

1 **Development and application of the Branched and**
2 **Isoprenoid GDGT Machine learning Classification**
3 **algorithm (BIGMaC) for paleoenvironmental**
4 **reconstruction**

5 **Pablo Martínez-Sosa¹, Jessica E. Tierney¹, Lina C. Pérez-Angel², Ioana C.**
6 **Stefanescu³, Jingjing Guo⁴, Frédérique Kirkels⁴, Julio Sepúlveda⁵, Francien**
7 **Peterse⁴, Bryan N. Shuman³, Alberto V. Reyes⁶**

8 ¹Department of Geosciences, The University of Arizona, Tucson, Arizona, USA

9 ²Institute at Brown for Environment and Society (IBES), Brown University, Rhode Island, USA

10 ³Department of Geology and Geophysics, University of Wyoming, Wyoming, USA

11 ⁴Department of Earth Sciences, Utrecht University, Utrecht, Netherlands

12 ⁵Department of Geological Sciences and Institute of Arctic and Alpine Research (INSTAAR), University
13 of Colorado, Colorado, USA

14 ⁶Department of Earth and Atmospheric Sciences, University of Alberta, Edmonton, Alberta, Canada

15 **Key Points:**

- 16 • The distribution of GDGTs is particular to each depositional environment and has
17 unique responses to environmental factors.
- 18 • The BIGMaC algorithm captures the relationship between both branched and iso-
19 prenoid GDGTs with depositional environments.
- 20 • Our approach can provide paleoclimatological and paleoenvironmental informa-
21 tion based only on GDGTs.

Corresponding author: Pablo Martínez-Sosa, pmartoza@arizona.edu

Abstract

Glycerol dialkyl glycerol tetraethers (GDGTs), both archaeal isoprenoid GDGTs (isoGDGTs) and bacterial branched GDGTs (brGDGTs), have been used in paleoclimate studies to reconstruct environmental conditions. Since GDGTs are produced in many types of environments, their relative abundances also depend on the depositional setting. This suggests that the distribution of GDGTs also preserves useful information that can be used more broadly to infer these depositional environments in the geological past. Here, we combined existing iso- and brGDGT relative abundance data with newly analyzed samples to generate a database of 1153 samples from several modern sedimentary settings. We observed a robust relationship between the depositional environment and the relative abundances of GDGTs in our samples. This dataset was used to train and test the **B**ranch**e**d and **I**soprenoid **G**DGT **M**achine learning **C**lassification (BIGMaC) algorithm, which identifies the environment a sample comes from based on the distribution of GDGTs with high precision and recall ($F1 = 0.95$). We tested the model on the sedimentary record from the Giraffe kimberlite pipe, an Eocene maar in subantarctic Canada, and found that the BIGMaC reconstruction agrees with independent stratigraphic and palynological information, provides new information about the paleoenvironment of this site, and helps improve its paleotemperature reconstruction. In contrast, we also include an example from the PETM-aged Cobham lignite as a cautionary example that illustrates the limitations of the algorithm. We propose that in cases where paleoenvironments are unknown or are changing, BIGMaC can be applied in concert with other proxies to generate more refined paleoclimate records.

1 Introduction

Glycerol dialkyl glycerol tetraethers (GDGTs) are membrane-spanning lipids found in sediments and soils around the world. There are two main types of these molecules, branched and isoprenoid. Branched glycerol dialkyl glycerol tetraethers (brGDGTs) are characterized by their branched alkyl chains, with a differing number (4 – 6) and position (5-methyl or 6-methyl) of methyl groups and cyclopentane moieties (0 – 2). This unique structure defies the classical evolutionary dichotomy of the lipid divide by combining traits of Bacteria and Archaeal cell membranes (Weijers et al., 2006). Based on evidence including the presence of alkyl chains, the stereochemistry of the glycerol group (Weijers et al., 2006), and most importantly, microbial culture studies (Y. Chen et al., 2022; Halamka et al., 2022, 2021; Sinninghe Damsté et al., 2011), brGDGTs have a bacterial source.

In contrast, isoprenoid glycerol dibiphytanyl glycerol tetraether GDGTs (isoGDGTs) are produced by Archaea (de Rosa et al., 1977; Sinninghe Damsté et al., 2002). Their structures contain two phytane chains (Langworthy, 1977) and vary in the number of cyclopentane moieties (0 – 8) (De Rosa et al., 1983). Crenarchaeol is a member of this group of particular importance as it has been shown to be specifically produced by Thaumarchaeota (Sinninghe Damsté et al., 2002). Crenarchaeol contains four cyclopentane rings, one cyclohexane ring, and has one identified stereoisomer known as crenarchaeol⁷ (Sinninghe Damsté et al., 2002, 2018).

Both isoprenoid and branched GDGTs are used in paleoclimate studies as their distribution is correlated with variables such as temperature and pH, and these molecules are relatively stable through the geological record. In marine sediments, the degree of cyclization of isoGDGTs is related to overlying water temperature, forming the basis of the TetraEther index of 86 carbons (TEX₈₆) proxy (Schouten et al., 2002, 2013). Similarly, the methylation, cyclization, and isomerization of brGDGTs have been shown to respond to temperature and pH in terrestrial environments, such as peats, soils, lakes, and rivers (Raberg et al., 2022; Martínez-Sosa et al., 2020; Dang et al., 2018; De Jonge, Stadnitskaia, et al., 2014; Tierney et al., 2010; Weijers et al., 2007). The Methylation

73 index of Branched Tetraethers (MBT'_{5Me}) proxy isolates the relationship between the
74 methylation of brGDGTs and temperature (Weijers et al., 2007; De Jonge, Hopmans,
75 et al., 2014) and has been widely used for terrestrial paleoclimate reconstructions (De Jonge,
76 Hopmans, et al., 2014; Zheng et al., 2017; Naafs et al., 2018; Lauretano et al., 2021; Zhao
77 et al., 2022).

78 Across environments, GDGT distributions broadly reflect the microbial commu-
79 nity present. This is, for example, the basis of the Methane Index, which measures the
80 contribution of methanotrophic archaea relative to marine Thaumarchaeota in the sed-
81 imentary isoGDGT pool (Zhang et al., 2011). Likewise, the distribution of isoGDGTs
82 in marine systems reflects not only sea-surface temperature (captured by the TEX₈₆ in-
83 dex), but also the water depth (and potentially, different archaeal communities) from which
84 the isoGDGTs derive (Rattanasriampaipong et al., 2022; Taylor et al., 2013). Further-
85 more, previous work has found that the ratio of crenarchaeol/crenarchaeol' can indicate
86 which *Thaumarchaeota* group (I.1a or I.1b) is responsible for the production of these lipids
87 in lake sediments (Li et al., 2016). In terrestrial settings, De Jonge et al. (2019) proposed
88 the Community Index for brGDGTs, which is based on the inference that brGDGTs are
89 produced by different communities of bacteria, each with a unique response to soil tem-
90 perature. The combined use of some of the GDGTs, through the Branched and Isoprenoid
91 Tetraether (BIT) index, is used to broadly discriminate between marine and terrestrial
92 environments based on the dominance of brGDGT-producing bacteria in most terres-
93 trial settings and crenarchaeol-producing Thaumarchaeota in marine settings (Hopmans
94 et al., 2004). However, BIT values in soils, lakes, and peats all tend to be high, which
95 limits the ability of this index to reliably distinguish between these different types of ter-
96 restrial settings.

97 Building on these observations, we posit that the full range of archaeal and bac-
98 terial GDGTs (isoprenoidal and branched) contains information about their biological
99 precursors and the overall composition of the microbial community. This information,
100 based on ecological interactions rather than a physiological response, can in turn be used
101 to discriminate between sediments deposited in terrestrial or marine environments, as
102 well as whether terrestrial sediments are derived from freshwater, soil, or peatland en-
103 vironments. This would provide an additional tool for the identification of ancient de-
104 positional conditions in instances when it is not clear what the environment was, and
105 therefore could inform the reconstruction of environmental variables, *i.e.*, which GDGT-
106 based temperature proxy and calibration is most appropriate to use. Machine learning
107 provides a way to model highly dimensional and nonlinear data with complex interac-
108 tions and missing values (El Boucheffy & de Souza, 2020). In this case it allows us to
109 investigate drivers of the abundance of 19 GDGT structures (high dimensionality) which
110 have a complex and non-linear relationship to the environment, making this an ideal ap-
111 proach to extract the environmental information present in the distribution of GDGTs.

112 Machine learning has previously been used in the Geosciences to discriminate be-
113 tween magma (Ueki et al., 2018) as well as identifying the source of water from oil wells
114 (Engle & Brunner, 2019). In the field of biomarker-based paleoclimatology, classifica-
115 tion algorithms have been applied to identify sources of alkenones (Zheng et al., 2019),
116 as well as plant waxes (Peuple et al., 2021). Machine learning regression algorithms, as
117 well, as deep neural network applications, have also been applied to GDGTs as in order
118 to generate temperature calibrations (Dunkley Jones et al., 2020; Véquaud et al., 2022;
119 Zheng et al., 2022). Here, we use a compilation of GDGT distributions in 1153 globally
120 distributed soils, peats, and sediments from diverse depositional environments to train
121 a classification algorithm which is capable of identifying the environment in which a sam-
122 ple was formed based on the distribution of both branched and isoprenoid GDGTs. We
123 then demonstrate the application of this algorithm by using it to interpret the paleoen-
124 vironment and the paleotemperature in a Paleogene deposit that records a transition from
125 a lacustrine to a peatland environment. We also highlight the limitations of this approach

126 in an application to a peatland dataset that spans the onset of the Paleocene-Eocene Ther-
 127 mal Maximum (PETM).

128 2 Materials and Methods

129 2.1 Global Dataset

130 We compiled a total of 1153 globally distributed (Fig. 1 and Table SII) samples
 131 from different depositional environments: marine, lake, peat, river, and soil. These sam-
 132 ples all have quantified relative abundances for the full suite of the most commonly used
 133 isoGDGTs (GDGT-0, GDGT-1, GDGT-2, GDGT-3, crenarchaeol, and crenarchaeol')
 134 and brGDGTs (IIIa, IIIa', IIIb, IIIb', IIa, IIa', IIb, IIb', IIc, IIc', Ia, Ib, and Ic) in pa-
 135 leoenvironmental reconstructions, and were all analyzed with the updated High Perform-
 136 ance Liquid Chromatography-Mass Spectrometry (HPLC-MS) method of Hopmans et
 137 al. (2016). From the 1153 samples, 475 are peat (Naafs et al., 2018), 215 are marine sed-
 138 iments (this study), 196 are soil (Guo, Ma, et al., 2022; Dearing Crampton-Flood et al.,
 139 2020; Guo et al., 2020; Pérez-Angel et al., 2020), 162 are lake sediments (Martínez-Sosa
 140 et al., 2021; Guo et al., 2020), and 105 are riverbed sediments (Kirkels, Usman, & Pe-
 141 terse, 2022). For the Colombian and Inner Mongolia soil samples (Guo, Ma, et al., 2022;
 142 Pérez-Angel et al., 2020) we include here newly reported isoGDGT values not part of
 143 the original dataset.

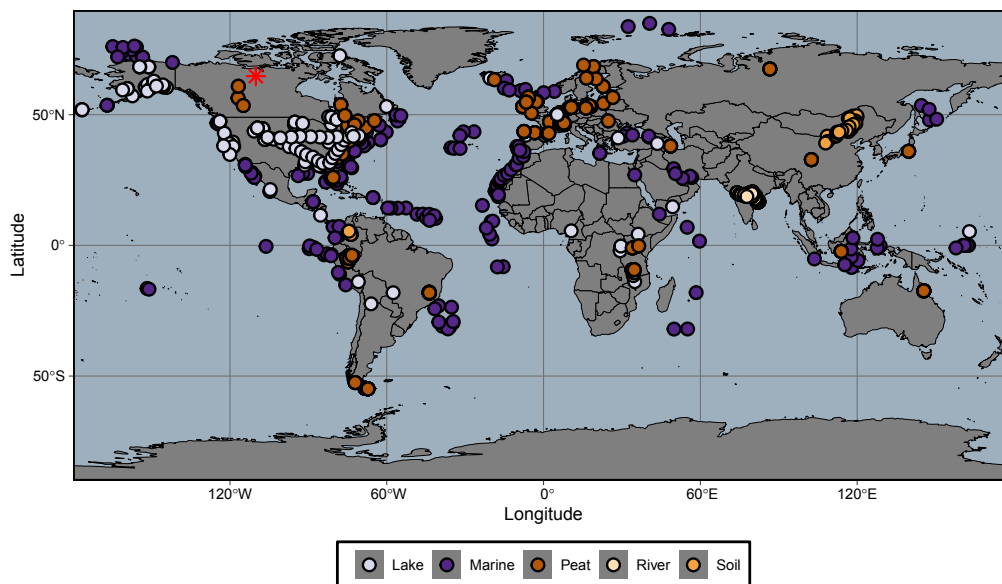


Figure 1. World map showing the distribution of the samples included in this work. Color code reflects the depositional environment which these samples were collected from. Red asterisk shows the modern location of the Giraffe pipe.

144 All marine sediments were processed at the University of Arizona following the method
 145 used in Martínez-Sosa et al. (2021). Briefly, sediments were freeze-dried, homogenized,
 146 and spiked with a C_{46} internal standard before being extracted with an Accelerated Sol-
 147 vent Extraction (ASE) system (run at 1500 psi, 100°C, with dichloromethane:methanol
 148 (DCM: MeOH, 9:1)). Total Lipid Extracts (TLEs) were eluted through a deactivated
 149 SiO_2 column with hexane:ethyl acetate (1:2) to obtain the polar fraction, and the elu-
 150 tant was dried under a N_2 stream. Polar fractions were redissolved in hexane:isopropanol
 151 (99:1), and then passed through a 0.45 μm PTFE filter prior to being analyzed by HPLC-

MS. GDGTs, isoprenoid and branched, were analyzed on an Agilent 1260 Infinity HPLC coupled to an Agilent 6120 single quadrupole mass spectrometer using two BEH HILIC silica columns (2.1×150 mm, 1.7 μm ; Waters) following the methodology of Hopmans et al. (2016). We calculated peak areas using the MATLAB package ORIGAmI (Fleming & Tierney, 2016) and quantified brGDGTs by comparing the obtained peaks with the internal standard (Huguet et al., 2006).

For all samples in this dataset we calculated the relative abundance of all brGDGTs (except IIIc and IIIc', due to their general low abundance), as well as isoGDGTs 0–3, crenarchaeol, and its isomer. For all the analyses we used the fractional abundance of each compound relative to the total sum of GDGTs (branched + isoprenoid), to account for the difference in the relative abundances of both GDGT types among the depositional environments. Although it is known that the ionization of isoGDGTs and brGDGTs in the MS might be different between laboratories (Schouten et al., 2013), the potential impact of this is minimized in our statistical approach because the data are normalized before applying the machine learning techniques (see Section 2.2).

We collected the environmental parameters associated with the samples using the data available in the source datasets (Table S1). For the marine sediments analyzed for this study, we obtained mean annual temperature of the top 200m of the water column from the World Ocean Atlas 2018 (Locarnini et al., 2018).

2.2 Unsupervised Machine Learning

For the unsupervised machine learning analysis we centered and scaled the fractional abundances of GDGTs across the whole dataset. The optimal number of clusters for this dataset was calculated through a silhouette analysis, which calculates how similar a data point is within-cluster compared to other clusters. This analysis was performed by using the the Partitioning Around Medoids method from the *cluster* R package (Maechler et al., 2019).

Samples were separated into clusters by applying the fuzzy version of the k-means clustering algorithm using the *e1071* R package (Meyer et al., 2020). This method calculates the degree to which each sample belongs to each of the clusters (membership value) instead of assigning a single classification; we consider this a useful tool to classify depositional environments, which have diffuse boundaries. The fuzzy k-means analysis was performed using the best performing number of clusters from the silhouette analysis, with all other parameters of the function at default values.

Following the cluster analysis, we compared the cluster assignment of each sample to the available information on their environmental data and depositional environment. Combining both the statistical results and the observation-based information we assigned one of four new labels to each of the samples.

2.3 Supervised Machine Learning

For the supervised machine learning we worked in the *tidymodels* and *tidyverse* R environments (Kuhn & Wickham, 2020; Wickham et al., 2019), where we used the fractional abundances of GDGTs as predictor variables and the statistical and observation based labels as the response variable. The dataset was split into a training and testing set in a 3:1 ratio. To avoid subsampling the dataset in a biased manner, we preserved the distribution of sample types in both sets. We tested the performance of four different algorithms commonly used for classification applications: Random Forest, eXtreme Gradient Boosting (XGBoost), K-Nearest Neighbour and Naive Bayes plus a control non-informative (null) model. Below we present a brief overview of each of them.

199 Random Forest is an ensemble classification algorithm, where classification of a sam-
200 ple is based on the voting results of an ensemble of independent decision trees. Impor-
201 tantly, each tree is presented with a randomly selected subset of samples and predictor
202 variables, ensuring that the trees generate independent results. This enables the result
203 from the voting to account for biases present in each individual decision tree. Random
204 Forest algorithms are considered reliable and fast (Parmar et al., 2019).

205 XGBoost is a gradient tree boosting algorithm. In this case, similar to Random For-
206 est, the algorithm uses a series of trees to classify samples. However, XGBoost iteratively
207 improves the performance of this process by improving on the results of each previous
208 iteration. This algorithm is particularly good for biased data and has been a preferred
209 algorithm for diverse applications due to its scalability (T. Chen & Guestrin, 2016).

210 Naive Bayes is a statistical algorithm, in which Bayes Theorem is applied to calcu-
211 late the posterior probability of a new sample belonging to each possible categorical
212 classification. The algorithm determines the classification by choosing that with the high-
213 est probability. Although this is a fast classification algorithm it is a relatively bad es-
214 timator (Sen et al., 2020).

215 K-nearest neighbor is another statistics-based method, where samples are classi-
216 fied as the voting result of their k closest neighbors. Neighbors are given priority based
217 on their closeness to the new data. While this is an effective algorithm for large datasets,
218 it can be computationally expensive as all the distances from the neighbors need to be
219 calculated for all new data (Sen et al., 2020).

220 For all algorithms the hyperparameter values—parameters whose values control the
221 learning process but are not part of the final model (*i.e.* number of independent trees
222 used for a Random Forest)—were selected (tuned) as those with the best performance
223 from a distribution of possible combination of values for each hyperparameter. Due to
224 the required computing power for the hyperparameter tuning step, this was run using
225 a High-Performance Computing cluster. Finally, the best hyperparameter values were
226 selected by comparing their Receiver Operating Characteristic Area Under the Curve
227 (ROC-AUC) score on the validation set (Table S1). ROC-AUC is a metric that measures
228 the tradeoff between the true positive rate and false positive rate of the model, for this
229 parameter higher values on the (0,1) range are desirable. An additional metric, F1, was
230 applied to evaluate the performance of all the trained algorithms. F1 is calculated as the
231 mean of the precision and recall—the ability to identify true positives—of the model. For
232 reproducibility, a detailed script with the packages and parameter values used for this
233 process is available on GitHub (Martínez-Sosa et al., 2023). Finally we can identify which
234 GDGTs contribute the most to the classification process through the the importance met-
235 ric. This value is calculated based on how much each GDGT contributes to decreasing
236 the probability of incorrectly classifying a sample across all decision trees (Gini impor-
237 tance) (Wright et al., 2019; Greenwell et al., 2020). This metric shows how often a GDGT
238 was selected to split the data and what was its discriminative power (Menze et al., 2009),
239 larger values indicate a better variable to separate the data.

240 For this work, all analyses were performed in R (v. 4.1.3) (R Core Team, 2022). Ad-
241 ditional dimensionality reduction analyses were done over the fractional abundances of
242 all GDGTs using Principal Component Analysis (PCA) through the `princomp()` func-
243 tion from base R. For the PCA analyses the loadings of variables were visually inspected
244 and corroborated by obtaining the eigenvector values. Correlation between environmen-
245 tal parameters for each sample and their scores on the principal components was per-
246 formed by applying the `cor.test()` function from the *psych* R package using a Spear-
247 man correlation (Revelle & Revelle, 2015).

2.4 Giraffe Kimberlite Pipe

We analyzed GDGTs from 83 samples from diamond exploration drill core BHP 99-01 from the Giraffe kimberlite pipe (paleolatitude $\sim 63^\circ\text{N}$) (Wolfe et al., 2017). This core is stored at the Geological Survey of Canada core repository (Calgary), and it contains ≥ 50 vertical-equivalent meters of lacustrine sediment topped with ~ 32 m of peat. The sediments were dated to 37.84 ± 1.99 Ma by glass fission-track dated rhyolitic tephra beds (Wolfe et al., 2017). Our dataset spans 83.5 vertical-equivalent meters and includes 19 samples from the peat section and 64 from the lacustrine section. For each sample, between 0.5 and 1 g of sediment was processed to obtain TLEs in the same manner as for the marine samples. For these samples, the GDGTs were isolated using a two-layer chromatography column filled with a 1:1 mix of LC-NH₂ (bottom layer) and 5% deactivated silica (top layer) gels as the solid phase (Windler et al., 2019). The GDGTs were recovered using dichloromethane:isopropanol (2:1) as the solvent. Branched and isoprenoid GDGTs were analyzed in all samples using the same HPLC-MS method described for the marine samples in section 2.1.

2.5 Cobham Lignite Bed

The Cobham lignite bed, Kent, UK ($\sim 48^\circ\text{N}$ palaeolatitude) is composed of a sand and mud unit at the base, overlain, in succession, by a charcoal-rich lower laminated lignite, a charcoal-poor upper laminated lignite, a middle clay layer, and a charcoal-poor blocky lignite. The Woolwich Shell Beds overly the Cobham Lignite (Collinson et al., 2009). A carbon isotope excursion is present near the top of the charcoal-poor upper laminated lignite, which is interpreted as being the characteristic excursion from the Paleocene Eocene Thermal Maximum (PETM, ~ 56 million years ago). Collinson et al. (2009) interpreted the units above this as representing the early part of the PETM. We tested our algorithm on the 27 samples obtained from this site previously analyzed by Inglis et al. (2019) and publicly available at the PANGAEA data repository (Inglis et al., 2019).

3 Results

3.1 Fuzzy K-means Classification

Our silhouette analysis showed that the global GDGT data is best separated into four clusters, a value which was then used to perform a fuzzy k-means classification. The four identified groups consist of between 219 and 465 samples each. When we compare the composition of each cluster using PCA, there is a clear differences between depositional environments (Fig. 2a and b, and Table 1). 86.9% of the peat samples fall within Group 1, while 84.6% of the lacustrine samples are assigned to Group 2. In turn, 92.4% of the river samples are assigned to Group 3, and 91.6% of the marine samples are assigned to Group 4 (Fig. 2a and b). Soil samples are more spread across the different groups, with the majority assigned to Group 3 (43.9%).

3.2 Within-Group Analyses

We analyzed the GDGT distribution within each of the clusters identified in the unsupervised machine learning step to assess its influence on the clustering results and how well it correlated with environmental parameters.

3.2.1 GDGT Distribution

Across the entire dataset crenarchaeol', GDGT-1–GDGT-3, Ib, Ic, IIc, IIc', IIIb, and IIIb' have the smallest proportion (< 0.1 fractional abundance) of all GDGTs (Fig. 3). There are, however, characteristic patterns associated with the different clusters. Samples from cluster 4 have a higher proportion of crenarchaeol, GDGT-0, and to a lesser

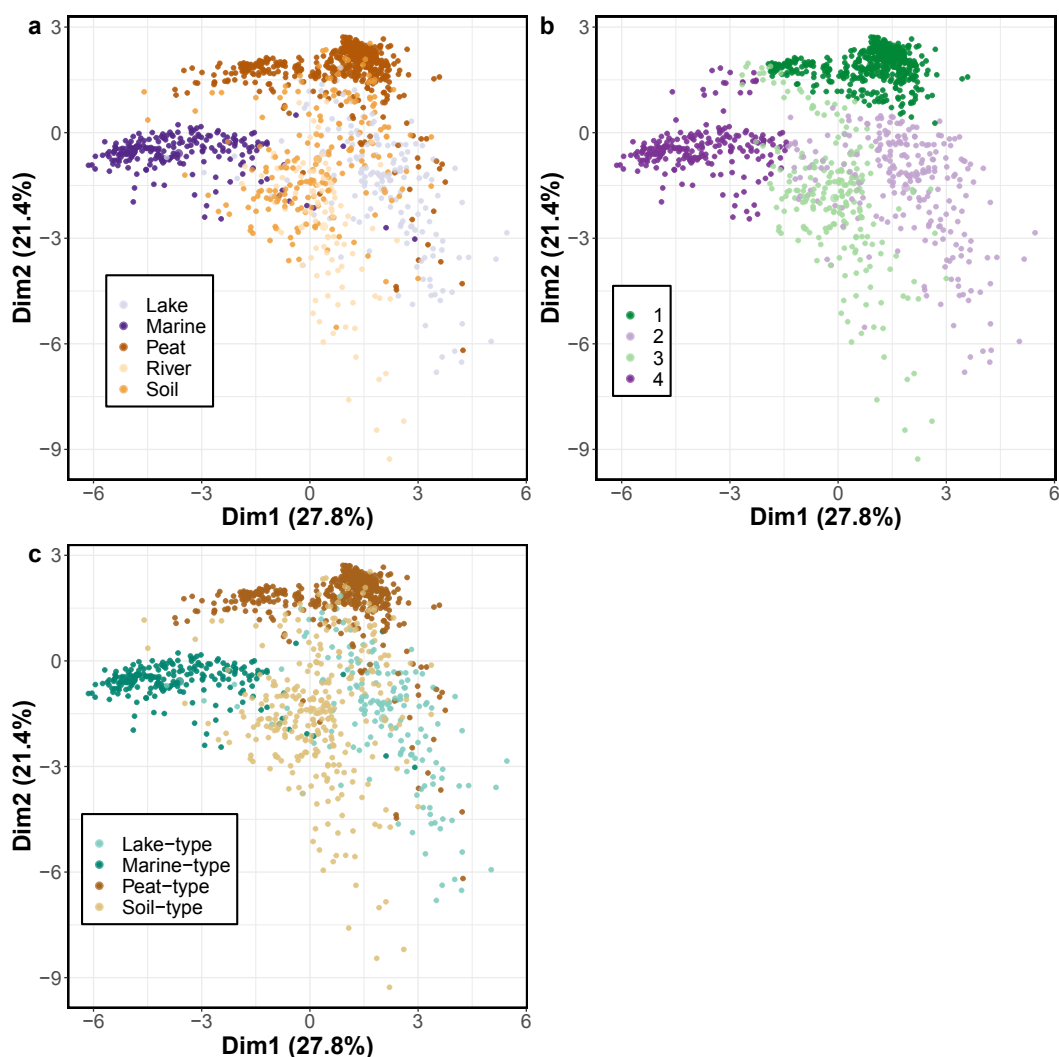


Figure 2. Samples from the dataset plotted in reduced dimensional space based on the fractional abundance of GDGTs. Plots show the same analysis with samples colored based on the depositional environment (a), their assigned group based on the fuzzy k-means analysis (b), and the curated clusters (c)

294 extent GDGT-1 and GDGT-2 compared with the other clusters (Fig. 3a). For crenar-
 295 chaerol, cluster 3 is the next group with the highest proportion, while cluster 2 has the
 296 second highest proportion of GDGT-0. Cluster 1 consistently shows the lowest isoGDGT
 297 values. On the other hand clusters 1, 2, and 3 have a higher proportion of brGDGTs than
 298 cluster 4. Cluster 1 has the highest proportion of brGDGTs Ia and IIa, while cluster 2
 299 has the highest proportion of IIIa and IIIa'. Although cluster 3 has lower proportion val-
 300 ues than cluster 1, it is also dominated by the penta- and tetramethylated brGDGTs,
 301 and it shows the highest proportion of GDGT Ib. Both clusters 2 and 3 show compar-
 302 able levels of IIa', the most abundant 6-methyl brGDGT.

303 3.2.2 *Environmental influence on GDGTs*

304 To better understand the effect that environmental parameters (such as temper-
 305 ature, elevation, and pH) might have had on the classification of the samples, we per-

Table 1. Percentage of the total of each sample type assigned to each of the four clusters determined by fuzzy k-means analysis (left) as well as the four manually curated clusters (right). At the right is the total number of samples from each type. The highest percentage for each type of sample in the fuzzy k-means clusters is indicated in bold.

Type	C. 1	C. 2	C. 3	C. 4		Peat-type	Lake-type	Soil-type	Marine-type	Total
Lake	7.4	84.6	5.6	2.5		0.6	97.5	1.2	0.6	162
Marine	0	5.6	2.8	91.6		0	0	0	100	215
Peat	86.9	5.7	4.4	2.9		100	0	0	0	475
River	0	7.6	92.4	0		0	0	100	0	105
Soil	20.4	30.6	43.9	5.1		0	0	100	0	196

306 formed a principal component analysis for each of the sample types separately (Fig. 4,
307 Fig. SI1, and Table SI2).

308 For peats (Fig. 4a), primarily classified in cluster 1, the main component (64.6%
309 of variability) is positively associated with GDGT Ia, while negatively related to GDGT
310 IIa. This component is strongly associated with MAAT ($\rho = 0.77$, Spearman's corre-
311 lation and Fig. SI1a). Peats classified in cluster 1 plot throughout the first component,
312 however, those peats classified as part of cluster 3 and 4 are associated with higher tem-
313 peratures, while those classified as part of cluster 2 are associated with lower MAATs
314 (Fig. 4a).

315 For lake sediments (Fig. 4b), which are mostly classified as part of cluster 2, the
316 first component (36.6% of variability) is associated negatively with GDGT Ia, while pos-
317 itively related to GDGT IIIa. The second component (32.4% variability) is positively
318 associated with GDGT-0. Both components have a strong correlation with MAAT ($\rho =$
319 -0.65 and 0.51 , respectively, and Fig. SI1b). While the majority of the distribution cloud
320 is classified in cluster 2, samples with more negative values on PC1, which have higher
321 MAAT values, are classified as cluster 3 and 1 (Fig. 4b and Fig. SI1b). Additionally, four
322 samples are classified in cluster 4, which were identified as sediments from Lake Kivu,
323 Mono Lake, and two from Lake Malawi.

324 Soil samples are classified into a wide range of clusters (Fig. 4c), with 44% of the
325 samples labeled as part of cluster 3, 31% as cluster 2, 20% as cluster 1, and 10% as clus-
326 ter 4. The first component of the PCA for this sample type (47% of variability) is neg-
327 atively associated with brGDGT Ia, while the second component is positively related with
328 GDGT IIa and negatively with crenarchaeol. The environmental parameters show strong
329 correlations to both the first and second components: pH ($\rho = 0.47$ and -0.74), MAAT
330 ($\rho = -0.45$ and -0.42), and elevation ($\rho = -0.21$ and -0.56 , respectively). Additionally, the
331 assigned cluster strongly correlates with the sample location (Fig. SI1c). Most samples
332 from the Godavari river catchment (98%) and Inner Mongolia (69%) were classified as
333 cluster 3, while 70% of samples from Carminowe creek are classified as cluster 2. Finally,
334 47% of samples from Colombia are classified as cluster 1.

335 Similar to soils, the river sediments are mostly classified as cluster 3 (Fig. 4d). For
336 these samples the first principal component explains 51.4% of the variance and is neg-
337 atively associated with brGDGT Ia. For this sample type, the most determinant vari-
338 able for their cluster classification was location site (Fig. SI1d), with all but one out of

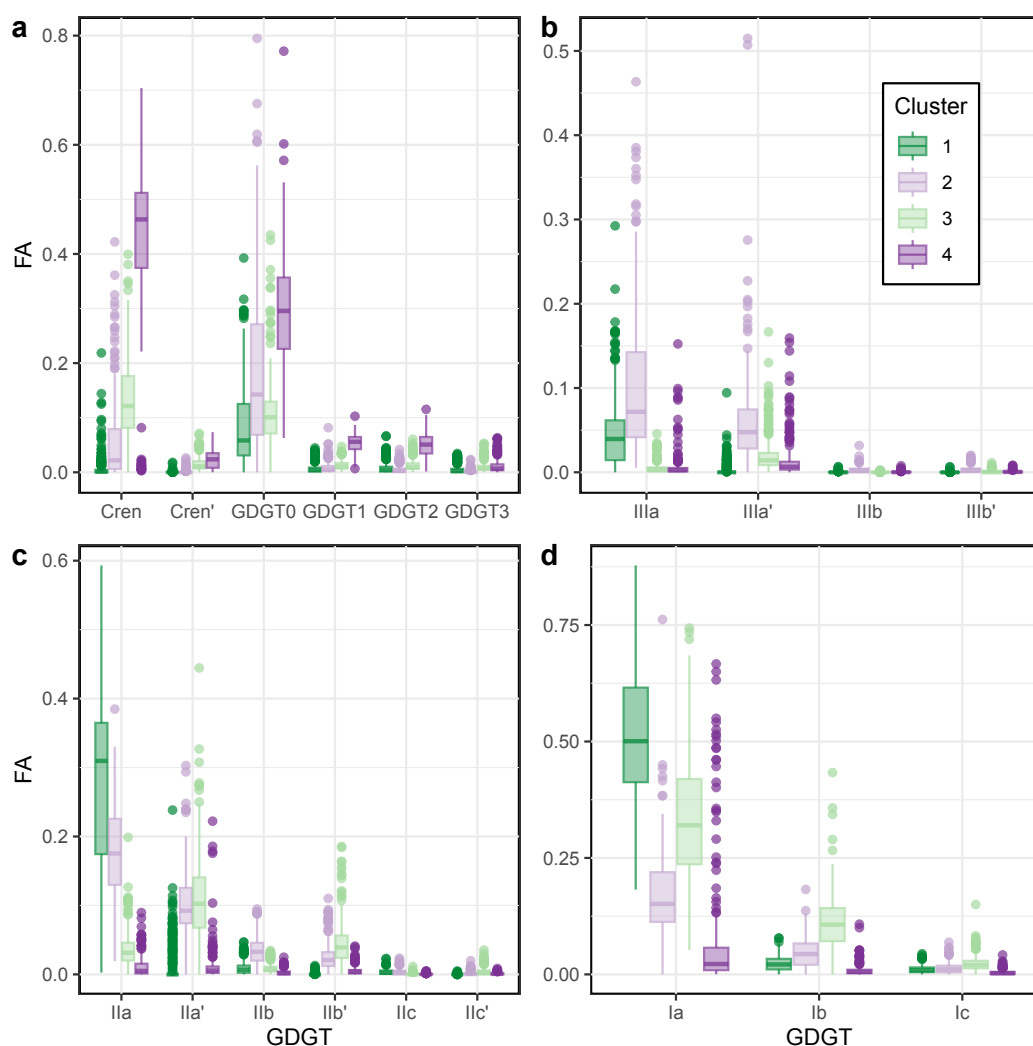


Figure 3. Box plots showing the distribution of the fractional abundance (FA) of all GDGTs in each of the fuzzy k-means clusters, following the color code of Figure 2b. GDGTs separated by isoGDGTs (a), hexamethylated brGDGTs (b), pentamethylated brGDGTs (c), and tetramethylated brGDGTs (d).

339 98 samples from the Godavari classified in cluster 3, while all samples from Carminowe
 340 Creek are classified as cluster 2.

341 An additional PCA with both river and soil samples (data not shown) showed that
 342 none of the first three components (88% of combined variance) separated these samples
 343 by type, but rather soil and river samples cluster together based on the location (Go-
 344 davari or Carminowe).

345 Finally, for the marine sediments, primarily classified as cluster 4 (Fig. 4e), the first
 346 component of the PCA (52.6% of variability) is positively associated with GDGT Ia and
 347 negatively associated with crenarchaeol. The second component (37.2%) is positively as-
 348 sociated with GDGT-0. The first component is associated with the sample's distance to
 349 the coast ($\rho = 0.43$ and Fig. SI1e), while the second component shows a strong cor-
 350 relation with the mixed layer temperature ($\rho = -0.79$). We, however, do not find a strong

351 association between samples classified as cluster 2 or 3 and environmental parameters
352 considered here.

353 Following the aim of this study, we considered the description of each sample, the
354 effects of environmental parameters in it, as well as the K-means classification. With this
355 information we generated new classification labels, which we prioritized in cases where
356 the fuzzy K-means classification and these labels disagreed. The informed labels are named
357 according to the dominant depositional environment. Group 1 was renamed as *Peat-type*,
358 Group 2 as *Lake-type*, Group 3 as *Soil-type*, and finally Group 4 as *Marine-type* (Fig. 2c).
359 While *Soil-type* combines samples from rivers and soils, we name it as such since we have
360 a much larger representation of soils in the dataset, but also mechanistically it is more
361 likely that GDGTs from the soils are influencing nearby rivers rather than the other way
362 around.

363 3.3 Supervised Machine Learning

364 The new informed labels generated through the unsupervised machine learning phase
365 were used for the supervised classification. We tested the performance of all four clas-
366 sification algorithms against each other and compared them with the null model using
367 both the F1 and ROC-AUC parameters. Our results suggest that overall, all methods
368 performed significantly better than the noninformative control and relatively similar to
369 each other. For the F1 scores, Random Forest performed the best (0.95), followed by XG-
370 Boost (0.94), K-Nearest Neighbour (0.91), and Naive Bayes (0.87). In contrast, the null
371 model had a score of 0.58. Similarly, for the ROC-AUC parameter Random Forest, XG-
372 Boost, and K-Nearest Neighbour had the same performance (0.99), followed by Naive
373 Bayes (0.96), and the null model had a value of only 0.5. Based on these results we chose
374 the Random Forest algorithm for our study. The performance of this algorithm in the
375 test set is similar to the one obtained for the training set (0.94 and 0.99 for F1 and ROC-
376 AUC respectively, Fig. 5).

377 Finally, we diagnose the importance that each predictor variable has on the trained
378 classification algorithm. This analysis shows that brGDGT IIa' and crenarchaeol have
379 the highest importance scores (> 90), followed by IIb', IIIa', IIIb, Ia, and crenarchaeol'
380 (> 30). All other variables had importance values < 30 (Fig. 5b).

381 The finalized model, named **B**ranching and **I**soprenoid **G**DGT **M**achine learning
382 **C**lassification algorithm (BIGMaC), is available on Github as an R object (Martínez-
383 Sosa et al., 2023).

384 4 Discussion

385 4.1 Unsupervised Machine Learning

386 While our fuzzy k-means clusters show strong patterns that reflect relationships
387 with depositional environments (Fig. 2a), some samples whose context was unequivocally
388 documented cluster in with samples unrelated to their depositional environment
389 (*i.e.* soils plotting as peats).

390 For both peats and lakes, temperature has a strong influence on the GDGT dis-
391 tribution, causing samples at either ends to be clustered with other sample types (Fig.
392 4 a, b, and Fig. S11 a and b). At high temperatures, samples from these depositional
393 environments tend to have a higher proportion of brGDGT Ia, and lower proportion of
394 IIa and IIIa (Weijers et al., 2007), which causes the lake sediments to be classified as clus-
395 ter 1 and 3 (high Ia), or peats as cluster 3 and 4 (low IIa and IIIa). At lower temper-
396 atures, the opposite effect causes some of the peats to be classified as cluster 2. In ad-
397 dition, particularly deep (*i.e.* Lake Malawi) or hypersaline alkaline lakes also showed an
398 increased proportion of isoGDGTs compared with other lakes, causing them to be clas-

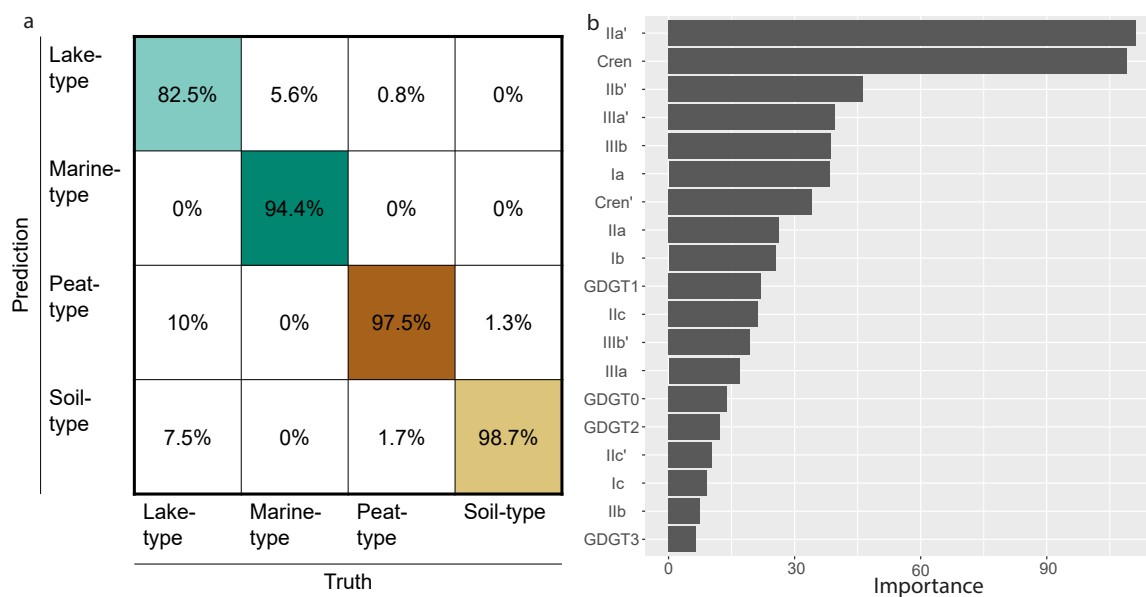


Figure 5. Confusion matrix and importance value for all the GDGTs considered in the classification model. The confusion matrix (a) shows the performance of the BIGMaC Random Forest algorithm in the test dataset. Columns show the true label of the samples and rows show the predicted label. Diagonal cells are color-coded based on Fig. 2. The bar plot (b) shows the importance value of each GDGT, calculated based on Gini impurity, used in the BIGMaC classification algorithm.

399 sified with marine samples. These distinct distributions could be due to specific niches
 400 in the water column associated with water chemistry, stratification, and/or nutrient con-
 401 tent, as previous work has suggested (Sinninghe Damsté et al., 2022; Baxter et al., 2021;
 402 Kumar et al., 2019).

403 Soil samples cluster into several different groups, which may reflect the fact that
 404 soils are highly diverse environments with diffuse boundaries and are often in contact
 405 with other depositional settings. Furthermore, studies have shown that chemical prop-
 406 erties of soils (i.e. pH, metal concentrations) have great spatial heterogeneity even at small
 407 scales (Yavitt et al., 2009). Indeed, the main factor that separated soils in this study is
 408 their location (Fig. 4c and Fig. SI1c), which overrides any other environmental signal.
 409 However, this may simply be a feature of the limited existence of soil datasets with both
 410 branched and isoprenoid GDGTs (represented by only four locations). It is clear from
 411 our results that soils require a more in-depth analysis, with the use of more extensive
 412 datasets.

413 While there is some debate regarding the relative influence that soil input and in
 414 situ production have on the GDGT pool in river organic matter (Kirkels et al., 2020; Zell
 415 et al., 2013; De Jonge, Stadnitskaia, et al., 2014), our analysis shows that the river sed-
 416 iments more closely resemble soils rather than peats or lake sediments, and similarly to
 417 soils, the sample location shows the strongest correlation with how these samples are clas-
 418 sified by fuzzy k-means (Fig. 4d and Fig. SI1d). We do find a difference in the 5-methyl/6-
 419 methyl proportion in these groups, where soils have a relatively similar proportion of these
 420 isomers, while river sediments contain relatively more 6-methyl brGDGTs (Fig. SI2 b
 421 and c). Although this could be interpreted as soil-derived GDGTs dominating river in-
 422 puts, with some autochthonous production of 6-methyl brGDGTs in rivers, our river data
 423 come from only two locations and primarily from only one system (the Godavari river),

424 so this interpretation could be particular to this watershed. Notably, within the Godavari
 425 River, the membership value for the samples, which measures the degree of belonging
 426 to each cluster, varies with their location and collection season (Fig. 6). Membership
 427 to the soil-dominated Group 3 is higher in the lower Godavari basin, as well as from the
 428 wet (post-monsoon) season (Fig. 6 c and d). In contrast, membership to the lake-dominated
 429 Group 2 is overall higher in the dry season, and in the upper basin year-round (Fig. 6
 430 a and b). These results are in line with those presented in the original study by Kirkels,
 431 Zwart, et al. (2022), where it was noted that GDGTs from soils have a stronger influ-
 432 ence on the river's GDGT content during the wet season and within the lower basin, which
 433 experiences higher precipitation. In contrast, in-situ production of brGDGTs, charac-
 434 terized by a high proportion of 6-methyl isomers, has a stronger influence on the GDGT
 435 content of samples from the dry season as well as those from the upper basin.

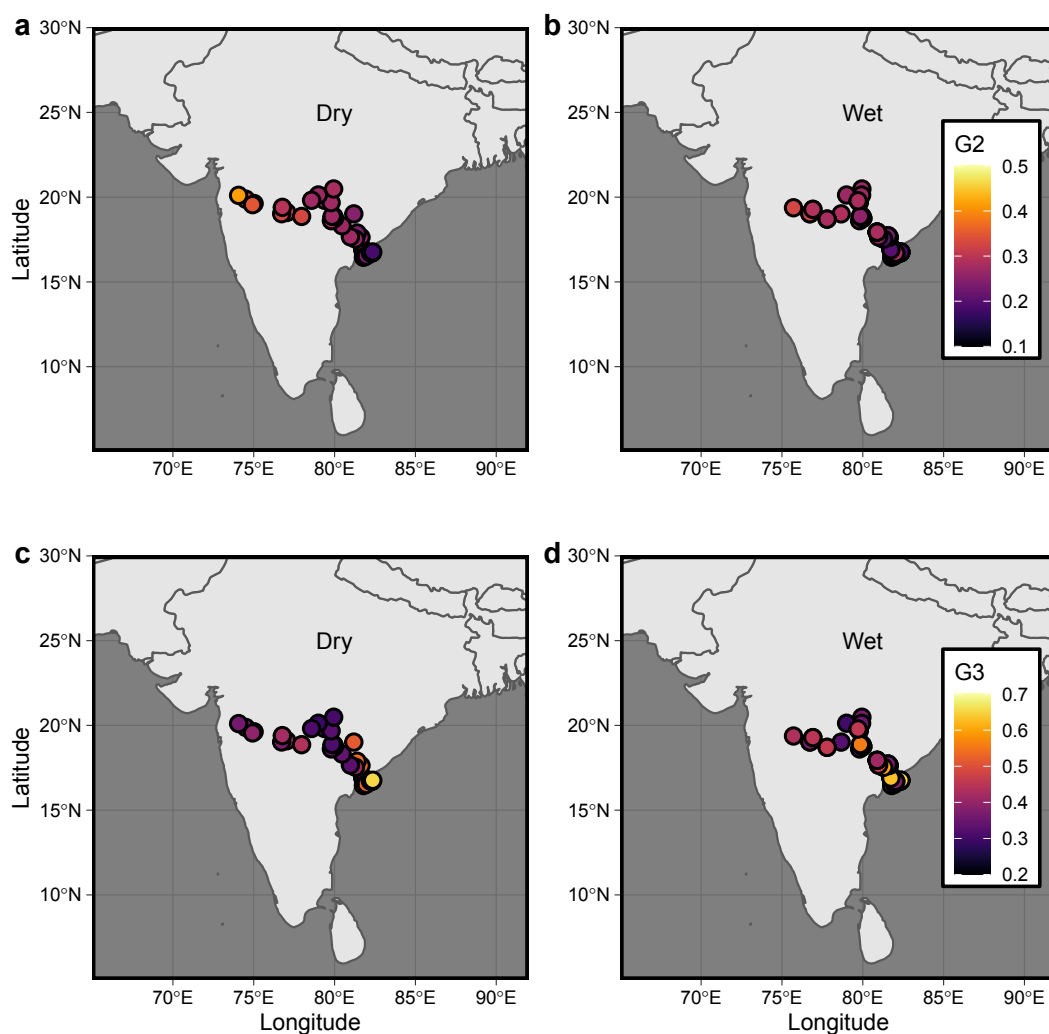


Figure 6. Maps for the Godavari River sample locations for the dry (left column) and wet (right column) seasons. Maps show the sample memberships, calculated through fuzzy k-means analysis, to the lake-dominated Group 2 (a and b), and to the soil-dominated Group 3 (c and d).

436 Notably, for marine sediments the first dimension of their PCA is dominated by
 437 a positive relation with brGDGT Ia and a negative one with crenarchaeol and it is as-

sociated with distance from the coast, although this is true only for samples more than 10 km away from the nearest coast, suggesting that more coastal samples may be affected by other GDGT sources (Fig. 4e and Fig. S11e). The second dimension, which is positively related to GDGT-0, more closely follows the mixed layer temperature, as shown by the Spearman correlation between this variable and the principal component ($\rho = -0.79$). Although GDGT-0 is traditionally omitted from the TEX_{86} calculation because it is a generic isoGDGT produced by many types of archaea (including methanotrophs and methanogens) (Kim et al., 2010; Schouten et al., 2002), our analysis shows that it is strongly influenced by temperature. Furthermore, the PCA shows no relation between GDGT-0 and brGDGTs (Fig. 4e), which suggests that GDGT-0 is not influenced by terrestrial sources. Our results suggest that temperature strongly influences the abundance of GDGT-0 and, unlike previously thought (Guo, Yuan, et al., 2022; Kim et al., 2010), other environmental parameters may not be as important in open marine settings. This supports the observation of Cramwinckel et al. (2018) that, at higher temperatures, the ratio of crenarchaeol to GDGT-0 might be more sensitive to temperature changes than TEX_{86} .

Since our intention with the supervised machine learning was to test whether GDGT distributions can be used to identify the depositional environment, we generated four new groups which broadly follow the fuzzy k-means clusters, but considered the actual depositional environmental that the samples came from (Fig. 2 a to c). For example, although some of the peat samples fell into cluster 2 (lakes) or cluster 4 (marine), since they were derived from peatlands, we re-assigned them as such. As noted above, some of these false assignments appear to be associated with the effect of temperature on GDGT distributions in peats. Similarly, although soils fell into several different clusters, this seems to be because their GDGT distributions were influenced by the site location, so we assigned them all to a single soil group. By manually reassigning samples to match their true environment, we feed the classification algorithm a more realistic dataset that includes some of the uncertainties associated with the relationship between GDGT distributions and their depositional settings. While we cannot rule out that the effects observed here are due to artifacts of the clustering technique used, as we only tested the performance of fuzzy k-means based on Euclidian distances, which can be affected by high dimensionalities, the environmentally relevant results obtained give us confidence in the approach used.

4.2 Curated clusters

The manual curation of these clusters does not alter the general composition of the groups compared with the statistically-derived ones (Table 1). *Peat-type* and *Marine-type* are very similar in composition and size to Group 1 and 4 respectively. While Group 1, with 465 samples, had 87% of the peats and 20% of the soils; *Peat-type*, with 476 samples, has all of the peats and only one lake sediment. Similarly, Group 4, with 225 samples, had 92% of the marine sediments, while *Marine-type* includes all of them and has a total of 216 samples. The reduction in size from Group 4 to *Marine-type* is mostly due to the reassignment of lake sediments, peats and soils. The largest change observed is between Group 2 and *Lake-type* (86 sample difference), and Group 3 and *Soil-type* (84 sample difference). Most of this comes from the reassignment of 60 soils from Group 2 to *Soil-type*.

The sample reassignment also does not alter the general GDGT distribution patterns of the groups when we compare them before and after the reassignment (Fig. 3). It does, however, preserve the cluster selection for samples where the depositional environments may not be as clear (i.e., humic-rich lakes which are classified as *Peat-type*).

487

4.3 Supervised Classification

488

489

490

491

492

493

494

495

496

In general, all of the machine learning algorithms exhibited good performance in the training phase, with F1 and ROC-AUC scores above 0.85 and 0.95 respectively. Nevertheless, we chose the Random Forest algorithm since it was the best performing one across all parameters, in addition to being widely used in the field of Geosciences (People et al., 2021; El Bouchefry & de Souza, 2020). This algorithm also performed well in the testing phase (0.94 and 0.99, for F1 and ROC-AUC respectively, and Fig. 5). This result suggests that the algorithm is not overfitting the data, as the algorithm’s performance would have significantly decreased if it had been trained to only classify samples it had previously been exposed to.

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

When we apply the BIGMaC algorithm to the complete dataset, we can investigate the importance of each GDGT in the model. This analysis shows that the two compounds that contribute the most to the classification are I Ia’ and crenarchaeol, although all GDGTs contribute to some extent to the classification, and thus relying only on these two compounds may not be as effective as the complete model. While these compounds have not been specifically associated with particular environments, they have very specific distributions in the identified clusters (Fig. 3), and also could be associated with characteristics of the depositional environments. BrGDGT I Ia’ is the 6-methyl GDGT with the highest abundance in lakes and soils, however, these isomers are less abundant in peats although this could also be due to most of the peats in our dataset being acidic (Naafs et al., 2017) and the inclusion of less acidic samples could change the variables used. In addition, generally all brGDGTs have a lower abundance in marine environments (Hopmans et al., 2004). In contrast, crenarchaeol is generally the most abundant isoGDGT in marine sediments (Liu et al., 2011), and it has been shown that it is also present in relatively high proportion in soils, but not in lakes (Naeher et al., 2014), and it has generally a low proportion in peats (Naafs et al., 2017). These specific distributions form four distinct patterns, which likely explain their selection as the most informative variables: low I Ia’—low crenarchaeol (*Peat-type*), high I Ia’—low crenarchaeol (*Lake-type*), low I Ia’—high crenarchaeol (*Marine-type*), and high I Ia’—high crenarchaeol (*Soil-type*).

516

4.4 Applications

517

518

519

520

521

We demonstrate that our model can be successfully used to analyze changes in depositional environments through time by testing the BIGMaC algorithm on GDGTs measured in two different sites: the Eocene-aged post-eruption peat and lacustrine sediments recovered from the Giraffe kimberlite pipe in the subarctic; and the Cobham lignite bed, dated to the beginning of the PETM.

522

4.4.1 Giraffe Kimberlite Pipe

523

524

525

526

527

528

529

530

531

532

533

534

535

536

When we apply the BIGMaC algorithm to the Giraffe kimberlite pipe core we see that the samples are generally correctly classified with the general stratigraphy previously described for the core (Wolfe et al., 2017; Hamblin et al., 2003) (Fig. 7 and GDGT distributions shown in Fig. SI4). All samples from the top peatland section are classified as *Peat-type*, and all samples from the lacustrine section below 85 m are classified as *Lake-type*. However, we also identified a section, between 76.5 and 85 m, within the lacustrine facies that is classified as *Peat-type*. Furthermore, the samples immediately above the excursion oscillate between *Lake-type* and *Soil-type* for at least one meter (Fig. 7), before returning to being classified as *Lake-type* for the last 4.8 m of the lake section. The results from our classification algorithm are in contrast to other approaches previously used to determine the origin of GDGTs in sediments. For example, applying the BIT index to this core shows that for the majority of the core, the mean value is 0.999 ± 0.001 , suggesting a terrestrial setting. The BIT index record only deviates from these values in the one-meter section where the BIGMaC classifies samples as either *Lake-type*

537 or *Soil-type*, although even in this section the BIT index values are only reduced to 0.97.
 538 Since the BIT index is unable to distinguish between soil, peat and lake deposits, this
 539 showcases the advantage of using the classification algorithm where all GDGTs are con-
 540 sidered over a singular index.

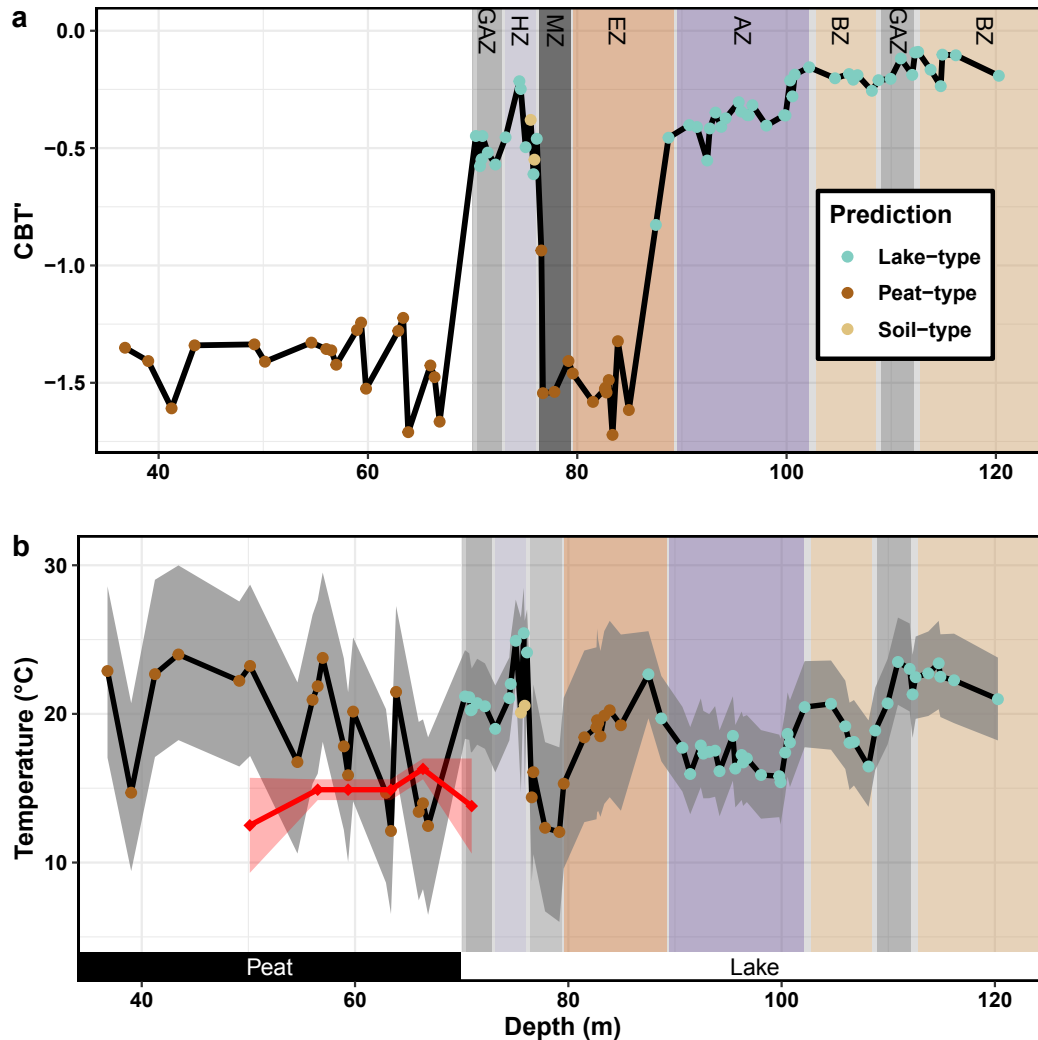


Figure 7. CBT' values (a) and inferred temperature (b) calculated from Giraffe core brGDGTs plotted against vertical-equivalent depth. The temperature reconstruction was generated by applying the Dearing Crampton-Flood et al. (2020) Bayesian calibration for *Peat* and *Soil-type* samples, and Martínez-Sosa et al. (2021) calibration for *Lake-type* samples. Palynological estimates of mean annual temperature with their associated error from Wolfe et al. (2017) are shown in red diamonds in (b). Samples are color-coded based on the predicted groups. White and gray shading indicates peat and lacustrine sediments in the core, respectively. Zones in the lake section as described by Siver and Lott (2023) are shown as colored areas: *Botryococcus* Zone (BZ) in light orange, Going Acidic Zone (GAZ) in light gray, *Aulacoseira* Zone (AZ) in purple, Eunotiid Zone (EZ) in orange, Mixing Zone (MZ) in gray, and Heterotrophic Zone (HZ) in light purple. The peat and lake sections of the core are indicated at the bottom of the plot with black and white rectangles, respectively

541 By using the CBT' index, which has been shown to be strongly associated with pH
 542 in peats (Naafs et al., 2017) and mildly correlated to pH in lakes (Martínez-Sosa et al.,
 543 2021), we estimate that in general the peat section has much lower CBT' values (asso-
 544 ciated with lower pH) than those in the lacustrine section. While this trend is maintained
 545 for most of the core, a marked decrease in CBT' values occurs in the section within the
 546 lacustrine facies that is classified as *Peat-type*.

547 Our results show a close relationship with the independently developed microfossil-
 548 based ecological reconstruction done by Siver and Lott (2023) for this site. While Hamblin
 549 et al. (2003) had previously speculated this site to be a shallow lacustrine setting with
 550 intermittent wet and dry periods, the microfossil ensemble suggest that the lake was ini-
 551 tially a shallow (~ 1 m) slightly acidic lake, this is represented by the *Botryococcus* Zone
 552 (BZ) between 102.8 and 124.6 m (Fig. 7). All samples in this section are classified as *Lake-*
 553 *type* with generally stable CBT' values. While Siver and Lott (2023) reported a shift in
 554 the microfossil ensemble within this region between 112.1 and 109 m, which was inter-
 555 preted as an acidifying section, we do not find evidence of this in the CBT' values; how-
 556 ever, this subsection does roughly align with an estimated reduction in MBT'_{5Me} val-
 557 ues. Following this initial section, the microfossils suggest that the lake deepened to be-
 558 tween 3 — 5 m in the *Aulacoseira* Zone (AZ) (Siver & Lott, 2023), between 102.1 and
 559 89.5m. BIGMaC still classifies all samples from this section as *Lake-type*, however there
 560 is a declining trend in CBT' values throughout this region, shifting from values of -0.2 to
 561 -0.55. The following section located between 89.2 to 79.6 m, identified by the microfos-
 562 sil ensemble as Eunotid Zone (EZ), is thought to represent a period characterized by
 563 enhanced acidity and increase in dissolved humic material. This section corresponds to
 564 the large CBT' excursion in our record, with values from -0.5 to below -1.5. Microfos-
 565 sils in this section are associated with low pH environments, such as lakes and bogs. This
 566 is in agreement with our BIGMaC results, where all but one sample in this section are
 567 classified as *Peat-type*, with the exception being a sample right at the bottom that BIG-
 568 MaC still classified as *Lake-type*. We note that the CBT' values are lower than those of
 569 previously reported lakes with a pH of 4.3 (Martínez-Sosa et al., 2021). In contrast, these
 570 values correlate to a pH between 4 and 5 in peats, which is more in line with what is as-
 571 sociated with the described organisms in this section (Siver & Lott, 2023). Between 79.4
 572 and 76.4 m the microfossil ensemble suggests the presence of a transitional zone that is
 573 not associated with any of the other sections. Our results classify all samples within this
 574 section as *Peat-type*, although the CBT' values have a steep increase, going from -1.5 to
 575 -0.5, which could suggest a transitory state towards a higher pH environment. The fol-
 576 lowing section, between 76 and 73.1 m, identified as the Heterotrophic Zone (HZ) by the
 577 microfossil ensemble, is interpreted as a period with higher pH and in line with condi-
 578 tions of an eutrophic body of water. Our analysis shows a transitory period at the be-
 579 ginning of this section where BIGMaC interprets two samples as *Soil-type*, while the rest
 580 of the samples are classified as *Lake-type*. The CBT' values are on par with those of the
 581 previous sections that were classified as *Lake-type*. Finally, the top-most part of the lake
 582 section is interpreted as a body of water with increased acidity and levels of dissolved
 583 humic matter. While the CBT' values do show a small decrease, the BIGMaC algorithm
 584 still classifies all samples in this section as *Lake-type*.

585 Overall, there is a good agreement in the interpretation of both independent proxy-
 586 based reconstructions for the lake section of the core. Differences between the interpre-
 587 tations could be due to constraints of the specific proxies. For example, while the mi-
 588 crofossil ensemble suggests periods of increased acidification throughout the BZ region,
 589 this is not reflected in the CBT' values. In contrast, while the microfossil ensemble sug-
 590 gests a sudden transition from the AZ to the EZ sections, our results show a decreas-
 591 ing trend in CBT' values throughout AZ leading to the much steeper acidification in EZ.
 592 It is possible that the decreases in pH in both cases have different origins, which are only
 593 captured by one of the proxies. This underscores the advantage of combining indepen-
 594 dent proxies for paleoenvironmental reconstructions.

595 To generate a temperature reconstruction for the Giraffe core, we applied the pre-
 596 viously published BayMBT calibration for lakes or soil/peat, depending on the results
 597 of BIGMaC for each sample. These particular calibrations were chosen as they are con-
 598 sistent with each other and allow us to generate a continuous confidence interval for the
 599 reconstruction. However, we emphasize that the choice of temperature calibration is in-
 600 dependent from BIGMaC and any GDGT-based calibration can be used. For example,
 601 we also applied the peat-specific calibration (Naafs, 2017) to the sections of the core clas-
 602 sified as *Peat-type* but it generated only marginally different results than BayMBT, so
 603 we chose to use the latter.

604 Our reconstruction suggests a relatively stable climate with no clear trend (Fig. 7a).
 605 The mean temperature of our reconstruction ($19 \pm 3.2^\circ\text{C}$) agrees with independent stud-
 606 ies. A pollen reconstruction at this site (red diamonds in Fig. 7a), suggests a MAAT of
 607 $14.5 \pm 1.3^\circ\text{C}$, with a warmest month mean temperature of $24.5 \pm 0.8^\circ\text{C}$ (Wolfe et al.,
 608 2017). In addition, Jahren and Sternberg (2003) estimated a mean annual temperature
 609 of $13.2 \pm 2^\circ\text{C}$ for the middle Eocene Arctic based on oxygen isotopes measured in cal-
 610 cite preserved in fossil *Metasequoia*. While our estimate is at the upper end of both es-
 611 timates, they fall within the confidence interval of our reconstruction (Fig. 7a). More-
 612 over, both the peat/soil and lake calibrations predict mean annual temperatures above
 613 freezing (MAF) rather than strictly MAAT, so if there were freezing temperatures dur-
 614 ing the winter, the GDGT estimates are expected to be higher. Conversely, if we had
 615 used only the lakes or soil/peat calibration for the entire core, there would be large tem-
 616 perature swings of more than 6°C associated with changes in core lithology. In partic-
 617 ular, the excursion to *Peat-type* samples within the lacustrine section would be estimated
 618 to be 5.7°C higher without the BIGMaC-based correction.

619 4.4.2 Cobham Lignite Bed

620 While the application of the BIGMaC algorithm in the Giraffe pipe showcases its
 621 strengths, our analysis of the Cobham lignite illustrates that there are some limitations
 622 of the approach. Inglis et al. (2019) previously showed that increased precipitation dur-
 623 ing the PETM in this area caused changes in the hydrology of the site, and that this po-
 624 tentially caused the brGDGTs to become unreliable as temperature proxies. Namely, while
 625 several lines of evidence suggest an increase in temperature during the PETM, the tem-
 626 perature reconstructions based on brGDGTs suggest cooling. We applied BIGMaC to
 627 this site to investigate whether changes in the depositional settings could explain the dis-
 628 crepancy (Fig. 8). Almost all samples preceding 54.15 cm, identified as the start of the
 629 PETM, are predicted to be *Peat-type*, with the exception of one sample from the upper
 630 laminated lignite unit that is classified as *Soil-type*. In contrast, there is a wider vari-
 631 ation in the sample classification during the PETM, where samples are classified as *Peat-*
 632 *type* (10), *Soil-type* (3) and *Lake-type* (1). Besides one sample classified as *Peat-type* from
 633 the PETM upper laminated lignite, all other PETM samples are located in the blocky
 634 lignite unit. The variations in predicted depositional environments do not coincide with
 635 changes in either MBT'_{5Me} or CBT' values, nor are they organized in any evident pat-
 636 tern within the unit. Moreover, the PETM samples are primarily classified as *Peat-type*
 637 and *Soil-type*, suggesting that adjusting for the predicted depositional environment would
 638 not significantly increase the temperature reconstructed by Inglis et al. (2019), as would
 639 be the case if the samples were classified as *Lake-type*. Vegetation and charcoal records
 640 suggest that the Cobham site became waterlogged and may have even developed areas
 641 of open water during the PETM (Inglis et al., 2019). Given this perspective, the oscil-
 642 lating results from BIGMaC likely point to an unstable, dynamically changing deposi-
 643 tional environment with mixed sources of brGDGTs. Since BIGMaC is a categorical clas-
 644 sification algorithm, it cannot detect mixed signatures. This underlines the need to in-
 645 corporate mixing models in studies where input from different sources is expected, and
 646 suggests that BIGMaC would benefit from incorporating this capability in future updates.

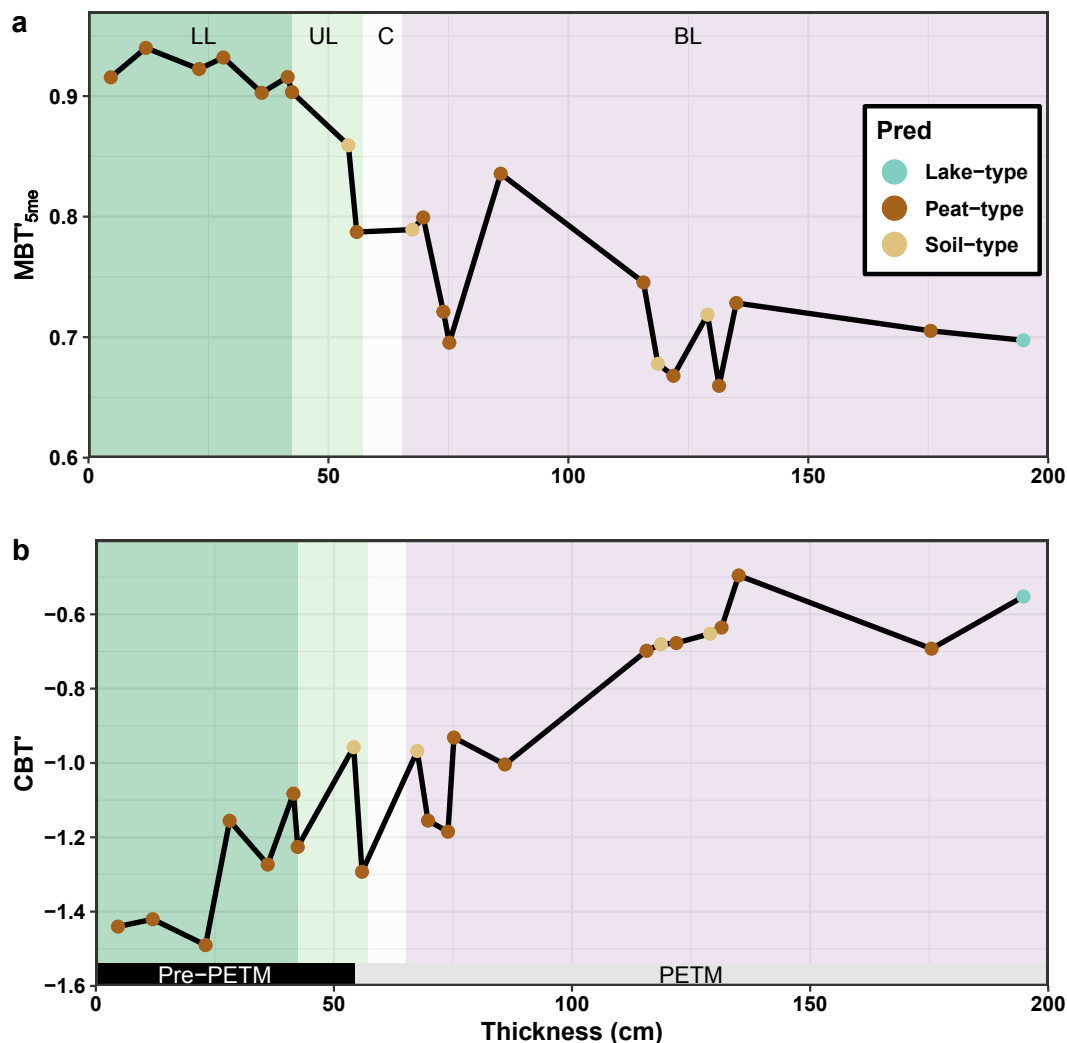


Figure 8. Calculated MBT'_{5Me} (a) and CBT' (b) values of the Cobham lignite bed across the site thickness (cm). Samples are color coded based on the BIGMaC predicted groups. Different units are colored and labeled on the top as: lower laminated lignite (LL, dark green), upper laminated lignite (UL, light green), clay (C, white), and blocky lignite (BL, purple).

647 Overall, this application suggests that other environmental parameters, not con-
 648 sidered in this work, could still affect the distribution of GDGTs, and we include the Cob-
 649 ham lignite example as a reminder that applying the BIGMaC algorithm is not always
 650 a panacea and is best done in concert with other independent proxies to accurately in-
 651 terpret the results.

652 5 Conclusions

653 Our analyses of GDGTs in 1153 globally distributed samples from soils, lakes, rivers,
 654 and marine sediments show that the depositional environment from which samples were
 655 obtained has a significant and measurable impact on the combined distribution of iso-
 656 prenyl and branched GDGTs, which allows us to cluster the samples from our dataset
 657 into environmentally relevant groups. Furthermore, we find that the distribution of GDGTs

658 in each cluster is uniquely impacted by the environment. There is a strong association
659 between temperature and the *Lake-type* and *Peat-type* groups. *Marine-type* samples are
660 also clearly influenced by temperature and their distance from the coast, at least for sam-
661 ples more than 10km away from the shore. The cause of this distinction is an observa-
662 tion that deserves further study. Although they are represented by a limited sample set,
663 soils show variability that is strongly influenced by their location. Additionally, while
664 our analysis groups soil and river samples together into the *Soil-type* cluster, river sys-
665 tems seem to have more 6-methyl brGDGTs and their GDGT distributions reflect lo-
666 cal changes within the catchment.

667 We used the dataset presented here to train the Random Forest classification al-
668 gorithm BIGMaC, which is capable of identifying the environment in which a sample was
669 formed based on the distribution of GDGTs. Our results show that GDGTs Iia' and cre-
670 narchaeol are the most influential compounds in the classification algorithm, due to their
671 combined unique changes between the four depositional groups. As a demonstration, we
672 apply the BIGMaC model to an independent record from the Giraffe kimberlite, which
673 was stratigraphically shown to record a transition from a lacustrine environment to peat-
674 land. Our BIGMaC algorithm is not only able to recreate the transition, but further sug-
675 gests an excursion to peatland conditions within the upper lacustrine section of the core,
676 which is consistent with independent evidence for more acidic conditions. This result is
677 encouraging for the application of our classification algorithm, as it comes from a dataset
678 not included in the training or testing sets, thus providing an independent testing case.
679 Using the BIGMaC results as a guide, we apply brGDGT-derived calibrations specific
680 to lakes or soils and peats as needed downcore and obtain a relatively stable tempera-
681 ture estimate for this area that is in general agreement with the pollen record.

682 While our Giraffe pipe results showcase the usefulness of our approach when ap-
683 plied to clear changes in depositional environments; the application of BIGMaC in the
684 Cobham site shows that this approach may not be suitable in cases where the deposi-
685 tional environment is changing rapidly and thereby results in mixed sources of GDGTs.
686 It is possible that the future integration of a mixing model in the BIGMaC workflow could
687 improve its performance in this type of scenario.

688 Ultimately, we show that the combined set of branched and isoprenoid GDGTs is
689 an effective tool for identifying depositional environments that can be used in combina-
690 tion with more established proxies to gain a better understanding of past environments.

691 **Open Research**

692 The GDGT fractional abundance data used for training the BIGMaC algorithm
693 in the study are directly available at Pangaea via Naafs (2017), Guo et al. (2020), Dearing
694 Crampton-Flood et al. (2019), Guo et al. (2021), and Inglis et al. (2019); as well as on
695 Zenodo via Martínez-Sosa et al. (2023b), Martínez-Sosa et al. (2023a) and Pérez-Angel
696 et al. (2020). V1.0 of the BIGMaC algorithm used for the classification of samples based
697 on GDGT fractional abundances is preserved at Martínez-Sosa et al. (2023), available
698 via MIT license and developed openly in the `tidymodels` environment in R.

699 **Acknowledgments**

700 We would like to thank Patrick Murphy for his assistance with the lipid analysis,
701 Dr. Jeffrey Donnelly and the Woods Hole Oceanographic Institution Seafloor Samples
702 Laboratory for access to marine sediment samples, and Dr. Cody Routson for contribut-
703 ing Alaskan lake samples. The hyperparameter tuning of the models was performed us-
704 ing the Ocelote cluster from the University of Arizona. This research was funded by the
705 American Chemical Society Petroleum Research Fund, grant 60772-ND2, and by CONA-
706 CYT through the student scholarship 440897. Ioana Stefanescu and Bryan Shuman ac-

707 knowledge support from the Microbial Ecology Collaborative Project through the Na-
708 tional Science Foundation grant EPS-1655726. Francien Peterse acknowledges funding
709 from the Nederlandse Organisatie voor Wetenschappelijk Onderzoek (NWO) through
710 Veni grant no. 863.13.016 and Vidi grant no. 192.074. Lina Pérez-Ángel and Julio Sepúlveda
711 acknowledge support from NSF Sedimentary Geology and Paleobiology grant 1929199.
712 We also thank Serhiy Buryak for assisting with the sampling of the Giraffe pipe sediments.

References

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

- Baxter, A., van Bree, L., Peterse, F., Hopmans, E., Villanueva, L., Verschuren, D., & Sinninghe Damsté, J. S. (2021). Seasonal and multi-annual variation in the abundance of isoprenoid GDGT membrane lipids and their producers in the water column of a meromictic equatorial crater lake (Lake Chala, East Africa). *Quaternary Science Reviews*, *273*, 107263.
- Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785–794).
- Chen, Y., Zheng, F., Yang, H., Yang, W., Wu, R., Liu, X., . . . others (2022). The production of diverse brGDGTs by an Acidobacterium providing a physiological basis for paleoclimate proxies. *Geochimica et Cosmochimica Acta*, *337*, 155–165.
- Collinson, M. E., Steart, D. C., Harrington, G. J., Hooker, J. J., Scott, A. C., Allen, L. O., . . . Gibbons, S. J. (2009). Palynological evidence of vegetation dynamics in response to palaeoenvironmental change across the onset of the Paleocene-Eocene Thermal Maximum at Cobham, Southern England. *Grana*, *48*(1), 38–66.
- Cramwinckel, M. J., Huber, M., Kocken, I. J., Agnini, C., Bijl, P. K., Bohaty, S. M., . . . others (2018). Synchronous tropical and polar temperature evolution in the Eocene. *Nature*, *559*(7714), 382–386.
- Dang, X., Ding, W., Yang, H., Pancost, R. D., Naafs, B. D. A., Xue, J., . . . Xie, S. (2018, May). Different temperature dependence of the bacterial brGDGT isomers in 35 Chinese lake sediments compared to that in soils. *Org. Geochem.*, *119*, 72–79.
- Dearing Crampton-Flood, E., Tierney, J. E., Peterse, F., Kirkels, F. M. S. A., & Sinninghe Damsté, J. S. (2019). *Global soil and peat branched GDGT compilation dataset* [Dataset]. PANGAEA. doi: 10.1594/PANGAEA.907818
- Dearing Crampton-Flood, E., Tierney, J. E., Peterse, F., Kirkels, F. M., & Sinninghe Damsté, J. S. (2020). BayMBT: A Bayesian calibration model for branched glycerol dialkyl glycerol tetraethers in soils and peats. *Geochimica et Cosmochimica Acta*, *268*, 142–159.
- De Jonge, C., Hopmans, E. C., Zell, C. I., Kim, J.-H., Schouten, S., & Damsté, J. S. S. (2014). Occurrence and abundance of 6-methyl branched glycerol dialkyl glycerol tetraethers in soils: Implications for palaeoclimate reconstruction. *Geochimica et Cosmochimica Acta*, *141*, 97–112.
- De Jonge, C., Radujković, D., Sigurdsson, B. D., Weedon, J. T., Janssens, I., & Peterse, F. (2019). Lipid biomarker temperature proxy responds to abrupt shift in the bacterial community composition in geothermally heated soils. *Organic Geochemistry*, *137*, 103897.
- De Jonge, C., Stadnitskaia, A., Hopmans, E. C., Cherkashov, G., Fedotov, A., & Sinninghe Damsté, J. S. (2014). In situ produced branched glycerol dialkyl glycerol tetraethers in suspended particulate matter from the Yenisei River, Eastern Siberia. *Geochim. Cosmochim. Acta*, *125*, 476–491.
- de Rosa, M., de Rosa, S., Gambacorta, A., Minale, L., & Bu'lock, J. D. (1977). Chemical structure of the ether lipids of thermophilic acidophilic bacteria of the Caldariella group. *Phytochemistry*, *16*(12), 1961–1965.
- De Rosa, M., Gambacorta, A., Nicolaus, B., Chappe, B., & Albrecht, P. (1983). Isoprenoid ethers; backbone of complex lipids of the archaeobacterium *Sulfolobus solfataricus*. *Biochimica et Biophysica Acta (BBA)-Lipids and Lipid Metabolism*, *753*(2), 249–256.
- Dunkley Jones, T., Eley, Y. L., Thomson, W., Greene, S. E., Mandel, I., Edgar, K., & Bendle, J. A. (2020). OPTiMAL: a new machine learning approach for GDGT-based palaeothermometry. *Climate of the Past*, *16*(6), 2599–2617. Retrieved from <https://cp.copernicus.org/articles/16/2599/2020/> doi:

- 10.5194/cp-16-2599-2020
- 768 El Bouchefry, K., & de Souza, R. S. (2020). Learning in big data: Introduction
769 to machine learning. In *Knowledge discovery in big data from astronomy and*
770 *earth observation* (pp. 225–249). Elsevier.
- 771
- 772 Engle, M. A., & Brunner, B. (2019). Considerations in the application of machine
773 learning to aqueous geochemistry: Origin of produced waters in the northern
774 US Gulf Coast Basin. *Applied Computing and Geosciences*, *3*, 100012.
- 775 Fleming, L. E., & Tierney, J. E. (2016). An automated method for the determina-
776 tion of the TEX_{86} and paleotemperature indices. *Org. Geochem.*, *92*, 84–91.
- 777 Greenwell, B., Boehmke, B., & Gray, B. (2020). Package ‘vip’. *Variable Importance*
778 *Plots*, *12*(1), 343–66.
- 779 Guo, J., Glendell, M., Meersmans, J., Kirkels, F., Middelburg, J. J., & Peterse, F.
780 (2020). Assessing branched tetraether lipids as tracers of soil organic car-
781 bon transport through the Carminowe Creek catchment (southwest England).
782 *Biogeosciences*, *17*(12), 3183–3201.
- 783 Guo, J., Glendell, M., Meersmans, J., Kirkels, F. M. S. A., Middelburg, J. J., & Pe-
784 terse, F. (2020). *Branched tetraether lipids in Carminowe Creek catchment*
785 *(southwest England)* [Dataset]. PANGAEA. doi: 10.1594/PANGAEA.918523
- 786 Guo, J., Ma, T., Liu, N., Zhang, X., Hu, H., Ma, W., . . . Peterse, F. (2021).
787 *Branched tetraether lipids and bacterial communities along an aridity soil*
788 *transect in Inner Mongolia, northern China* [Dataset]. PANGAEA. doi:
789 10.1594/PANGAEA.938067
- 790 Guo, J., Ma, T., Liu, N., Zhang, X., Hu, H., Ma, W., . . . Peterse, F. (2022). Soil pH
791 and aridity influence distributions of branched tetraether lipids in grassland
792 soils along an aridity transect. *Organic Geochemistry*, 104347.
- 793 Guo, J., Yuan, H., Song, J., Li, X., Duan, L., Li, N., & Wang, Y. (2022). Influ-
794 ence of bottom seawater oxygen on archaeal tetraether lipids in sediments:
795 Implications for archaeal lipid-based proxies. *Marine Chemistry*, 104138.
- 796 Halamka, T. A., McFarlin, J. M., Younkin, A. D., Depoy, J., Dildar, J., & Kopf,
797 S. H. (2021). Oxygen limitation can trigger the production of branched
798 GDGTs in culture. *Geochemical Perspectives Letters*, *19*, 36 – 39.
- 799 Halamka, T. A., Raberg, J. H., McFarlin, J. M., Younkin, A. D., Mulligan, C., Liu,
800 X.-L., & Kopf, S. H. (2022). Production of diverse brGDGTs by *Acidobac-*
801 *terium Solibacter usitatus* in response to temperature, pH, and O_2 provides a
802 culturing perspective on br GDGT proxies and biosynthesis. *Geobiology*.
- 803 Hamblin, A., Stasiuk, L., Sweet, A., Lockhart, G., Dyck, D., Jagger, K., & Snow-
804 don, L. (2003). Post-kimberlite Eocene strata within a crater basin, Lac de
805 Gras, Northwest Territories, Canada. In *International kimberlite conference:*
806 *Extended abstracts* (Vol. 8).
- 807 Hopmans, E. C., Schouten, S., & Damsté, J. S. S. (2016). The effect of improved
808 chromatography on GDGT-based palaeoproxies. *Organic Geochemistry*, *93*, 1–
809 6.
- 810 Hopmans, E. C., Weijers, J. W., Schefuß, E., Herfort, L., Damsté, J. S. S., &
811 Schouten, S. (2004). A novel proxy for terrestrial organic matter in sedi-
812 ments based on branched and isoprenoid tetraether lipids. *Earth and Planetary*
813 *Science Letters*, *224* (1-2), 107–116.
- 814 Huguet, C., Hopmans, E. C., Febo-Ayala, W., Thompson, D. H., Sinninghe Damsté,
815 J. S., & Schouten, S. (2006). An improved method to determine the absolute
816 abundance of glycerol dibiphytanyl glycerol tetraether lipids. *Org. Geochem.*,
817 *37*(9), 1036–1041.
- 818 Inglis, G. N., Farnsworth, A., Collinson, M. E., Carmichael, M. J., Naafs, B. D. A.,
819 Lunt, D. J., . . . Pancost, R. D. (2019). Terrestrial environmental change across
820 the onset of the PETM and the associated impact on biomarker proxies: A
821 cautionary tale. *Global and Planetary Change*, *181*, 102991.
- 822 Inglis, G. N., Farnsworth, A., Collinson, M. E., Carmichael, M. J., Naafs, B. D. A.,

- 823 Lunt, D. J., ... Pancost, R. D. (2019). *Terrestrial environmental change*
824 *across the onset of the PETM and the associated impact on biomarker proxies:*
825 *a cautionary tale* [Dataset]. PANGAEA. Retrieved from [https://doi.org/](https://doi.org/10.1594/PANGAEA.901285)
826 [10.1594/PANGAEA.901285](https://doi.org/10.1594/PANGAEA.901285) doi: 10.1594/PANGAEA.901285
- 827 Jahren, A. H., & Sternberg, L. S. L. (2003). Humidity estimate for the middle
828 Eocene Arctic rain forest. *Geology*, *31*(5), 463–466.
- 829 Kim, J.-H., Van der Meer, J., Schouten, S., Helmke, P., Willmott, V., Sangiorgi, F.,
830 ... Sinninghe Damsté, J. S. J. (2010). New indices and calibrations derived
831 from the distribution of crenarchaeal isoprenoid tetraether lipids: Implications
832 for past sea surface temperature reconstructions. *Geochimica et Cosmochimica*
833 *Acta*, *74*(16), 4639–4654.
- 834 Kirkels, F. M., Ponton, C., Galy, V., West, A. J., Feakins, S. J., & Peterse, F.
835 (2020). From Andes to Amazon: Assessing branched tetraether lipids as
836 tracers for soil organic carbon in the Madre de Dios River system. *Journal of*
837 *Geophysical Research: Biogeosciences*, *125*(1), e2019JG005270.
- 838 Kirkels, F. M., Usman, M. O., & Peterse, F. (2022). Distinct sources of bacte-
839 rial branched GMGTs in the Godavari River basin (India) and Bay of Bengal
840 sediments. *Organic Geochemistry*, *167*, 104405.
- 841 Kirkels, F. M., Zwart, H. M., Usman, M. O., Hou, S., Ponton, C., Giosan, L., ...
842 others (2022). From soil to sea: sources and transport of organic carbon traced
843 by tetraether lipids in the monsoonal godavari river, india. *Biogeosciences*,
844 *19*(17), 3979–4010.
- 845 Kuhn, M., & Wickham, H. (2020). Tidymodels: a collection of packages for mod-
846 eling and machine learning using tidyverse principles. [Computer software
847 manual]. Retrieved from <https://www.tidymodels.org>
- 848 Kumar, D. M., Woltering, M., Hopmans, E. C., Damste, J. S. S., Schouten, S., &
849 Werne, J. P. (2019). The vertical distribution of Thaumarchaeota in the water
850 column of Lake Malawi inferred from core and intact polar tetraether lipids.
851 *Organic Geochemistry*, *132*, 37–49.
- 852 Langworthy, T. A. (1977). Long-chain diglycerol tetraethers from Thermo-
853 plasma acidophilum. *Biochimica et Biophysica Acta (BBA)-Lipids and Lipid*
854 *Metabolism*, *487*(1), 37–50.
- 855 Lauretano, V., Kennedy-Asser, A. T., Korasidis, V. A., Wallace, M. W., Valdes,
856 P. J., Lunt, D. J., ... Naafs, B. D. A. (2021). Eocene to oligocene terrestrial
857 southern hemisphere cooling caused by declining p co2. *Nature Geoscience*,
858 *14*(9), 659–664.
- 859 Li, J., Pancost, R. D., Naafs, B. D. A., Yang, H., Zhao, C., & Xie, S. (2016). Distri-
860 bution of glycerol dialkyl glycerol tetraether (gdgt) lipids in a hypersaline lake
861 system. *Organic Geochemistry*, *99*, 113–124.
- 862 Liu, X., Lipp, J. S., & Hinrichs, K.-U. (2011). Distribution of intact and core
863 GDGTs in marine sediments. *Organic Geochemistry*, *42*(4), 368–375.
- 864 Locarnini, M., Mishonov, A., Baranova, O., Boyer, T., Zweng, M., Garcia, H., ...
865 others (2018). World ocean atlas 2018, volume 1: Temperature.
- 866 Maechler, M., et al. (2019). Finding groups in data”: Cluster analysis extended
867 Rousseeuw et al. *R package version*, *2*(0).
- 868 Martínez-Sosa, P., Tierney, J. E., & Meredith, L. K. (2020). Controlled lacustrine
869 microcosms show a brGDGT response to environmental perturbations. *Org.*
870 *Geochem.*, 104041.
- 871 Martínez-Sosa, P., Tierney, J. E., Stefanescu, I. C., Crampton-Flood, E. D., Shu-
872 man, B. N., & Routsom, C. (2021). A global Bayesian temperature calibration
873 for lacustrine brGDGTs. *Geochimica et Cosmochimica Acta*, *305*, 87–105.
- 874 Martínez-Sosa, P., Tierney, J., Pérez-Angel, L., Stefanescu, I. C., Guo, J., Kierkels,
875 F., ... Reyes, A. V. (2023). *BIGMaC GDGT algorithm* [Software]. Zenodo.
876 (V. 1.0) doi: 10.5281/zenodo.7513557
- 877 Martínez-Sosa, P., Tierney, J., Pérez-Angel, L., Stefanescu, I. C., Guo, J., Kirkels,

- 878 F., ... Reyes, A. V. (2023b). *Giraffe kimberlite pipe core GDGT fractional*
879 *abundance* [Dataset]. Zenodo. doi: 10.5281/zenodo.7540094
- 880 Martínez-Sosa, P., Tierney, J., Pérez-Angel, L. C., Stefanescu, I. C., Guo, J., Kirkels,
881 F., ... Reyes, A. V. (2023a). *Environmental data and fractional abundance of*
882 *iso and branched GDGT data used to train the BIGMaC algorithm* [Dataset].
883 Zenodo. doi: 10.5281/zenodo.7522415
- 884 Menze, B. H., Kelm, B. M., Masuch, R., Himmelreich, U., Bachert, P., Petrich, W.,
885 & Hamprecht, F. A. (2009). A comparison of random forest and its gini im-
886 portance with standard chemometric methods for the feature selection and
887 classification of spectral data. *BMC bioinformatics*, 10, 1–16.
- 888 Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., & Leisch, F. (2020). *e1071:*
889 *Misc Functions of the Department of Statistics, Probability Theory Group*
890 *(Formerly: E1071), TU Wien, 2018, R package version 1.7-0.*
- 891 Naafs, B., Rohrsen, M., Inglis, G. N., Lähteenoja, O., Feakins, S. J., Collinson,
892 M. E., ... others (2018). High temperatures in the terrestrial mid-latitudes
893 during the early Palaeogene. *Nature Geoscience*, 11(10), 766–771.
- 894 Naafs, B. D. A. (2017). *Global biomarker (GDGT) database for peatlands* [Dataset].
895 PANGAEA. Retrieved from <https://doi.org/10.1594/PANGAEA.883765>
896 doi: 10.1594/PANGAEA.883765
- 897 Naafs, B. D. A., Inglis, G. N., Zheng, Y., Amesbury, M., Biester, H., Bindler, R., ...
898 others (2017). Introducing global peat-specific temperature and pH calibra-
899 tions based on brGDGT bacterial lipids. *Geochimica et Cosmochimica Acta*,
900 208, 285–301.
- 901 Naeher, S., Peterse, F., Smittenberg, R. H., Niemann, H., Zigah, P. K., & Schubert,
902 C. J. (2014). Sources of glycerol dialkyl glycerol tetraethers (GDGTs) in
903 catchment soils, water column and sediments of Lake Rotsee (Switzerland)–
904 Implications for the application of GDGT-based proxies for lakes. *Organic*
905 *Geochemistry*, 66, 164–173.
- 906 Parmar, A., Katariya, R., & Patel, V. (2019). A review on random forest: An en-
907 semble classifier. In *International conference on intelligent data communication*
908 *technologies and internet of things (icici) 2018* (pp. 758–763).
- 909 People, M. D., Tierney, J. E., McGee, D., Lowenstein, T. K., Bhattacharya, T., &
910 Feakins, S. J. (2021). Identifying plant wax inputs in lake sediments using
911 machine learning. *Organic Geochemistry*, 156, 104222.
- 912 Pérez-Angel, L. C., Sepúlveda, J., Molnar, P., Montes, C., Rajagopalan, B., Snell,
913 K., ... Dildar, N. (2020). Soil and air temperature calibrations using branched
914 GDGTs for the Tropical Andes of Colombia: Toward a pan-tropical calibra-
915 tion. *Geochemistry, Geophysics, Geosystems*, 21(8), e2020GC008941.
- 916 Pérez-Angel, L. C., Sepúlveda, J., Molnar, P., Montes, C., Rajagopalan, B., Snell,
917 K., ... Dildar, N. (2020). *In situ temperature and brGDGTs measurements*
918 *in soils from the Tropical Andes of Colombia and a tropical soil brGDGT*
919 *compilation dataset* [Dataset]. Zenodo. doi: 10.5281/zenodo.3939270
- 920 R Core Team. (2022). R: A Language and Environment for Statistical Computing
921 [Computer software manual]. Vienna, Austria. Retrieved from [https://www.R-](https://www.R-project.org/)
922 [project.org/](https://www.R-project.org/)
- 923 Raberg, J. H., Miller, G. H., Geirsdóttir, Á., & Sepúlveda, J. (2022). Near-universal
924 trends in brGDGT lipid distributions in nature. *Science Advances*, 8(20),
925 eabm7625.
- 926 Rattanasriampaipong, R., Zhang, Y. G., Pearson, A., Hedlund, B. P., & Zhang,
927 S. (2022). Archaeal lipids trace ecology and evolution of marine ammonia-
928 oxidizing archaea. *Proceedings of the National Academy of Sciences*, 119(31),
929 e2123193119.
- 930 Revelle, W., & Revelle, M. W. (2015). Package ‘psych’. *The comprehensive R*
931 *archive network*, 337, 338.
- 932 Schouten, S., Hopmans, E. C., & Damsté, J. S. S. (2013). The organic geochemistry

- of glycerol dialkyl glycerol tetraether lipids: A review. *Organic geochemistry*, 54, 19–61.
- Schouten, S., Hopmans, E. C., Schefuß, E., & Damste, J. S. S. (2002). Distributional variations in marine crenarchaeotal membrane lipids: a new tool for reconstructing ancient sea water temperatures? *Earth and Planetary Science Letters*, 204(1-2), 265–274.
- Sen, P. C., Hajra, M., & Ghosh, M. (2020). Supervised classification algorithms in machine learning: A survey and review. In *Emerging technology in modelling and graphics: Proceedings of iem graph 2018* (pp. 99–111).
- Sinninghe Damsté, J. S., Rijpstra, W. I. C., Hopmans, E. C., den Uijl, M. J., Weijers, J. W., & Schouten, S. (2018). The enigmatic structure of the crenarchaeol isomer. *Organic Geochemistry*, 124, 22–28.
- Sinninghe Damsté, J. S., Rijpstra, W. I. C., Hopmans, E. C., Weijers, J. W., Foessel, B. U., Overmann, J., & Dedysh, S. N. (2011). 13, 16-Dimethyl octacosanedioic acid (iso-diabolic acid), a common membrane-spanning lipid of Acidobacteria subdivisions 1 and 3. *Applied and Environmental Microbiology*, 77(12), 4147–4154.
- Sinninghe Damsté, J. S., Schouten, S., Hopmans, E. C., Van Duin, A. C., & Geenevasen, J. A. (2002). Crenarchaeol: the characteristic core glycerol dibiphytanyl glycerol tetraether membrane lipid of cosmopolitan pelagic crenarchaeota. *Journal of lipid research*, 43(10), 1641–1651.
- Sinninghe Damsté, J. S., Weber, Y., Zopfi, J., Lehmann, M. F., & Niemann, H. (2022). Distributions and sources of isoprenoidal GDGTs in Lake Lugano and other central European (peri-) alpine lakes: Lessons for their use as paleotemperature proxies. *Quaternary Science Reviews*, 277, 107352.
- Siver, P. A., & Lott, A. M. (2023). History of the Giraffe Pipe locality inferred from microfossil remains: a thriving freshwater ecosystem near the Arctic Circle during the warm Eocene. *Journal of Paleontology*, 97(2), 271–291.
- Taylor, K. W., Huber, M., Hollis, C. J., Hernandez-Sanchez, M. T., & Pancost, R. D. (2013). Re-evaluating modern and Palaeogene GDGT distributions: Implications for SST reconstructions. *Global and Planetary Change*, 108, 158–174.
- Tierney, J. E., Russell, J. M., Eggermont, H., Hopmans, E., Verschuren, D., & Sinninghe Damsté, J. S. (2010). Environmental controls on branched tetraether lipid distributions in tropical East African lake sediments. *Geochim. Cosmochim. Acta*, 74(17), 4902–4918.
- Ueki, K., Hino, H., & Kuwatani, T. (2018). Geochemical discrimination and characteristics of magmatic tectonic settings: A machine-learning-based approach. *Geochemistry, Geophysics, Geosystems*, 19(4), 1327–1347.
- Véquaud, P., Thibault, A., Derenne, S., Anquetil, C., Collin, S., Contreras, S., ... Huguet, A. (2022). FROG: A global machine-learning temperature calibration for branched GDGTs in soils and peats. *Geochimica et Cosmochimica Acta*, 318, 468–494.
- Weijers, J. W., Schouten, S., Hopmans, E. C., Geenevasen, J. A., David, O. R., Coleman, J. M., ... Sinninghe Damsté, J. S. (2006). Membrane lipids of mesophilic anaerobic bacteria thriving in peats have typical archaeal traits. *Environmental Microbiology*, 8(4), 648–657.
- Weijers, J. W., Schouten, S., van den Donker, J. C., Hopmans, E. C., & Damsté, J. S. S. (2007). Environmental controls on bacterial tetraether membrane lipid distribution in soils. *Geochimica et Cosmochimica Acta*, 71(3), 703–713.
- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ... Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. doi: 10.21105/joss.01686
- Windler, G., Tierney, J. E., DiNezio, P. N., Gibson, K., & Thunell, R. (2019). Shelf exposure influence on Indo-Pacific Warm Pool climate for the last 450,000

- 988 years. *Earth and Planetary Science Letters*, 516, 66–76.
- 989 Wolfe, A. P., Reyes, A. V., Royer, D. L., Greenwood, D. R., Doria, G., Gagen,
990 M. H., . . . Westgate, J. A. (2017). Middle Eocene CO_2 and climate recon-
991 structed from the sediment fill of a subarctic kimberlite maar. *Geology*, 45(7),
992 619–622.
- 993 Wright, M. N., Wager, S., Probst, P., & Wright, M. M. N. (2019). Package ‘ranger’.
994 *Version 0.11*, 2.
- 995 Yavitt, J., Harms, K., Garcia, M., Wright, S., He, F., & Mirabello, M. (2009). Spa-
996 tial heterogeneity of soil chemical properties in a lowland tropical moist forest,
997 Panama. *Soil Research*, 47(7), 674–687.
- 998 Zell, C., Kim, J.-H., Moreira-Turcq, P., Abril, G., Hopmans, E. C., Bonnet, M.-P.,
999 . . . Damsté, J. S. S. (2013). Disentangling the origins of branched tetraether
1000 lipids and crenarchaeol in the lower Amazon River: Implications for GDGT-
1001 based proxies. *Limnology and Oceanography*, 58(1), 343–353.
- 1002 Zhang, Y. G., Zhang, C. L., Liu, X.-L., Li, L., Hinrichs, K.-U., & Noakes, J. E.
1003 (2011). Methane Index: A tetraether archaeal lipid biomarker indicator for
1004 detecting the instability of marine gas hydrates. *Earth and Planetary Science*
1005 *Letters*, 307(3-4), 525–534.
- 1006 Zhao, B., Castañeda, I. S., Salacup, J. M., Thomas, E. K., Daniels, W. C., Schnei-
1007 der, T., . . . Bradley, R. S. (2022). Prolonged drying trend coincident with the
1008 demise of Norse settlement in southern Greenland. *Science advances*, 8(12),
1009 eabm4346.
- 1010 Zheng, Y., Heng, P., Conte, M. H., Vachula, R. S., & Huang, Y. (2019). System-
1011 atic chemotaxonomic profiling and novel paleotemperature indices based on
1012 alkenones and alkenoates: Potential for disentangling mixed species input.
1013 *Organic Geochemistry*, 128, 26–41.
- 1014 Zheng, Y., Liu, H., Yang, H., Wang, H., Zhao, W., Zhang, Z., . . . Liu, W. (2022).
1015 Decoupled Asian monsoon intensity and precipitation during glacial-
1016 interglacial transitions on the Chinese Loess Plateau. *Nature Communications*,
1017 13(1), 5397.
- 1018 Zheng, Y., Pancost, R. D., Liu, X., Wang, Z., Naafs, B., Xie, X., . . . Yang, H.
1019 (2017). Atmospheric connections with the north Atlantic enhanced the
1020 deglacial warming in northeast China. *Geology*, 45(11), 1031–1034.