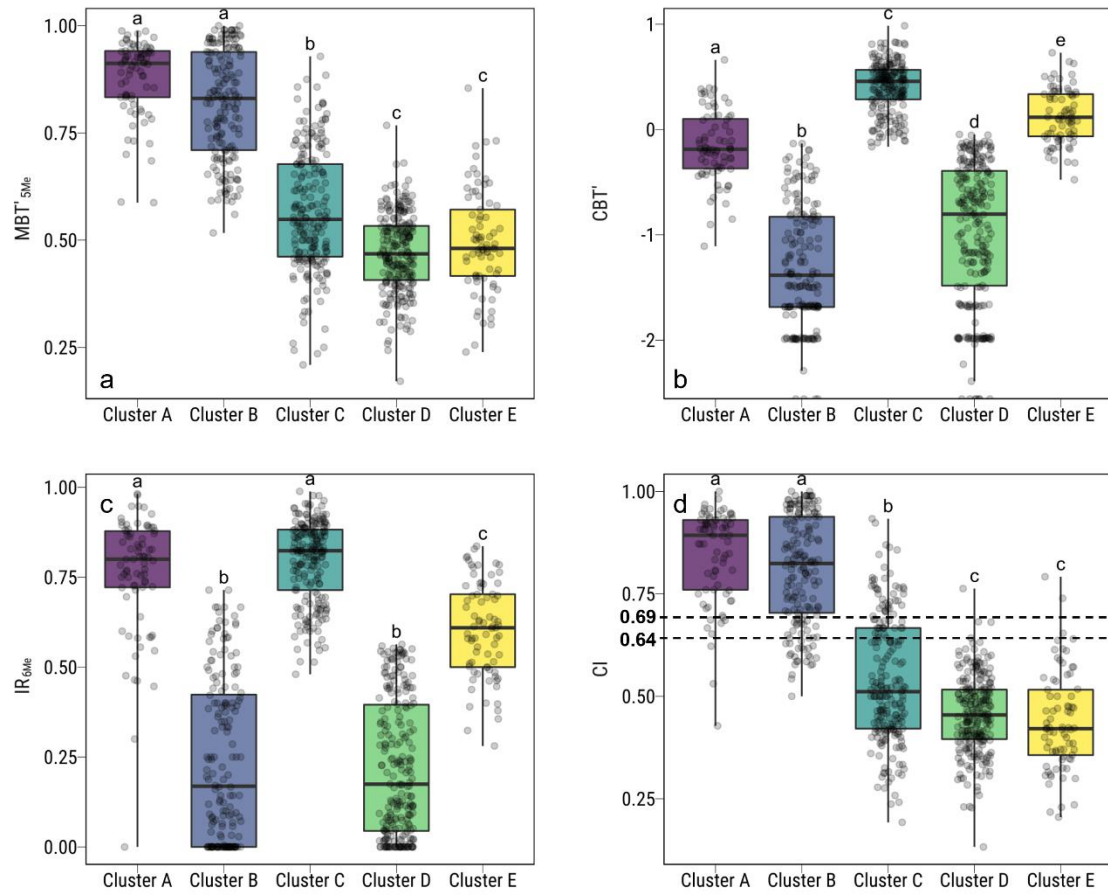


Figure 5

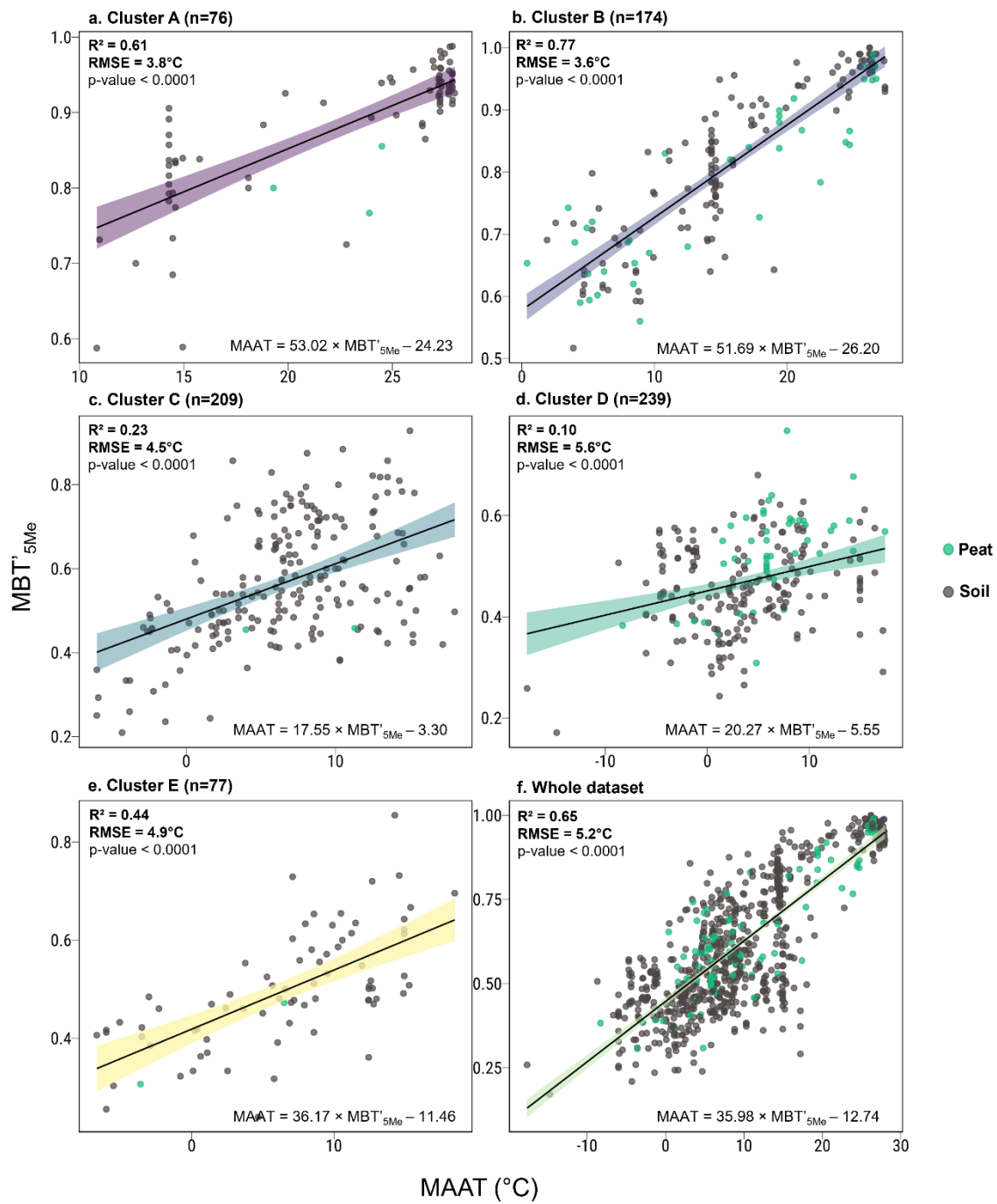


Figure 6

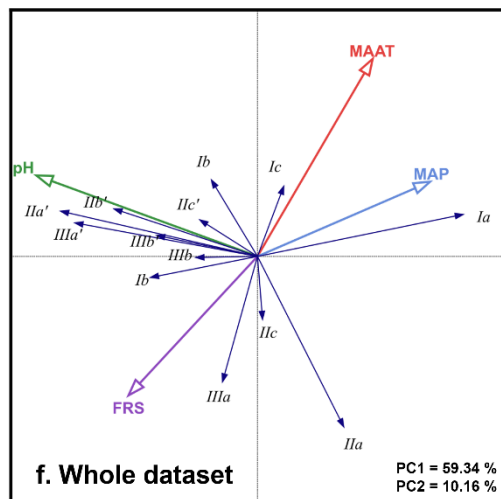
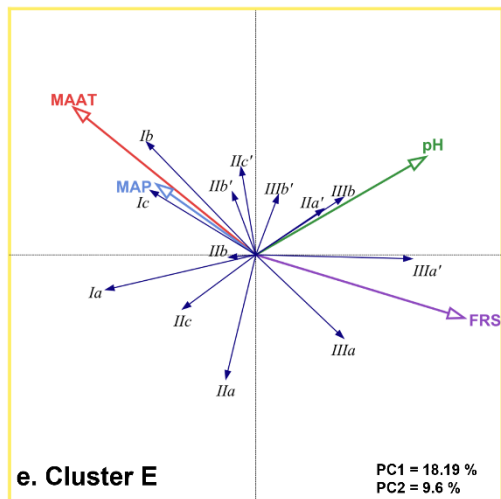
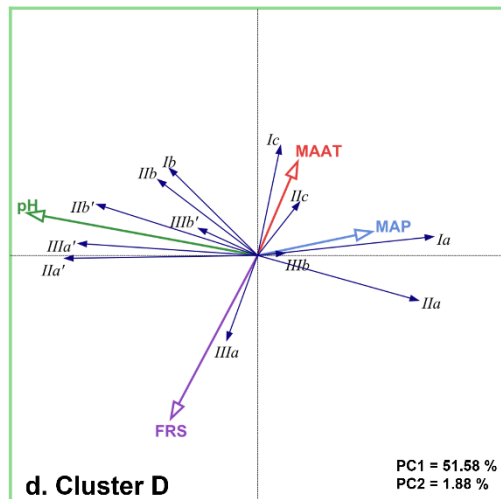
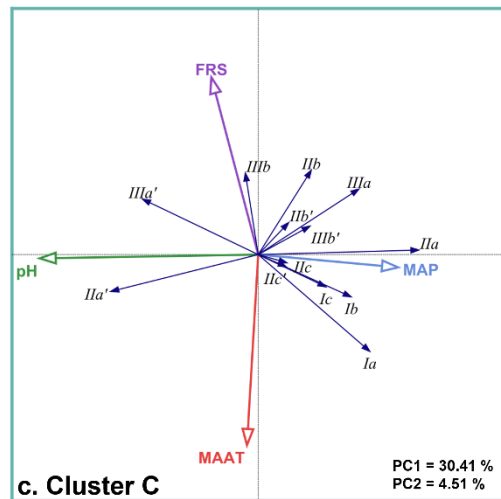
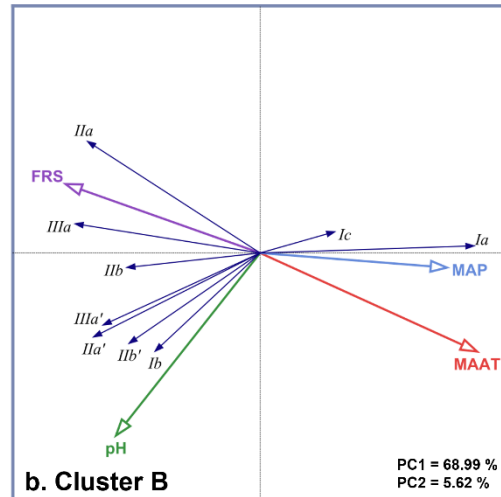
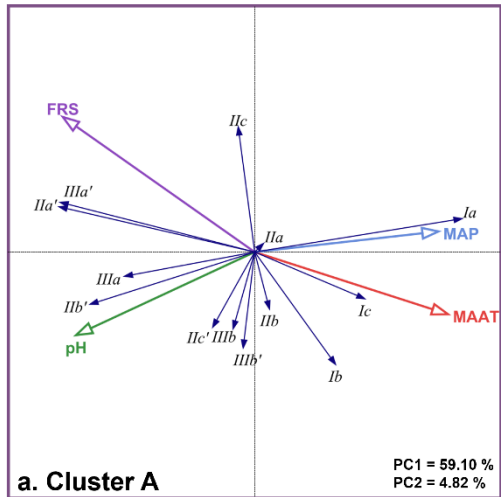


Figure 7

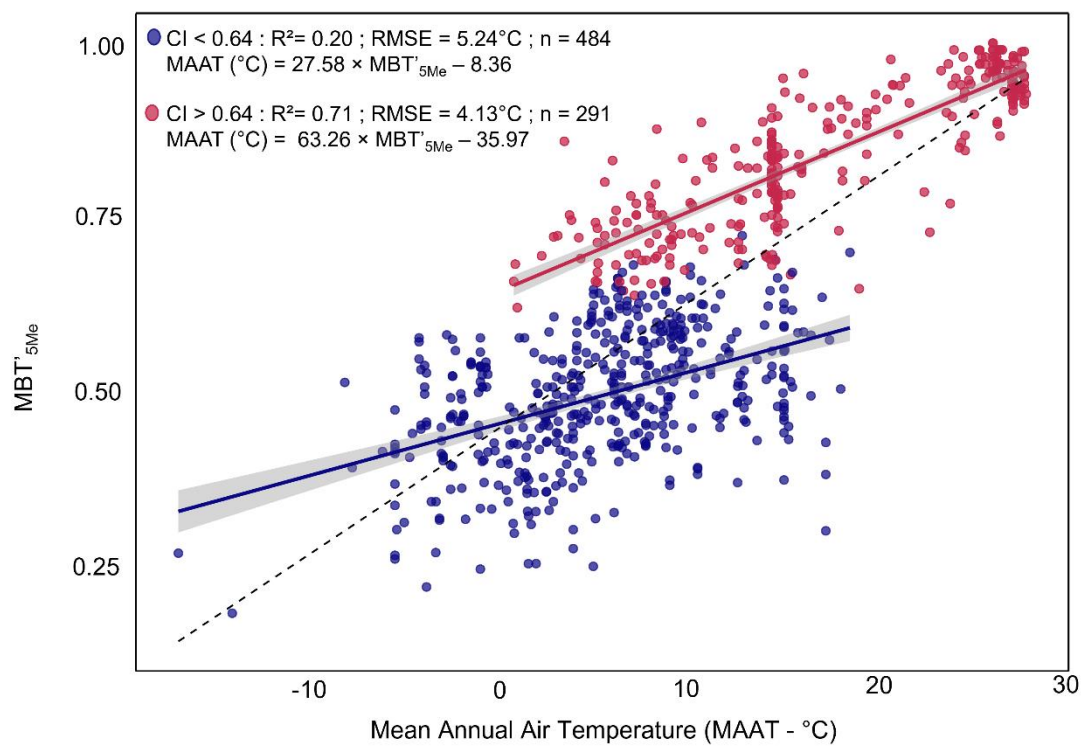


Figure 8

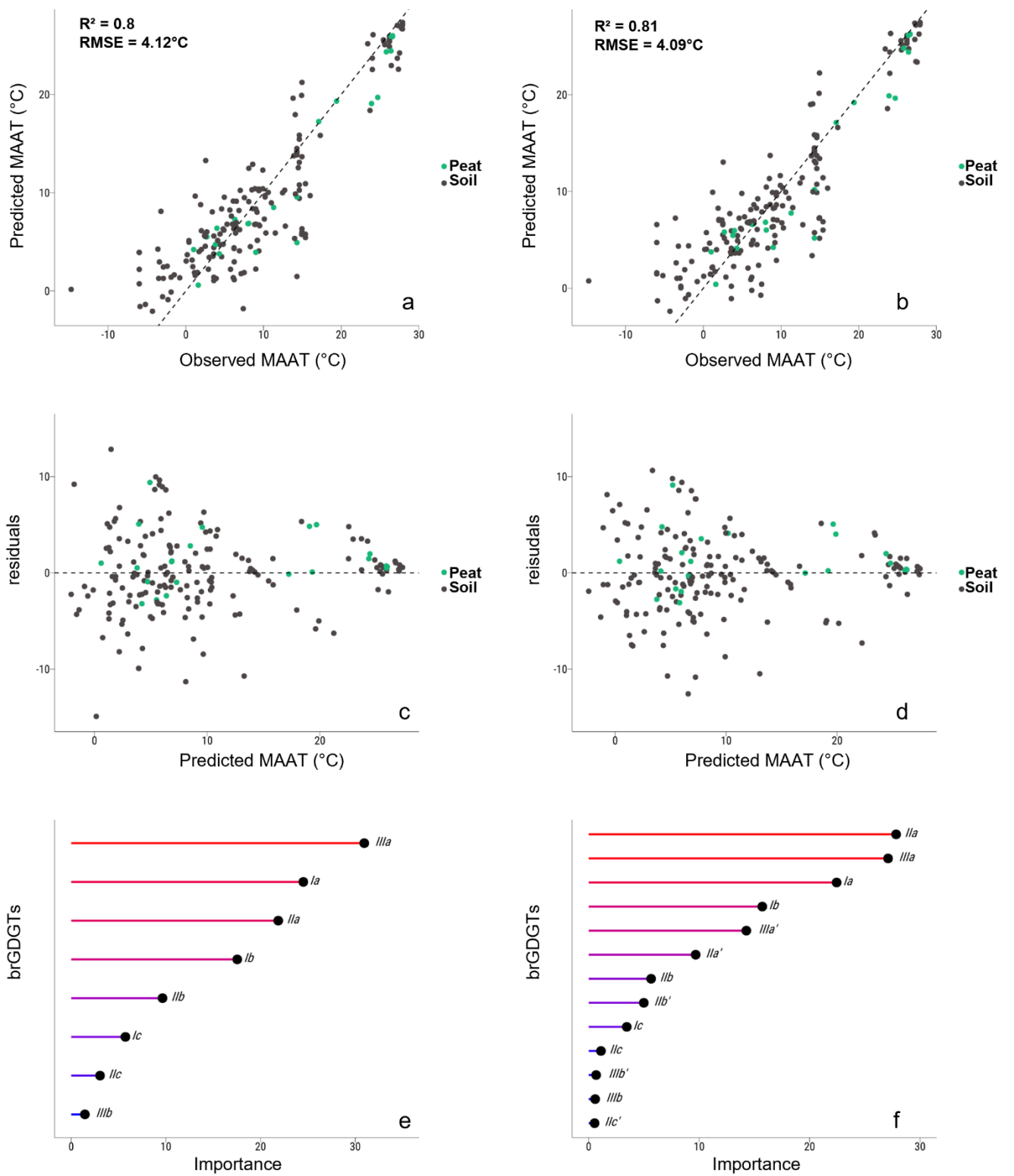


Figure 9

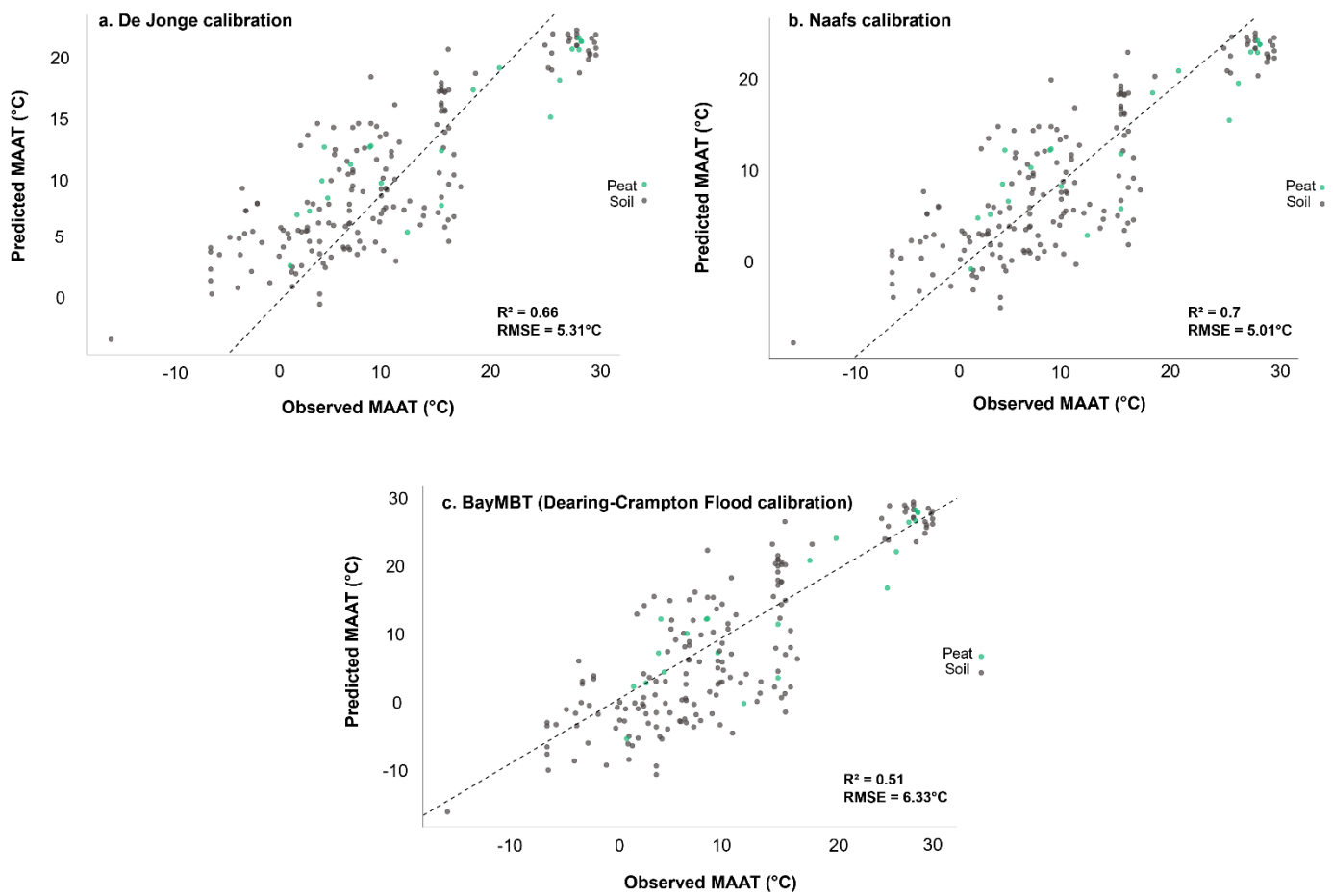


Figure 10

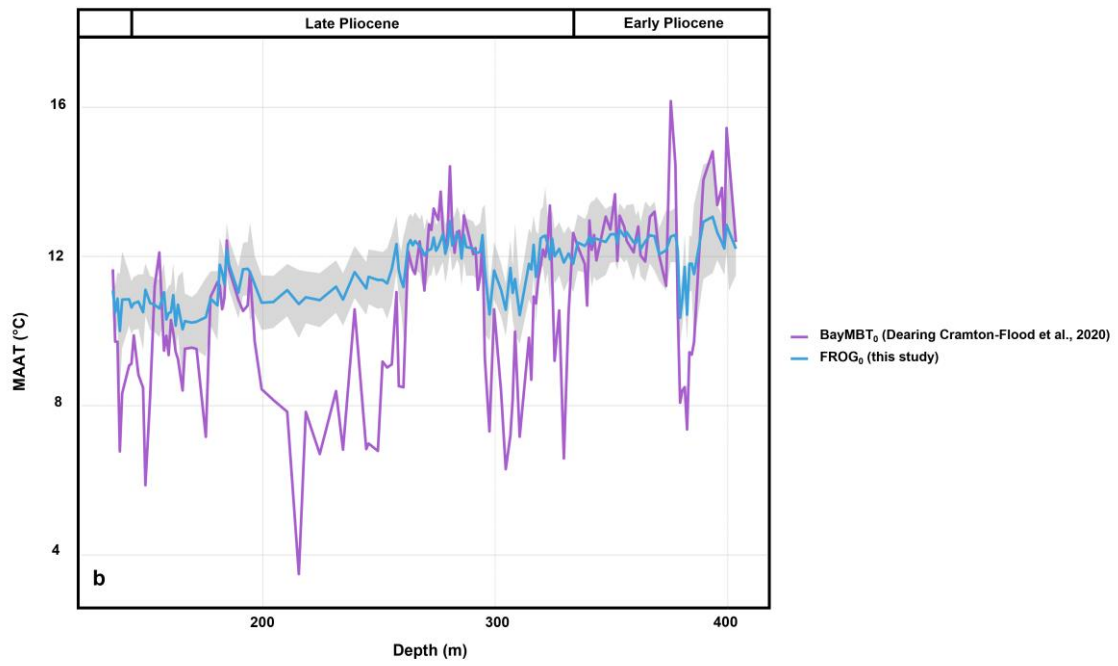
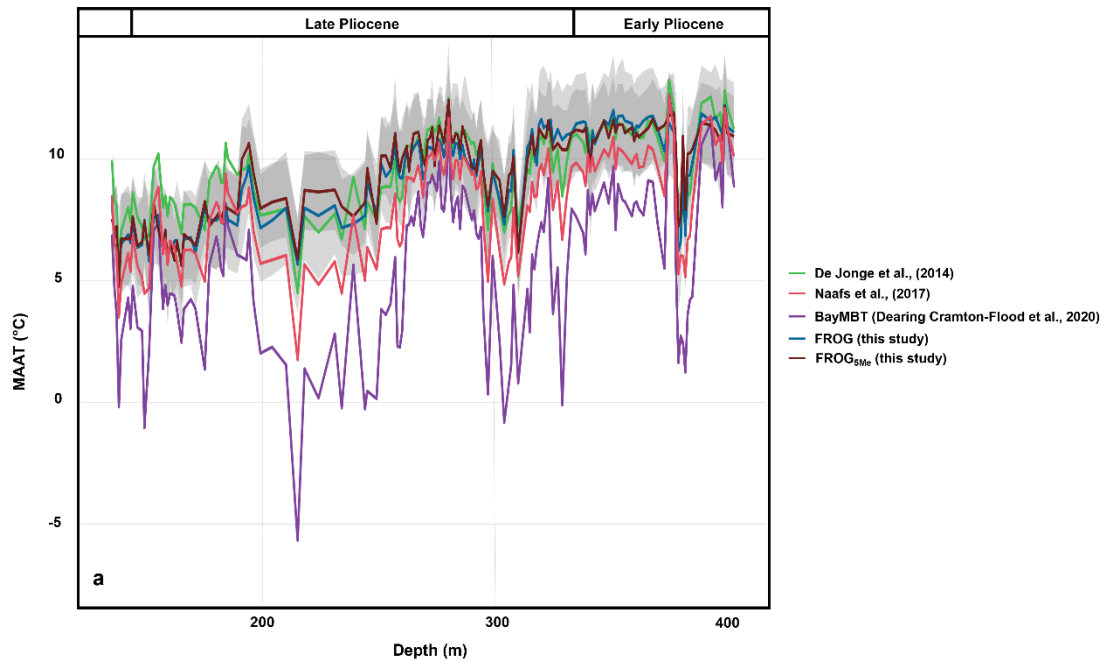


Figure 11

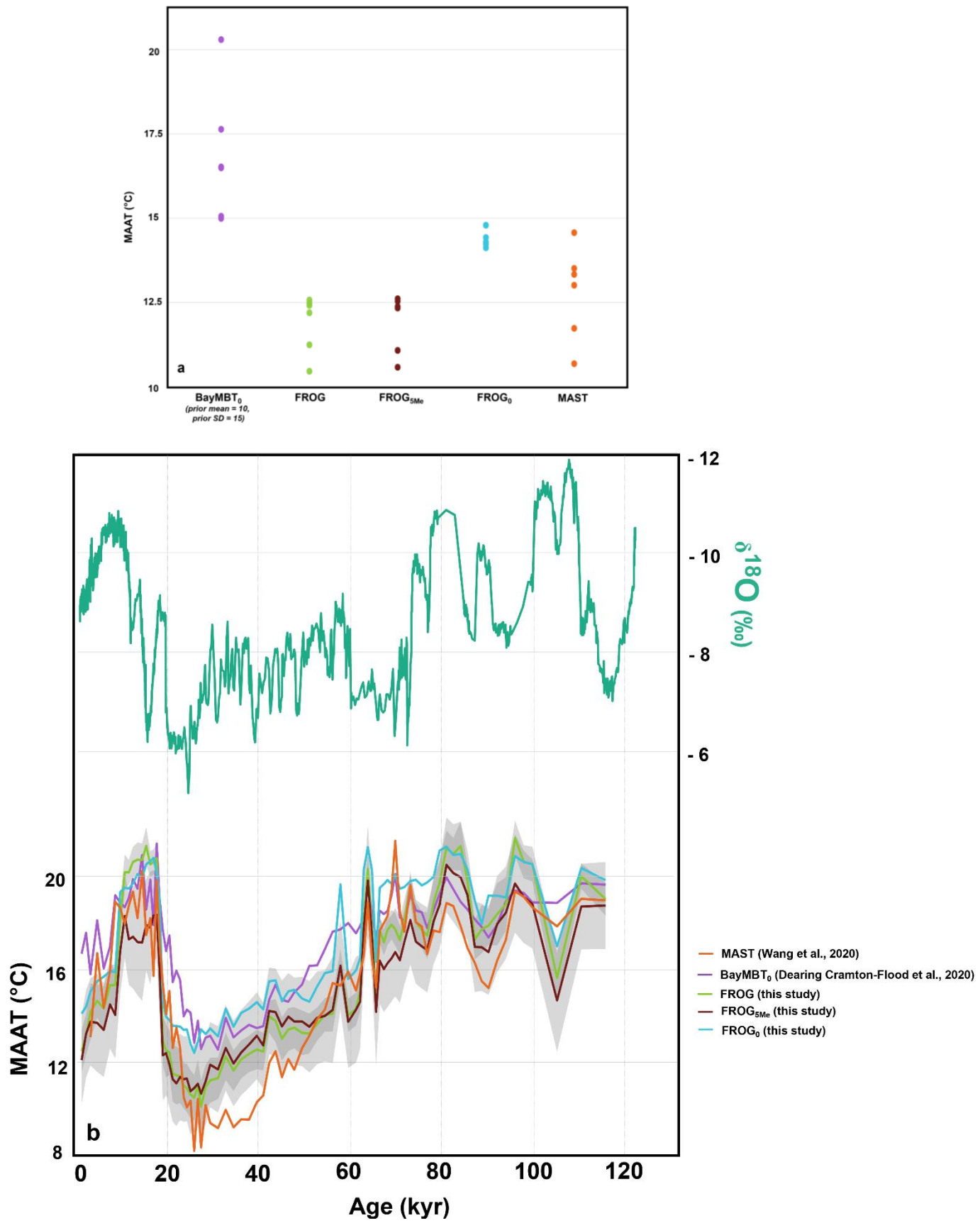


Figure 12

Location	N (samples)	Reference	Available parameters
Italy*	13	Véquaud et al., (2020)	MAAT, pH
Italy*	11	Huguet et al., (2019), Véquaud et al., (2020)	MAAT, pH
Tibet*	17	Véquaud et al., (2020)	MAAT, pH, MAP
Peru*	14	Véquaud et al., (2020)	MAAT, pH, MAP
Chile*	8	Véquaud et al., (2021)	MAAT, pH
France*	49	Véquaud et al., (2021)	MAAT, pH
Globally distributed	229	Weijers et al. (2007a), Peterse et al. (2012), De Jonge et al. (2014b)	MAAT, pH, MAP, FRS
India	46	Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
Russia/Siberia	4	Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
New Zealand	1	Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
China	15	Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
China	27	Xiao et al. (2015), Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
China	26	Yang et al. (2015), Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
China	27	Ding et al. (2015), Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
China	44	Lei et al. (2016), Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
China	148	Wang et al. (2016), Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
Globally distributed	96	Naafs et al. (2017b), Dearing Crampton-Flood et al., (2020)	MAAT, pH, MAP, FRS
TOTAL	775		

Table 1

	Cluster A	Cluster B	Cluster C	Cluster D	Cluster E	All
Soil	73	131	207	193	75	679
Peat	3	43	2	46	2	96
Samples	76	174	209	239	77	775
Min pH	5.6	3.0	4.9	3.4	4.0	3.0
Max pH	8.1	7.9	8.7	7.5	9.3	9.3
Mean pH	7.1	5.2	7.5	5.4	7.2	6.3
Min MAAT (°C)	10.8	0.4	-6.0	-17.7	-6.7	-17.7
Max MAAT (°C)	28.0	27.4	18.0	17.4	18.5	28.0
Mean MAAT (°C)	22.4	16.0	6.7	3.9	6.5	9.5
Min MAP (mm/yr)	106	383	168	128	156	106
Max MAP (mm/yr)	2770	3584	1055	2177	1191	3584
Mean MAP (mm/yr)	1069	1237	453	784	756	855
Min FRS	0	0	3	0	2	0
Max FRS	11	19	29	26	27	29
Mean FRS	2	5	15	14	13	11

Table 2

Variables	RDA correlation coefficients		Conditionnal effects (%)	
	Axis 1	Axis 2		
Cluster A	pH	-0.75	-0.34	18.6*
	MAAT (°C)	0.8	-0.26	5.3*
	MAP (mm/year)	0.77	0.08	2.2
	FRS	-0.8	0.56	39.6*
	Expl. variation (%)	59.10	4.82	
	Expl. fitted variation (%)	89.97	7.34	
Cluster B	pH	-0.6	-0.76	16.7*
	MAAT (°C)	0.9	-0.41	57.3*
	MAP (mm/year)	0.78	-0.06	0.2
	FRS	-0.81	0.29	0.5
	Expl. variation (%)	68.99	5.62	
	Expl. fitted variation (%)	92.42	7.52	
Cluster C	pH	-0.91	-0.02	25.6*
	MAAT (°C)	-0.05	-0.78	4.5*
	MAP (mm/year)	0.59	-0.05	5.7*
	FRS	-0.2	0.74	0.7
	Expl. variation (%)	30.41	4.51	
	Expl. fitted variation (%)	83.36	12.37	
Cluster D	pH	-0.96	0.18	47.4*
	MAAT (°C)	0.17	0.39	1.5*
	MAP (mm/year)	0.48	0.1	3.4*
	FRS	-0.36	-0.68	1.5*
	Expl. variation (%)	51.58	1.88	
	Expl. fitted variation (%)	95.9	3.5	
Cluster E	pH	0.71	0.41	7.7*
	MAAT (°C)	-0.76	0.62	6.6*
	MAP (mm/year)	-0.41	0.29	0.6
	FRS	0.87	-0.26	15.2*
	Expl. variation (%)	18.19	9.6	
	Expl. fitted variation (%)	62.5	31.9	
Whole dataset	pH	-0.92	0.34	51.8*
	MAAT (°C)	0.48	0.83	15.7*
	MAP (mm/year)	0.72	0.31	1.1*
	FRS	-0.54	-0.58	1.6*
	Expl. variation (%)	59.34	10.16	
	Expl. fitted variation (%)	84.52	14.46	

Table 3

	n (samples)	R ²	RMSE (°C)	Variance in residuals (°C)	Lower limit (°C)	Upper limit (°C)
Global Calibration						
FROG _{5Me}	192 (775)	0.8 (0.78)	4.12 (4.14)	17.0 (17.1)	-2.06 (-3.50)	27.33 (27.47)
FROG	192 (775)	0.81 (0.79)	4.09 (4.01)	16.7 (16.1)	-2.38 (-4.23)	27.47 (27.58)
Naafs calibration*	192 (775)	0.70 (0.64)	5.01 (5.23)	25.1 (27.4)	-7.81 (-7.81)	24.59 (24.59)
De Jonge calibration*	192 (775)	0.66 (0.61)	5.31 (5.45)	28.2 (29.7)	-3.19 (-3.19)	22.88 (22.88)
BayMBT*	192 (775)	0.52 (0.43)	6.33 (6.58)	40.1 (43.3)	-15.89 (-15.89)	29.83 (29.83)
MBT' _{5Me} (this study)	192 (775)	0.70 (0.65)	5.01 (5.20)	25.1 (27.0)	-6.73 (-6.73)	23.29 (23.29)
Alternative Calibration						
FROG ₀	164 (661)	0.83 (0.85)	2.53 (2.34)	6.4 (5.5)	5.39 (4.58)	27.02 (27.72)
BayMBT ₀ *	164 (661)	0.56 (0.54)	4.07 (4.07)	16.6 (16.6)	2.68 (1.43)	27.02 (27.68)
FROG ₅₀₀	108 (442)	0.85 (0.82)	3.56 (3.66)	12.7 (13.4)	-1.27 (-1.82)	27.48 (27.79)

Table 4