

# Paleoceanography and Paleoclimatology\*

## RESEARCH ARTICLE

10.1029/2023PA004611

### Key Points:

- The distribution of glycerol dialkyl glycerol tetraethers (GDGTs) is particular to each depositional environment and has unique responses to environmental factors
- The BIGMaC algorithm captures the relationship between both branched and isoprenoid GDGTs (isoGDGTs) with depositional environments
- Our approach can provide paleoclimatological and paleoenvironmental information based only on GDGTs

### Supporting Information:

Supporting Information may be found in the online version of this article.

### Correspondence to:

P. Martínez-Sosa,  
[pmartoza@arizona.edu](mailto:pmartoza@arizona.edu)

### Citation:

Martínez-Sosa, P., Tierney, J. E., Pérez-Angel, L. C., Stefanescu, I. C., Guo, J., Kirkels, F., et al. (2023). Development and application of the branched and isoprenoid GDGT machine learning classification algorithm (BIGMaC) for paleoenvironmental reconstruction. *Paleoceanography and Paleoclimatology*, 38, e2023PA004611. <https://doi.org/10.1029/2023PA004611>

Received 16 JAN 2023

Accepted 3 JUL 2023

### Author Contributions:

**Conceptualization:** Pablo Martínez-Sosa, Jessica E. Tierney

**Data curation:** Pablo Martínez-Sosa, Jessica E. Tierney, Ioana C. Stefanescu, Jingjing Guo, Frédérique Kirkels, Julio Sepúlveda, Francien Peterse, Bryan N. Shuman, Alberto V. Reyes

**Formal analysis:** Pablo Martínez-Sosa

**Funding acquisition:** Jessica E. Tierney









**Methodology:** Pablo Martínez-Sosa, Jessica E. Tierney

**Project Administration:** Jessica E. Tierney

**Software:** Pablo Martínez-Sosa

**Supervision:** Jessica E. Tierney

## Development and Application of the Branched and Isoprenoid GDGT Machine Learning Classification Algorithm (BIGMaC) for Paleoenvironmental Reconstruction

Pablo Martínez-Sosa<sup>1</sup> , Jessica E. Tierney<sup>1</sup>, Lina C. Pérez-Angel<sup>2</sup> , Ioana C. Stefanescu<sup>3</sup> , Jingjing Guo<sup>4</sup> , Frédérique Kirkels<sup>4</sup> , Julio Sepúlveda<sup>5</sup>, Francien Peterse<sup>4</sup> , Bryan N. Shuman<sup>3</sup> , and Alberto V. Reyes<sup>6</sup> 

<sup>1</sup>Department of Geosciences, The University of Arizona, Tucson, AZ, USA, <sup>2</sup>Institute at Brown for Environment and Society (IBES), Brown University, Providence, RI, USA, <sup>3</sup>Department of Geology and Geophysics, University of Wyoming, Laramie, WY, USA, <sup>4</sup>Department of Earth Sciences, Utrecht University, Utrecht, The Netherlands, <sup>5</sup>Department of Geological Sciences and Institute of Arctic and Alpine Research (INSTAAR), University of Colorado Boulder, Boulder, CO, USA, <sup>6</sup>Department of Earth and Atmospheric Sciences, University of Alberta, Edmonton, AB, Canada

**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 1,153 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 data set was used to train and test the **Branched and isoGDGT Machine learning Classification (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 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 Archeal 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 crenarchaeol' (Sinninghe Damsté et al., 2002, 2018).

Both isoprenoid and brGDGTs 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

**Writing – original draft:** Pablo Martínez-Sosa

**Writing – review & editing:** Jessica E. Tierney, Ioana C. Stefanescu, Jingjing Guo, Julio Sepúlveda, Francien Peterse, Bryan N. Shuman, Alberto V. Reyes

sediments, the degree of cyclization of isoGDGTs is related to overlying water temperature, forming the basis of the TetraEther index of 86 carbons ( $TEX_{86}$ ) 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 (Dang et al., 2018; De Jonge, Stadnitskaia, et al., 2014; Martínez-Sosa et al., 2020; Raberg et al., 2022; Tierney et al., 2010; Weijers et al., 2007). The Methylation index of Branched Tetraethers ( $MBT'_{5Me}$ ) proxy isolates the relationship between the methylation of brGDGTs and temperature (De Jonge, Hopmans, et al., 2014; Weijers et al., 2007) and has been widely used for terrestrial paleoclimate reconstructions (De Jonge, Hopmans, et al., 2014; Lauretano et al., 2021; Naafs et al., 2018; Zhao et al., 2022; Zheng et al., 2017).

Across environments, GDGT distributions broadly reflect the microbial community present. This is, for example, the basis of the Methane Index, which measures the contribution of methanotrophic archaea relative to marine Thaumarchaeota in the sedimentary isoGDGT pool (Zhang et al., 2011). Likewise, the distribution of isoGDGTs in marine systems reflects not only sea-surface temperature (captured by the  $TEX_{86}$  index), but also the water depth (and potentially, different archeal communities) from which the isoGDGTs derive (Rattanasriampaipong et al., 2022; Taylor et al., 2013). Furthermore, previous work has found that the ratio of crenarchaeol/crenarchaeol' can indicate which *Thaumarchaeota* group (Figures I1a or I1b in Supporting Information S1) is responsible for the production of these lipids in lake sediments (Li et al., 2016). In terrestrial settings, De Jonge et al. (2019) proposed the Community Index for brGDGTs, which is based on the inference that brGDGTs are produced by different communities of bacteria, each with a unique response to soil temperature. The combined use of some of the GDGTs, through the Branched and Isoprenoid Tetraether (BIT) index, is used to broadly discriminate between marine and terrestrial environments based on the dominance of brGDGT-producing bacteria in most terrestrial settings and crenarchaeol-producing Thaumarchaeota in marine settings (Hopmans et al., 2004). However, BIT values in soils, lakes, and peats all tend to be high, which limits the ability of this index to reliably distinguish between these different types of terrestrial settings.

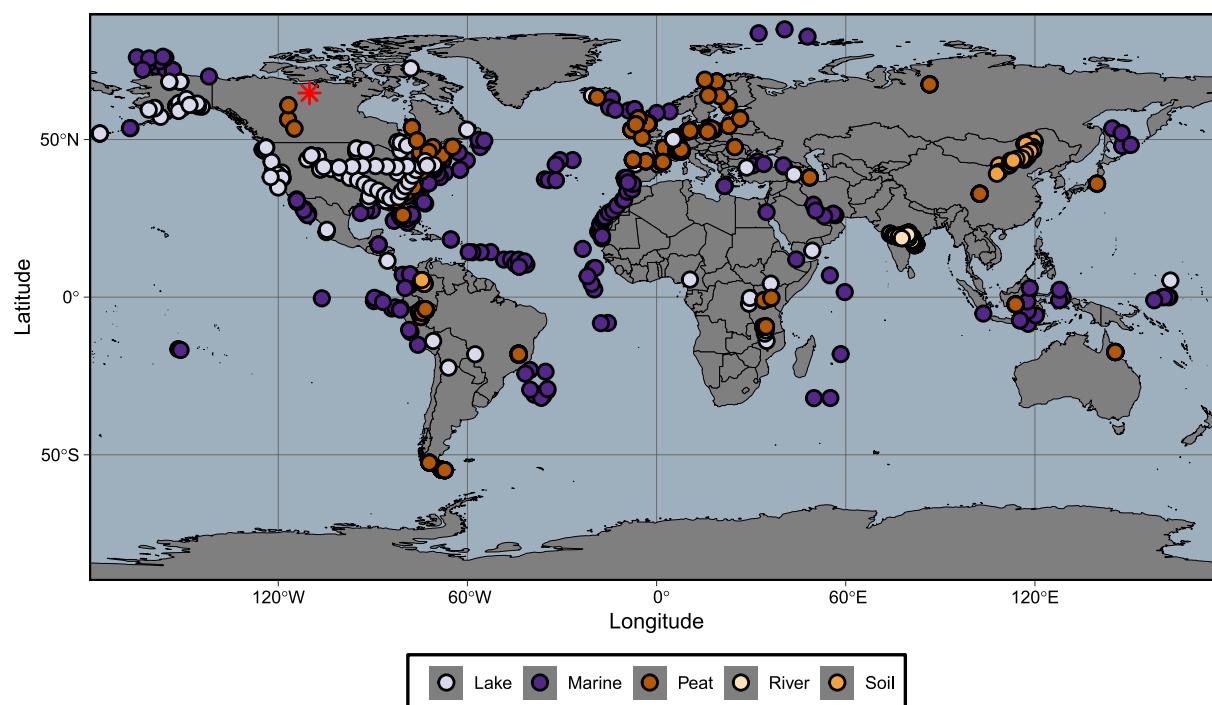
Building on these observations, we posit that the full range of archeal and bacterial GDGTs (isoprenoidal and branched) contains information about their biological precursors and the overall composition of the microbial community. This information, based on ecological interactions rather than a physiological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

Machine learning has previously been used in the Geosciences to discriminate between magma (Ueki et al., 2018) as well as identifying the source of water from oil wells (Engle & Brunner, 2019). In the field of biomarker-based paleoclimatology, classification algorithms have been applied to identify sources of alkenones (Zheng et al., 2019), as well as plant waxes (Peale et al., 2021). Machine learning regression algorithms, as well as deep neural network applications, have also been applied to GDGTs as in order to generate temperature calibrations (Dunkley Jones et al., 2020; Véquaud et al., 2022; Zheng et al., 2022). Here, we use a compilation of GDGT distributions in 1,153 globally distributed soils, peats, and sediments from diverse depositional environments to train a classification algorithm which is capable of identifying the environment in which a sample was formed based on the distribution of both branched and isoGDGTs. We then demonstrate the application of this algorithm by using it to interpret the paleoenvironment and the paleotemperature in a Paleogene deposit that records a transition from a lacustrine to a peatland environment. We also highlight the limitations of this approach in an application to a peatland data set that spans the onset of the Paleocene-Eocene Thermal Maximum (PETM).

## 2. Materials and Methods

### 2.1. Global Data Set

We compiled a total of 1,153 globally distributed (Figure 1 and Table S11 in Supporting Information S1) samples from different depositional environments: marine, lake, peat, river, and soil. These samples all have quantified



**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.

relative abundances for the full suite of the most commonly used isoGDGTs (GDGT-0, GDGT-1, GDGT-2, GDGT-3, crenarchaeol, and crenarchaeol') and brGDGTs (IIIa, IIIa', IIIb, IIIb', IIa, IIa', IIb, IIb', IIc, IIc', Ia, Ib, and Ic) in paleoenvironmental reconstructions, and were all analyzed with the updated High Performance Liquid Chromatography-Mass Spectrometry (HPLC-MS) method of Hopmans et al. (2016). From the 1,153 samples, 475 are peat (Naafs et al., 2018), 215 are marine sediments (this study), 196 are soil (Dearing Crampton-Flood et al., 2020; Guo et al., 2020a, 2020b; Guo, Ma, et al., 2022; Pérez-Angel et al., 2020a, 2020b), 162 are lake sediments (Guo et al., 2020a, 2020b; Martínez-Sosa et al., 2021), and 105 are riverbed sediments (Kirkels, Usman, & Peterse, 2022). For the Colombian and Inner Mongolia soil samples (Guo, Ma, et al., 2022; Pérez-Angel et al., 2020a, 2020b) we include here newly reported isoGDGT values not part of the original data set.

All marine sediments were processed at the University of Arizona following the method used in Martínez-Sosa et al. (2021). Briefly, sediments were freeze-dried, homogenized, and spiked with a C<sub>46</sub> internal standard before being extracted with an Accelerated Solvent Extraction system (run at 1,500 psi, 100°C, with dichloromethane:methanol (DCM: MeOH, 9:1)). Total Lipid Extracts (TLEs) were eluted through a deactivated SiO<sub>2</sub> column with hexane:ethyl acetate (1:2) to obtain the polar fraction, and the elutant was dried under a N<sub>2</sub> stream. Polar fractions were redissolved in hexane:isopropanol (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 1,260 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 data set 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 (FA) 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 data sets (Table S1 in Supporting Information S1). For the marine sediments analyzed for this study, we obtained mean annual temperature of the top 200 m 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 data set. The optimal number of clusters for this data set 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 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 data set was split into a training and testing set in a 3:1 ratio. To avoid subsampling the data set 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 Neighbor and Naive Bayes plus a control non-informative (null) model. Below we present a brief overview of each of them.

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

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

Naive Bayes is a statistical algorithm, in which Bayes Theorem is applied to calculate the posterior probability of a new sample belonging to each possible categorical classification. The algorithm determines the classification by choosing that with the highest probability. Although this is a fast classification algorithm it is a relatively bad estimator (Sen et al., 2020).

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

For all algorithms the hyperparameter values—parameters whose values control the learning process but are not part of the final model (i.e., number of independent trees used for a Random Forest)—were selected (tuned) as those with the best performance from a distribution of possible combination of values for each hyperparameter. Due to the required computing power for the hyperparameter tuning step, this was run using a High-Performance

Computing cluster. Finally, the best hyperparameter values were selected by comparing their Receiver Operating Characteristic Area Under the Curve (ROC-AUC) score on the validation set (Table S1 in Supporting Information S1). ROC-AUC is a metric that measures the tradeoff between the true positive rate and false positive rate of the model, for this parameter higher values on the (0,1) range are desirable. An additional metric, F1, was applied to evaluate the performance of all the trained algorithms. F1 is calculated as the mean of the precision and recall—the ability to identify true positives—of the model. For reproducibility, a detailed script with the packages and parameter values used for this process is available on GitHub (Martínez-Sosa et al., 2023). Finally we can identify which GDGTs contribute the most to the classification process through the importance metric. This value is calculated based on how much each GDGT contributes to decreasing the probability of incorrectly classifying a sample across all decision trees (Gini importance) (Greenwell et al., 2020; Wright et al., 2019). This metric shows how often a GDGT was selected to split the data and what was its discriminative power (Menze et al., 2009), larger values indicate a better variable to separate the data.

For this work, all analyses were performed in R (v. 4.1.3) (R Core Team, 2022). Additional dimensionality reduction analyses were done over the fractional abundances of all GDGTs using Principal Component Analysis (PCA) through the `princomp()` function from base R. For the PCA analyses the loadings of variables were visually inspected and corroborated by obtaining the eigenvector values. Correlation between environmental parameters for each sample and their scores on the principal components was performed by applying the `cor.test()` function from the *psych* R package using a Spearman 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  $\geq 50$  vertical-equivalent meters of lacustrine sediment topped with  $\sim 32$  m of peat. The sediments were dated to  $37.84 \pm 1.99$  Ma by glass fission-track dated rhyolitic tephra beds (Wolfe et al., 2017). Our data set 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-NH<sub>2</sub> (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 isoGDGTs 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

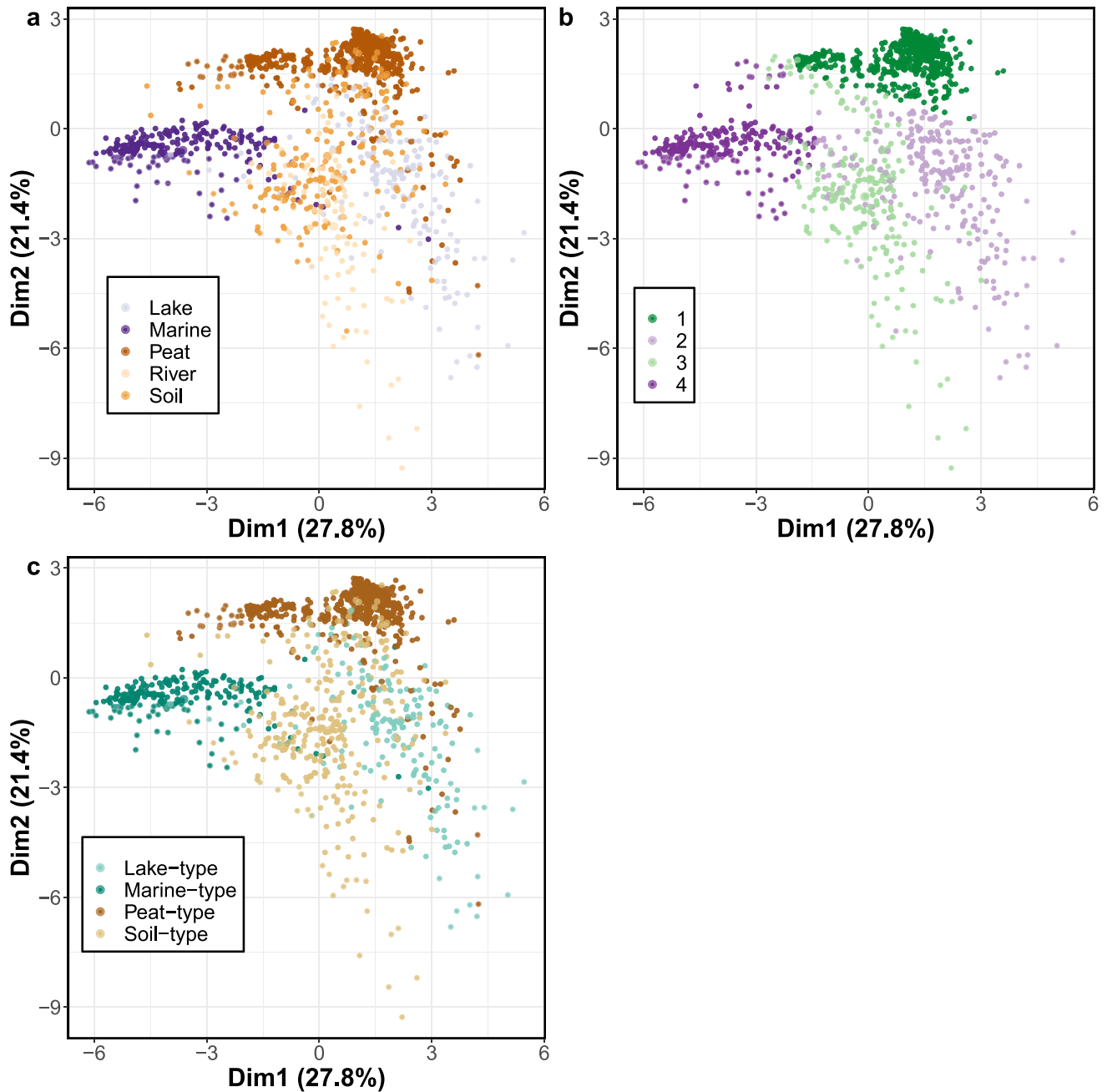
The Cobham lignite bed, Kent, UK ( $\sim 48^\circ\text{N}$  paleolatitude) is composed of a sand and mud unit at the base, overlain, in succession, by a charcoal-rich lower laminated lignite (LL), a charcoal-poor upper LL, a middle clay layer, and a charcoal-poor blocky lignite (BL). The Woolwich Shell Beds overly the Cobham Lignite (Collinson et al., 2009). A carbon isotope excursion is present near the top of the charcoal-poor upper LL, which is interpreted as being the characteristic excursion from the Paleocene Eocene Thermal Maximum (PETM,  $\sim 56$  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 (Figures 2a and 2b, 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 (Figures 2a and 2b). Soil samples are more spread across the different groups, with the majority assigned to Group 3 (43.9%).





**Figure 2.** Samples from the data set plotted in reduced dimensional space based on the fractional abundance of glycerol dialkyl glycerol tetraethers. 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).

### 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 data set crenarchaeol', GDGT-1–GDGT-3, Ib, Ic, Iic, Iic', IIIb, and IIIb' have the smallest proportion (<0.1 FA) of all GDGTs (Figure 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 extent GDGT-1 and GDGT-2 compared with the other clusters (Figure 3a). For crenarchaeol, cluster 3 is the

**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)

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

Note. 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.

next group with the highest proportion, while cluster 2 has the second highest proportion of GDGT-0. Cluster 1 consistently shows the lowest isoGDGT values. On the other hand clusters 1, 2, and 3 have a higher proportion of brGDGTs than cluster 4. Cluster 1 has the highest proportion of brGDGTs Ia and IIa, while cluster 2 has the highest proportion of IIIa and IIIa'. Although cluster 3 has lower proportion values than cluster 1, it is also dominated by the penta- and tetramethylated brGDGTs, and it shows the highest proportion of GDGT Ib. Both clusters 2 and 3 show comparable levels of IIa', the most abundant 6-methyl brGDGT.

### 3.2.2. Environmental Influence on Glycerol Dialkyl Glycerol Tetraethers

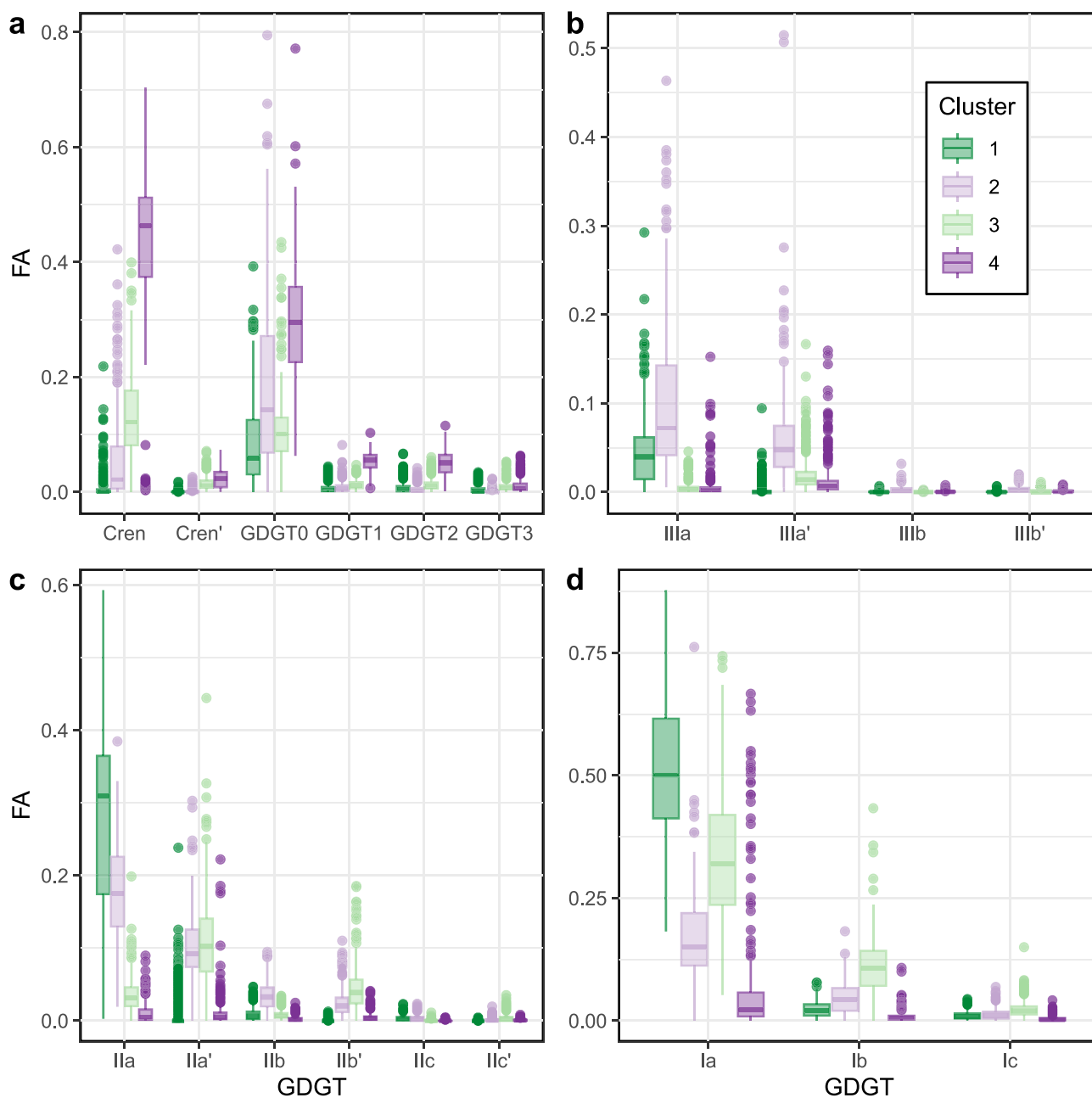
To better understand the effect that environmental parameters (such as temperature, elevation, and pH) might have had on the classification of the samples, we performed a PCA for each of the sample types separately (Figure 4 and Figure S11 and Table S12 in Supporting Information S1).

For peats (Figure 4a), primarily classified in cluster 1, the main component (64.6% of variability) is positively associated with GDGT Ia, while negatively related to GDGT IIa. This component is strongly associated with MAAT ( $\rho = 0.77$ , Spearman's correlation and Figure S11a in Supporting Information S1). Peats classified in cluster 1 plot throughout the first component, however, those peats classified as part of cluster 3 and 4 are associated with higher temperatures, while those classified as part of cluster 2 are associated with lower MAATs (Figure 4a).

For lake sediments (Figure 4b), which are mostly classified as part of cluster 2, the first component (36.6% of variability) is associated negatively with GDGT Ia, while positively related to GDGT IIIa. The second component (32.4% variability) is positively associated with GDGT-0. Both components have a strong correlation with MAAT ( $\rho = -0.65$  and  $0.51$ , respectively, and Figure S11b in Supporting Information S1). While the majority of the distribution cloud is classified in cluster 2, samples with more negative values on PC1, which have higher MAAT values, are classified as cluster 3 and 1 (Figure 4b and Figure S11b in Supporting Information S1). Additionally, four samples are classified in cluster 4, which were identified as sediments from Lake Kivu, Mono Lake, and two from Lake Malawi.

Soil samples are classified into a wide range of clusters (Figure 4c), with 44% of the samples labeled as part of cluster 3, 31% as cluster 2, 20% as cluster 1, and 10% as cluster 4. The first component of the PCA for this sample type (47% of variability) is negatively associated with brGDGT Ia, while the second component is positively related with GDGT IIa and negatively with crenarchaeol. The environmental parameters show strong correlations to both the first and second components: pH ( $\rho = 0.47$  and  $-0.74$ ), MAAT ( $\rho = -0.45$  and  $-0.42$ ), and elevation ( $\rho = -0.21$  and  $-0.56$ , respectively). Additionally, the assigned cluster strongly correlates with the sample location (Figure S11c in Supporting Information S1). Most samples from the Godavari river catchment (98%) and Inner Mongolia (69%) were classified as cluster 3, while 70% of samples from Carminowe creek are classified as cluster 2. Finally, 47% of samples from Colombia are classified as cluster 1.

Similar to soils, the river sediments are mostly classified as cluster 3 (Figure 4d). For these samples the first principal component explains 51.4% of the variance and is negatively associated with brGDGT Ia. For this sample type, the most determinant variable for their cluster classification was location site (Figure S11d in Supporting Information S1), with all but one out of 98 samples from the Godavari classified in cluster 3, while all samples from Carminowe Creek are classified as cluster 2.

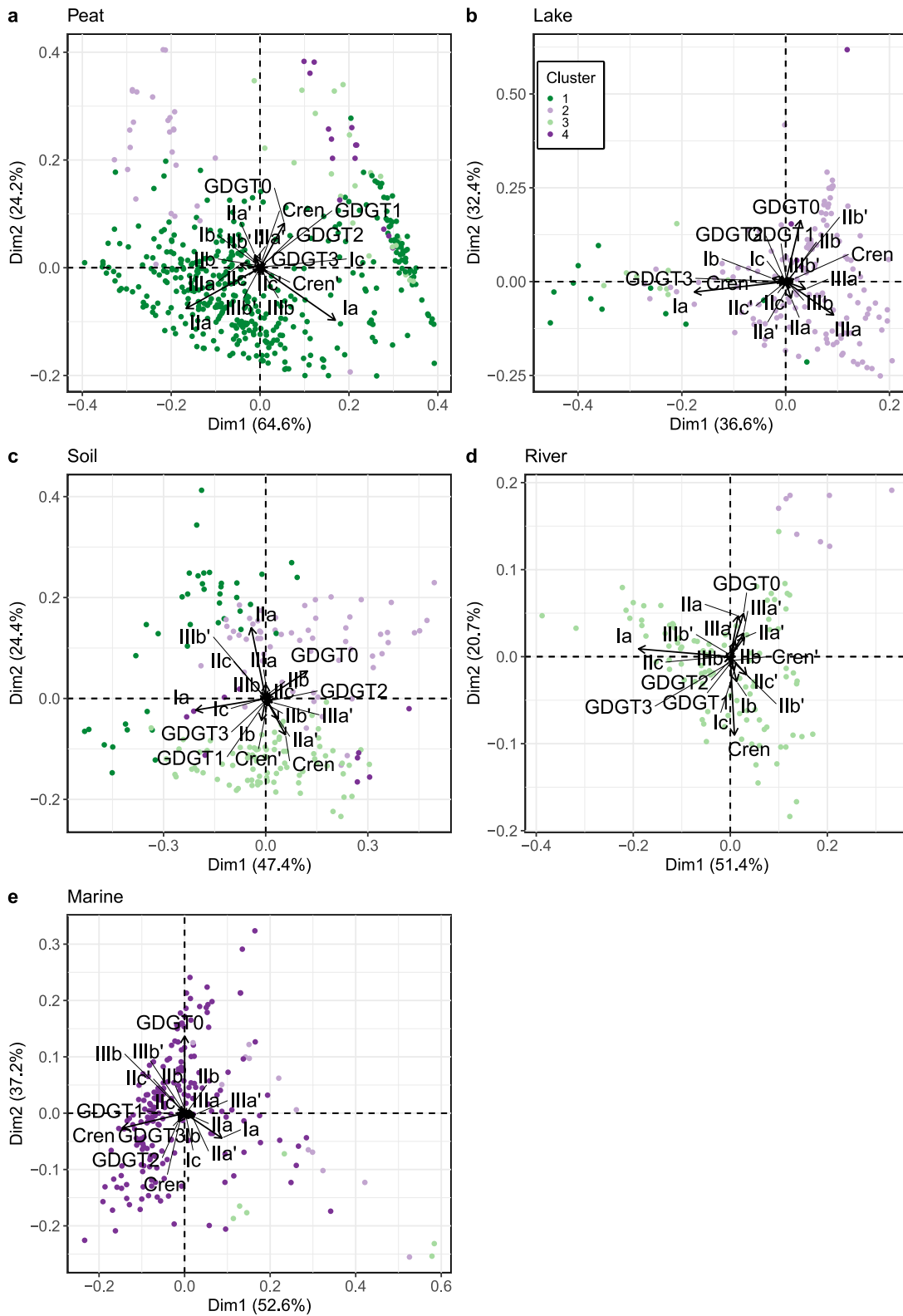


**Figure 3.** Box plots showing the distribution of the fractional abundance of all Glycerol dialkyl glycerol tetraethers (GDGTs) in each of the fuzzy  $k$ -means clusters, following the color code of Figure 2b. GDGTs separated by isoprenoid GDGTs (a), hexamethylated branched GDGTs (brGDGTs) (b), pentamethylated brGDGTs (c), and tetramethylated brGDGTs (d).

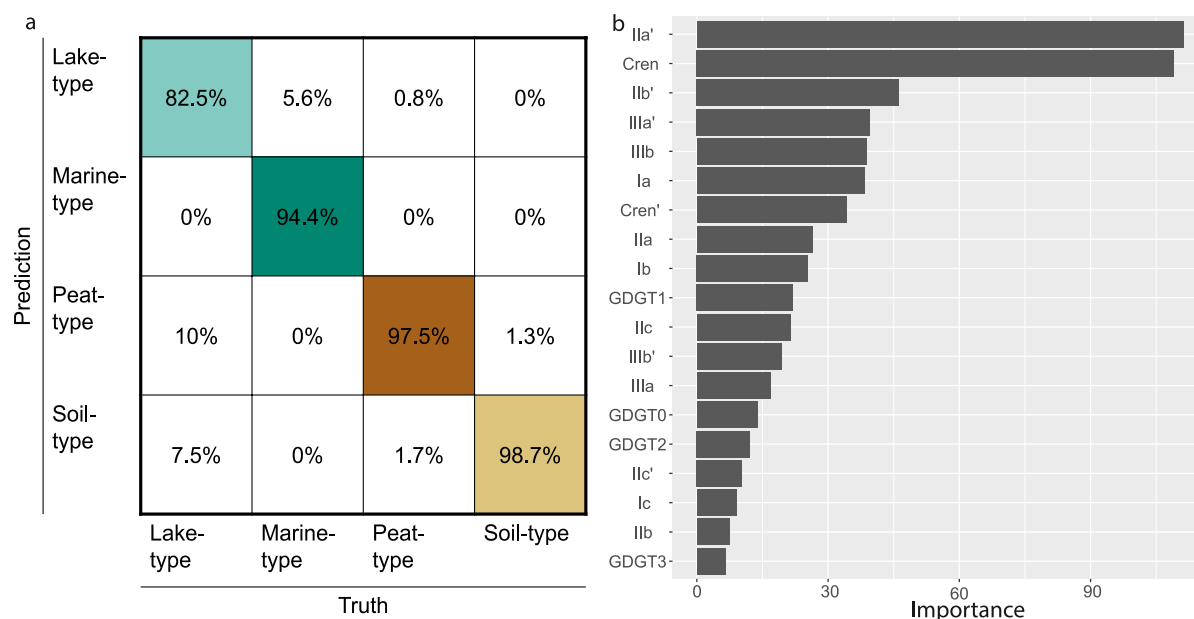
An additional PCA with both river and soil samples (data not shown) showed that none of the first three components (88% of combined variance) separated these samples by type, but rather soil and river samples cluster together based on the location (Godavari or Carminowe).

Finally, for the marine sediments, primarily classified as cluster 4 (Figure 4e), the first component of the PCA (52.6% of variability) is positively associated with GDGT Ia and negatively associated with crenarchaeol. The second component (37.2%) is positively associated with GDGT-0. The first component is associated with the sample's distance to the coast ( $\rho = 0.43$  and Figure S11e in Supporting Information S1), while the second component shows a strong correlation with the mixed layer temperature ( $\rho = -0.79$ ). We, however, do not find a strong association between samples classified as cluster 2 or 3 and environmental parameters considered here.





**Figure 4.** Principal component analysis (PCA) of each sample type with samples colored based on their assigned *K*-means cluster. Each panel shows the PCA for peat (a), lake sediments (b), soils (c), river sediments (d), and marine sediments (e).



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

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

### 3.3. Supervised Machine Learning

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

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

The finalized model, named **Branched and isoGDGT Machine learning Classification algorithm (BIGMaC)**, is available on Github as an R object (Martínez-Sosa et al., 2023).

## 4. Discussion

### 4.1. Unsupervised Machine Learning

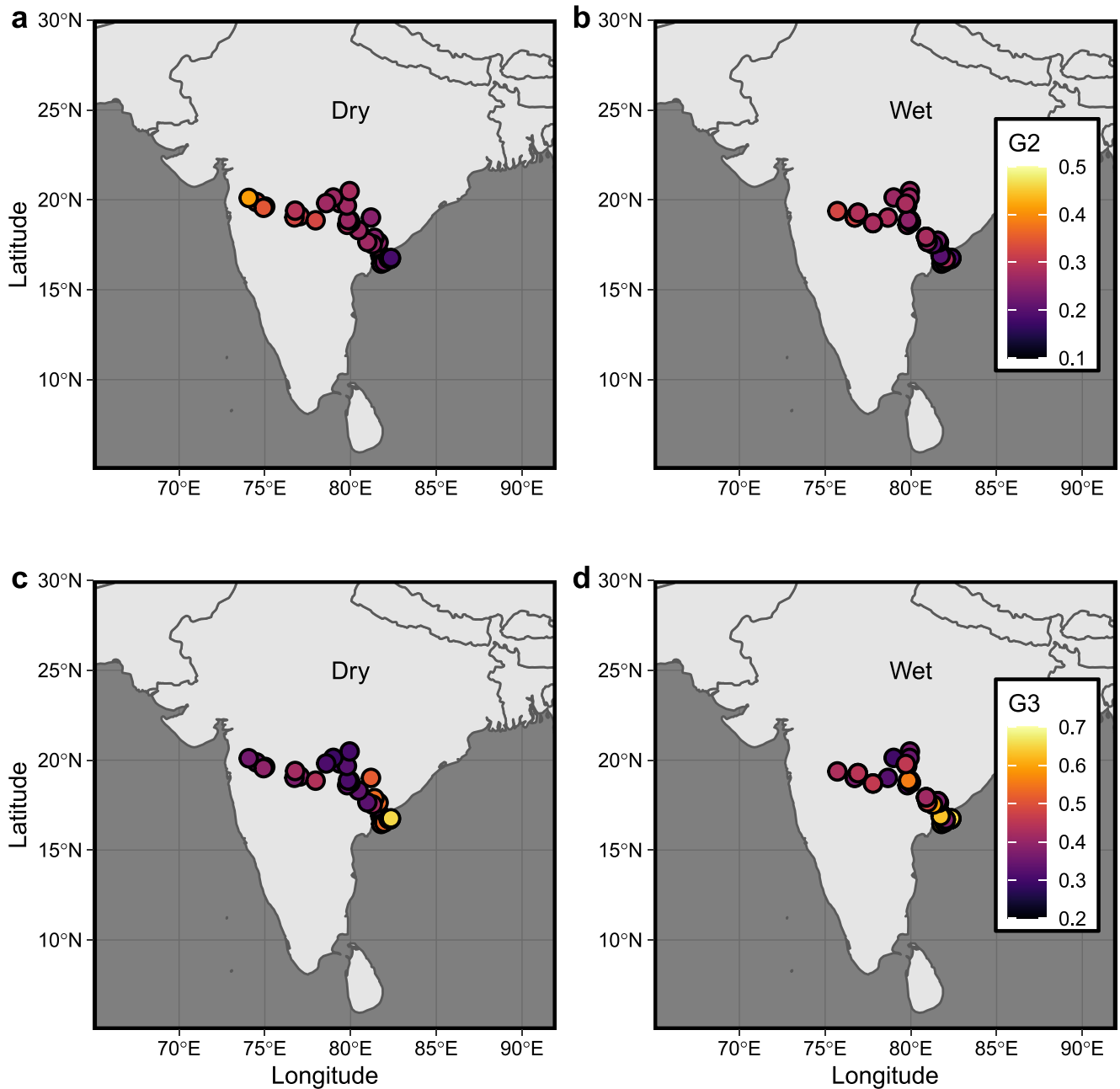
While our fuzzy  $k$ -means clusters show strong patterns that reflect relationships with depositional environments (Figure 2a), some samples whose context was unequivocally documented cluster in with samples unrelated to their depositional environment (i.e., soils plotting as peats).

For both peats and lakes, temperature has a strong influence on the GDGT distribution, causing samples at either ends to be clustered with other sample types (Figures 4a and 4b, and Figures SI1a and SI1b in Supporting Information S1). At high temperatures, samples from these depositional environments tend to have a higher proportion of brGDGT Ia, and lower proportion of IIa and IIIa (Weijers et al., 2007), which causes the lake sediments to be classified as cluster 1 and 3 (high Ia), or peats as cluster 3 and 4 (low IIa and IIIa). At lower temperatures, the opposite effect causes some of the peats to be classified as cluster 2. In addition, particularly deep (i.e., Lake Malawi) or hypersaline alkaline lakes also showed an increased proportion of isoGDGTs compared with other lakes, causing them to be classified with marine samples. These distinct distributions could be due to specific niches in the water column associated with water chemistry, stratification, and/or nutrient content, as previous work has suggested (Baxter et al., 2021; Kumar et al., 2019; Sinninghe Damsté et al., 2022).

Soil samples cluster into several different groups, which may reflect the fact that soils are highly diverse environments with diffuse boundaries and are often in contact with other depositional settings. Furthermore, studies have shown that chemical properties of soils (i.e., pH, metal concentrations) have great spatial heterogeneity even at small scales (Yavitt et al., 2009). Indeed, the main factor that separated soils in this study is their location (Figure 4c and Figure SI1c in Supporting Information S1), which overrides any other environmental signal. However, this may simply be a feature of the limited existence of soil data sets with both branched and isoGDGTs (represented by only four locations). It is clear from our results that soils require a more in-depth analysis, with the use of more extensive data sets.

While there is some debate regarding the relative influence that soil input and in situ production have on the GDGT pool in river organic matter (De Jonge, Stadnitskaia, et al., 2014; Kirkels et al., 2020; Zell et al., 2013), our analysis shows that the river sediments more closely resemble soils rather than peats or lake sediments, and similarly to soils, the sample location shows the strongest correlation with how these samples are classified by fuzzy  $k$ -means (Figure 4d and Figure SI1d in Supporting Information S1). We do find a difference in the 5-methyl/6-methyl proportion in these groups, where soils have a relatively similar proportion of these isomers, while river sediments contain relatively more 6-methyl brGDGTs (Figure SI2b and SI2c in Supporting Information S1). Although this could be interpreted as soil-derived GDGTs dominating river inputs, with some autochthonous production of 6-methyl brGDGTs in rivers, our river data come from only two locations and primarily from only one system (the Godavari river), so this interpretation could be particular to this watershed. Notably, within the Godavari River, the membership value for the samples, which measures the degree of belonging to each cluster, varies with their location and collection season (Figure 6). Membership to the soil-dominated Group 3 is higher in the lower Godavari basin, as well as from the wet (post-monsoon) season (Figures 6c and 6d). In contrast, membership to the lake-dominated Group 2 is overall higher in the dry season, and in the upper basin year-round (Figures 6a and 6b). These results are in line with those presented in the original study by Kirkels, Zwart, et al. (2022), where it was noted that GDGTs from soils have a stronger influence on the river's GDGT content during the wet season and within the lower basin, which experiences higher precipitation. In contrast, in-situ production of brGDGTs, characterized by a high proportion of 6-methyl isomers, has a stronger influence on the GDGT content of samples from the dry season as well as those from the upper basin.

Notably, for marine sediments the first dimension of their PCA is dominated by a positive relation with brGDGT Ia and a negative one with crenarchaeol and it is associated 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 (Figure 4e and Figure SI1e in Supporting Information S1). 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<sub>86</sub> 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 (Figure 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



**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).

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 (Figures 2a–2c). 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 fed the classification algorithm a more realistic data set 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 (Figure 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*).

#### 4.3. Supervised Classification

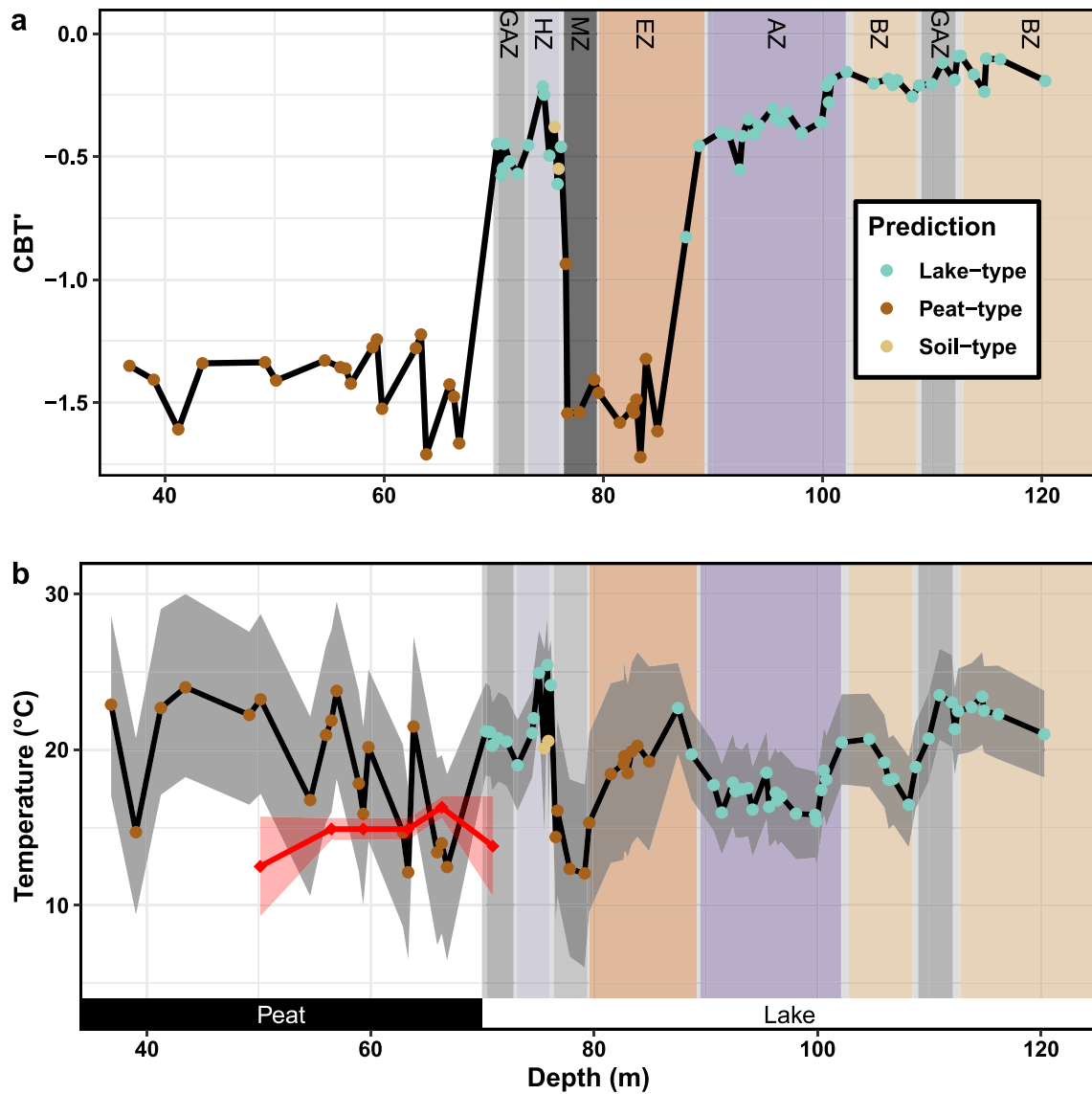
In general, all of the machine learning algorithms exhibited good performance in the training phase, with F1 and ROC-AUC scores above 0.85 and 0.95 respectively. Nevertheless, we chose the Random Forest algorithm since it was the best performing one across all parameters, in addition to being widely used in the field of Geosciences (El Boucheffy & de Souza, 2020; Peuple et al., 2021). This algorithm also performed well in the testing phase (0.94 and 0.99, for F1 and ROC-AUC respectively, and Figure 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.

When we apply the BIGMaC algorithm to the complete data set, 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 Ila' 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 (Figure 3), and also could be associated with characteristics of the depositional environments. BrGDGT Ila' 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 data set 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 Ila'—low crenarchaeol (*Peat-type*), high Ila'—low crenarchaeol (*Lake-type*), low Ila'—high crenarchaeol (*Marine-type*), and high Ila'—high crenarchaeol (*Soil-type*).

#### 4.4. Applications

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.





**Figure 7.** CBT' values (a) and inferred temperature (b) calculated from Giraffe core branched GDGTs 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* in light orange, *Going Acidic Zone* in light gray, *Aulacoseira Zone* in purple, *Eunotiod Zone* in orange, *Mixing Zone* in gray, and *Heterotrophic Zone* 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.

#### 4.4.1. Giraffe Kimberlite Pipe

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 (Hamblin et al., 2003; Wolfe et al., 2017) (Figure 7 and GDGT distributions shown in Figure SI4 in Supporting Information S1). 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 m (Figure 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  $0.999 \pm 0.001$ , suggesting a terrestrial setting. The BIT index record only deviates from these values in the one-m section where the BIGMaC classifies samples as either *Lake-type* or *Soil-type*, although even in this section the BIT index values are only reduced to 0.97. Since the BIT index is unable to distinguish between soil, peat and lake deposits, this showcases the advantage of using the classification algorithm where all GDGTs are considered over a singular index.

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

Our results show a close relationship with the independently developed microfossil-based ecological reconstruction done by Siver and Lott (2023) for this site. While Hamblin et al. (2003) had previously speculated this site to be a shallow lacustrine setting with intermittent wet and dry periods, the microfossil ensemble suggest that the lake was initially a shallow (~1m) slightly acidic lake, this is represented by the *Botryococcus* Zone (BZ) between 102.8 and 124.6 m (Figure 7). All samples in this section are classified as *Lake-type* with generally stable CBT' values. While Siver and Lott (2023) reported a shift in the microfossil ensemble within this region between 112.1 and 109 m, which was interpreted as an acidifying section, we do not find evidence of this in the CBT' values; however, this subsection does roughly align with an estimated reduction in  $MBT'_{5Me}$  values. Following this initial section, the microfossils suggest that the lake deepened to between 3 and 5 m in the *Aulacoseira* Zone (AZ) (Siver & Lott, 2023), between 102.1 and 89.5 m. BIGMaC still classifies all samples from this section as *Lake-type*, however there is a declining trend in CBT' values throughout this region, shifting from values of  $-0.2$  to  $-0.55$ . The following section located between 89.2 and 79.6 m, identified by the microfossil ensemble as Eunotid Zone (EZ), is thought to represent a period characterized by enhanced acidity and increase in dissolved humic material. This section corresponds to the large CBT' excursion in our record, with values from  $-0.5$  to below  $-1.5$ . Microfossils in this section are associated with low pH environments, such as lakes and bogs. This is in agreement with our BIGMaC results, where all but one sample in this section are classified as *Peat-type*, with the exception being a sample right at the bottom that BIGMaC still classified as *Lake-type*. We note that the CBT' values are lower than those of previously reported lakes with a pH of 4.3 (Martínez-Sosa et al., 2021). In contrast, these values correlate to a pH between 4 and 5 in peats, which is more in line with what is associated with the described organisms in this section (Siver & Lott, 2023). Between 79.4 and 76.4 m the microfossil ensemble suggests the presence of a transitional zone that is not associated with any of the other sections. Our results classify all samples within this section as *Peat-type*, although the CBT' values have a steep increase, going from  $-1.5$  to  $-0.5$ , which could suggest a transitory state toward a higher pH environment. The following section, between 76 and 73.1 m, identified as the Heterotrophic Zone by the microfossil ensemble, is interpreted as a period with higher pH and in line with conditions of an eutrophic body of water. Our analysis shows a transitory period at the beginning of this section where BIGMaC interprets two samples as *Soil-type*, while the rest of the samples are classified as *Lake-type*. The CBT' values are on par with those of the previous sections that were classified as *Lake-type*. Finally, the top-most part of the lake section is interpreted as a body of water with increased acidity and levels of dissolved humic matter. While the CBT' values do show a small decrease, the BIGMaC algorithm still classifies all samples in this section as *Lake-type*.

Overall, there is a good agreement in the interpretation of both independent proxy-based reconstructions for the lake section of the core. Differences between the interpretations could be due to constraints of the specific proxies. For example, while the microfossil ensemble suggests periods of increased acidification throughout the BZ region, this is not reflected in the CBT' values. In contrast, while the microfossil ensemble suggests a sudden transition from the AZ to the EZ sections, our results show a decreasing trend in CBT' values throughout AZ leading to the much steeper acidification in EZ. It is possible that the decreases in pH in both cases have different origins, which are only captured by one of the proxies. This underscores the advantage of combining independent proxies for paleoenvironmental reconstructions.

To generate a temperature reconstruction for the Giraffe core, we applied the previously published BayMBT calibration for lakes or soil/peat, depending on the results of BIGMaC for each sample. These particular calibrations were chosen as they are consistent with each other and allow us to generate a continuous confidence interval for the reconstruction. However, we emphasize that the choice of temperature calibration is independent from

BIGMaC and any GDGT-based calibration can be used. For example, we also applied the peat-specific calibration (Naafs, 2017) to the sections of the core classified as *Peat-type* but it generated only marginally different results than BayMBT, so we chose to use the latter.

Our reconstruction suggests a relatively stable climate with no clear trend (Figure 7a). The mean temperature of our reconstruction ( $19 \pm 3.2^\circ\text{C}$ ) agrees with independent studies. A pollen reconstruction at this site (red diamonds in Figure 7a), suggests a MAAT of  $14.5 \pm 1.3^\circ\text{C}$ , with a warmest month mean temperature of  $24.5 \pm 0.8^\circ\text{C}$  (Wolfe et al., 2017). In addition, Jahren and Sternberg (2003) estimated a mean annual temperature of  $13.2 \pm 2^\circ\text{C}$  for the middle Eocene Arctic based on oxygen isotopes measured in calcite preserved in fossil *Metasequoia*. While our estimate is at the upper end of both estimates, they fall within the confidence interval of our reconstruction (Figure 7a). Moreover, both the peat/soil and lake calibrations predict mean annual temperatures above freezing (MAF) rather than strictly MAAT, so if there were freezing temperatures during the winter, the GDGT estimates are expected to be higher. Conversely, if we had used only the lakes or soil/peat calibration for the entire core, there would be large temperature swings of more than  $6^\circ\text{C}$  associated with changes in core lithology. In particular, the excursion to *Peat-type* samples within the lacustrine section would be estimated to be  $5.7^\circ\text{C}$  higher without the BIGMaC-based correction.

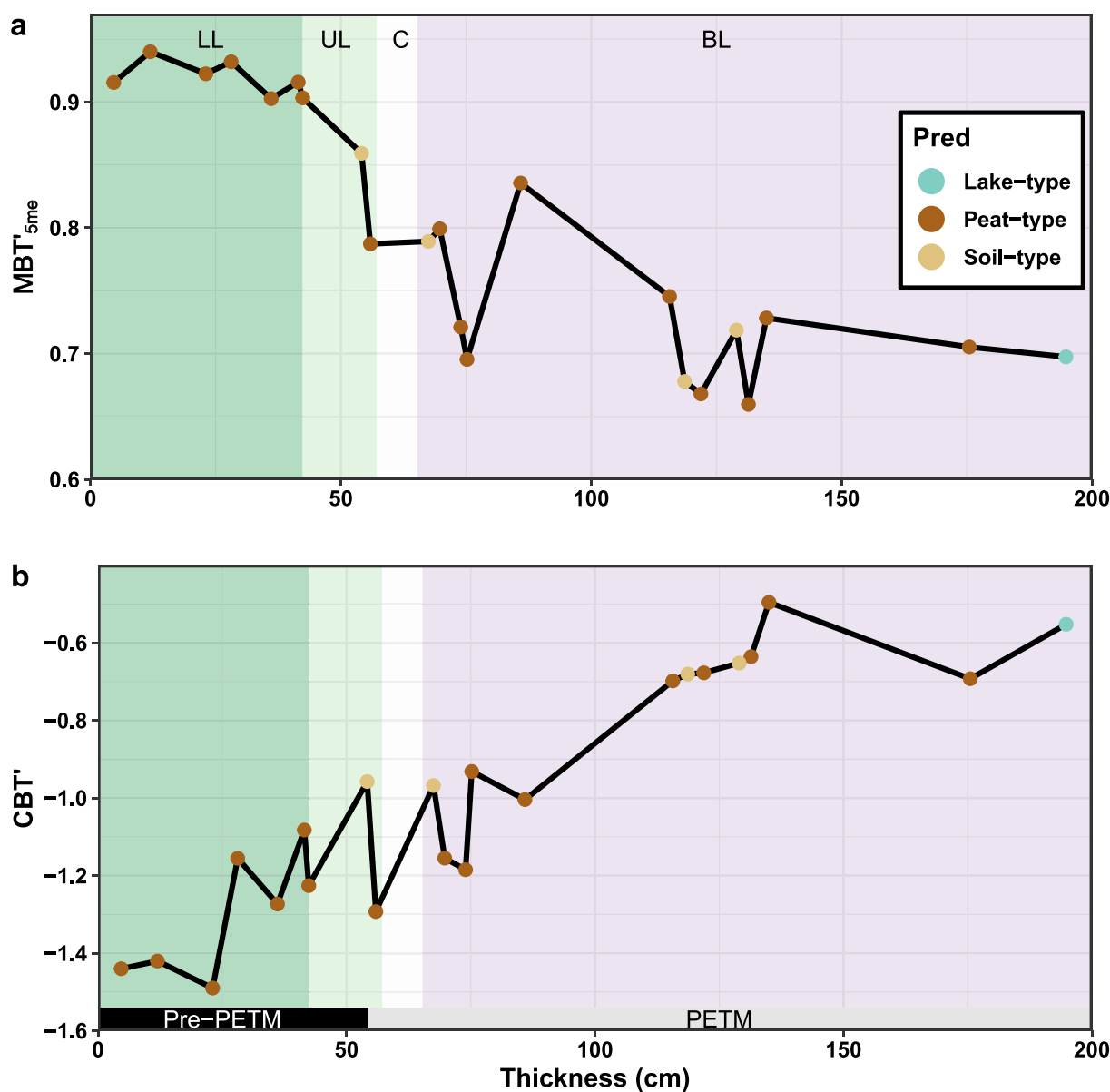
#### 4.4.2. Cobham Lignite Bed

While the application of the BIGMaC algorithm in the Giraffe pipe showcases its strengths, our analysis of the Cobham lignite illustrates that there are some limitations of the approach. Inglis et al. (2019) previously showed that increased precipitation during the PETM in this area caused changes in the hydrology of the site, and that this potentially caused the brGDGTs to become unreliable as temperature proxies. Namely, while several lines of evidence suggest an increase in temperature during the PETM, the temperature reconstructions based on brGDGTs suggest cooling. We applied BIGMaC to this site to investigate whether changes in the depositional settings could explain the discrepancy (Figure 8). Almost all samples preceding 54.15 cm, identified as the start of the PETM, are predicted to be *Peat-type*, with the exception of one sample from the upper LL unit that is classified as *Soil-type*. In contrast, there is a wider variation in the sample classification during the PETM, where samples are classified as *Peat-type* (10), *Soil-type* (3) and *Lake-type* (1). Besides one sample classified as *Peat-type* from the PETM upper LL, all other PETM samples are located in the BL unit. The variations in predicted depositional environments do not coincide with changes in either  $MBT'_{5Me}$  or  $CBT'$  values, nor are they organized in any evident pattern within the unit. Moreover, the PETM samples are primarily classified as *Peat-type* and *Soil-type*, suggesting that adjusting for the predicted depositional environment would not significantly increase the temperature reconstructed by Inglis et al. (2019), as would be the case if the samples were classified as *Lake-type*. Vegetation and charcoal records suggest that the Cobham site became waterlogged and may have even developed areas of open water during the PETM (Inglis et al., 2019). Given this perspective, the oscillating results from BIGMaC likely point to an unstable, dynamically changing depositional environment with mixed sources of brGDGTs. Since BIGMaC is a categorical classification algorithm, it cannot detect mixed signatures. This underlines the need to incorporate mixing models in studies where input from different sources is expected, and suggests that BIGMaC would benefit from incorporating this capability in future updates.

Overall, this application suggests that other environmental parameters, not considered in this work, could still affect the distribution of GDGTs, and we include the Cobham lignite example as a reminder that applying the BIGMaC algorithm is not always a panacea and is best done in concert with other independent proxies to accurately interpret the results.

## 5. Conclusions

Our analyses of GDGTs in 1,153 globally distributed samples from soils, lakes, rivers, and marine sediments show that the depositional environment from which samples were obtained has a significant and measurable impact on the combined distribution of isoprenoid and brGDGTs, which allows us to cluster the samples from our data set into environmentally relevant groups. Furthermore, we find that the distribution of GDGTs in each cluster is uniquely impacted by the environment. There is a strong association between temperature and the *Lake-type* and *Peat-type* groups. *Marine-type* samples are also clearly influenced by temperature and their distance from the coast, at least for samples more than 10 km away from the shore. The cause of this distinction is an observation that deserves further study. Although they are represented by a limited sample set, soils show variability that is



**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 LL (UL, light green), clay (C, white), and blocky lignite (BL, purple).

strongly influenced by their location. Additionally, while our analysis groups soil and river samples together into the *Soil-type* cluster, river systems seem to have more 6-methyl brGDGTs and their GDGT distributions reflect local changes within the catchment.

We used the data set presented here to train the Random Forest classification algorithm BIGMaC, which is capable of identifying the environment in which a sample was formed based on the distribution of GDGTs. Our results show that GDGTs IIa' and crenarchaeol are the most influential compounds in the classification algorithm, due to their combined unique changes between the four depositional groups. As a demonstration, we apply the BIGMaC model to an independent record from the Giraffe kimberlite, which was stratigraphically shown to record a transition from a lacustrine environment to peatland. Our BIGMaC algorithm is not only able to recreate the transition, but further suggests an excursion to peatland conditions within the upper lacustrine section of the core, which is consistent with independent evidence for more acidic conditions. This result is encouraging for the application of our classification algorithm, as it comes from a data set not included in the training or testing sets,

thus providing an independent testing case. Using the BIGMaC results as a guide, we apply brGDGT-derived calibrations specific to lakes or soils and peats as needed downcore and obtain a relatively stable temperature estimate for this area that is in general agreement with the pollen record.

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

Ultimately, we show that the combined set of branched and isoGDGTs is an effective tool for identifying depositional environments that can be used in combination with more established proxies to gain a better understanding of past environments.

### Data Availability Statement

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

### References

- 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. <https://doi.org/10.1016/j.quascirev.2021.107263>
- 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., et al. (2022). The production of diverse brGDGTs by an acidobacterium providing a physiological basis for paleoclimate proxies. *Geochimica et Cosmochimica Acta*, 337, 155–165. <https://doi.org/10.1016/j.gca.2022.08.033>
- Collinson, M. E., Steart, D. C., Harrington, G. J., Hooker, J. J., Scott, A. C., Allen, L. O., et al. (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. <https://doi.org/10.1080/00173130802707980>
- Cramwinckel, M. J., Huber, M., Kocken, I. J., Agnini, C., Bijl, P. K., Bohaty, S. M., et al. (2018). Synchronous tropical and polar temperature evolution in the Eocene. *Nature*, 559(7714), 382–386. <https://doi.org/10.1038/s41586-018-0272-2>
- Dang, X., Ding, W., Yang, H., Pancost, R. D., Naafs, B. D. A., Xue, J., et al. (2018). Different temperature dependence of the bacterial brGDGT isomers in 35 Chinese lake sediments compared to that in soils. *Organic Geochemistry*, 119, 72–79. <https://doi.org/10.1016/j.orggeochem.2018.02.008>
- 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. <https://doi.org/10.1016/j.gca.2014.06.013>
- 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. <https://doi.org/10.1016/j.gca.2019.09.043>
- 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. <https://doi.org/10.1594/PANGAEA.907818>
- 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. <https://doi.org/10.1016/j.orggeochem.2019.07.006>
- 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. *Geochimica et Cosmochimica Acta*, 125, 476–491. <https://doi.org/10.1016/j.gca.2013.10.031>
- 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. [https://doi.org/10.1016/0031-9422\(77\)80105-2](https://doi.org/10.1016/0031-9422(77)80105-2)
- 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. [https://doi.org/10.1016/0005-2760\(83\)90014-0](https://doi.org/10.1016/0005-2760(83)90014-0)
- 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. <https://doi.org/10.5194/cp-16-2599-2020>
- El Boucheffy, K., & de Souza, R. S. (2020). Learning in big data: Introduction to machine learning. In *Knowledge discovery in big data from astronomy and earth observation* (pp. 225–249). Elsevier.
- Engle, M. A., & Brunner, B. (2019). Considerations in the application of machine learning to aqueous geochemistry: Origin of produced waters in the northern US Gulf Coast Basin. *Applied Computing and Geosciences*, 3, 100012. <https://doi.org/10.1016/j.acags.2019.100012>

### Acknowledgments

We would like to thank Patrick Murphy for his assistance with the lipid analysis, Dr. Jeffrey Donnelly and the Woods Hole Oceanographic Institution Seafloor Samples Laboratory for access to marine sediment samples, and Dr. Cody Routson for contributing Alaskan lake samples. The hyperparameter tuning of the models was performed using the Ocelote cluster from the University of Arizona. This research was funded by the American Chemical Society Petroleum Research Fund, Grant 60772-ND2, and by CONACYT through the student scholarship 440897. Ioana Stefanescu and Bryan Shuman acknowledge support from the Microbial Ecology Collaborative Project through the National Science Foundation grant EPS-1655726. Francien Peterse acknowledges funding from the Nederlandse Organisatie voor Wetenschappelijk Onderzoek (NWO) through Veni Grant 863.13.016 and Vidi Grant 192.074. Lina Pérez-Ángel and Julio Sepúlveda acknowledge support from NSF Sedimentary Geology and Paleobiology Grant 1929199. We also thank Serhiy Buryak for assisting with the sampling of the Giraffe pipe sediments.



- Fleming, L. E., & Tierney, J. E. (2016). An automated method for the determination of the  $TEX_{86}$  and paleotemperature indices. *Organic Geochemistry*, 92, 84–91. <https://doi.org/10.1016/j.orggeochem.2015.12.011>
- Greenwell, B., Boehmke, B., & Gray, B. (2020). Variable Importance Plots—An Introduction to the vip Package. *The R Journal*, 12(1), 343–366. <https://doi.org/10.32614/rj-2020-013>
- Guo, J., Glendell, M., Meersmans, J., Kirkels, F., Middelburg, J. J., & Peterse, F. (2020a). Assessing branched tetraether lipids as tracers of soil organic carbon transport through the Carminowe Creek catchment (southwest England). *Biogeosciences*, 17(12), 3183–3201. <https://doi.org/10.5194/bg-17-3183-2020>
- Guo, J., Glendell, M., Meersmans, J., Kirkels, F. M. S. A., Middelburg, J. J., & Peterse, F. (2020b). Branched tetraether lipids in Carminowe Creek catchment (southwest England) [Dataset]. PANGAEA. <https://doi.org/10.1594/PANGAEA.918523>
- Guo, J., Ma, T., Liu, N., Zhang, X., Hu, H., Ma, W., et al. (2021). Branched tetraether lipids and bacterial communities along an aridity soil transect in Inner Mongolia, northern China. [Dataset]. PANGAEA. <https://doi.org/10.1594/PANGAEA.938067>
- Guo, J., Ma, T., Liu, N., Zhang, X., Hu, H., Ma, W., et al. (2022). Soil pH and aridity influence distributions of branched tetraether lipids in grassland soils along an aridity transect. *Organic Geochemistry*, 164, 104347. <https://doi.org/10.1016/j.orggeochem.2021.104347>
- Guo, J., Yuan, H., Song, J., Li, X., Duan, L., Li, N., & Wang, Y. (2022). Influence of bottom seawater oxygen on archaeal tetraether lipids in sediments: Implications for archaeal lipid-based proxies. *Marine Chemistry*, 244, 104138. <https://doi.org/10.1016/j.marchem.2022.104138>
- Halamka, T. A., McFarlin, J. M., Younkin, A. D., Depoy, J., Dildar, J., & Kopf, S. H. (2021). Oxygen limitation can trigger the production of branched GDGTs in culture. *Geochemical Perspectives Letters*, 19, 36–39. <https://doi.org/10.7185/geochemlet.2132>
- Halamka, T. A., Raberg, J. H., McFarlin, J. M., Younkin, A. D., Mulligan, C., Liu, X.-L., & Kopf, S. H. (2022). Production of diverse brGDGTs by *Acidobacterium Solibacter usitatus* in response to temperature, pH, and  $O_2$  provides a culturing perspective on br GDGT proxies and biosynthesis. *Geobiology*.
- Hamblin, A., Stasiuk, L., Sweet, A., Lockhart, G., Dyck, D., Jagger, K., & Snowdon, L. (2003). Post-kimberlite Eocene strata within a crater basin, Lac de Gras, Northwest Territories, Canada. *International Kimberlite Conference: Extended abstracts*, 8.
- Hopmans, E. C., Schouten, S., & Damsté, J. S. S. (2016). The effect of improved chromatography on GDGT-based palaeoproxies. *Organic Geochemistry*, 93, 1–6. <https://doi.org/10.1016/j.orggeochem.2015.12.006>
- Hopmans, E. C., Weijers, J. W., Schefuß, E., Herfort, L., Damsté, J. S. S., & Schouten, S. (2004). A novel proxy for terrestrial organic matter in sediments based on branched and isoprenoid tetraether lipids. *Earth and Planetary Science Letters*, 224(1–2), 107–116. <https://doi.org/10.1016/j.epsl.2004.05.012>
- Huguet, C., Hopmans, E. C., Febo-Ayala, W., Thompson, D. H., Sinninghe Damsté, J. S., & Schouten, S. (2006). An improved method to determine the absolute abundance of glycerol dibiphytanyl glycerol tetraether lipids. *Organic Geochemistry*, 37(9), 1036–1041. <https://doi.org/10.1016/j.orggeochem.2006.05.008>
- Inglis, G. N., Farnsworth, A., Collinson, M. E., Carmichael, M. J., Naafs, B. D. A., Lunt, D. J., et al. (2019). Terrestrial environmental change across the onset of the PETM and the associated impact on biomarker proxies: A cautionary tale [Dataset]. PANGAEA. <https://doi.org/10.1594/PANGAEA.901285>
- Inglis, G. N., Farnsworth, A., Collinson, M. E., Carmichael, M. J., Naafs, B. D. A., Lunt, D. J., et al. (2019). Terrestrial environmental change across the onset of the PETM and the associated impact on biomarker proxies: A cautionary tale. *Global and Planetary Change*, 181, 102991. <https://doi.org/10.1016/j.gloplacha.2019.102991>
- Jahren, A. H., & Sternberg, L. S. L. (2003). Humidity estimate for the middle Eocene Arctic rain forest. *Geology*, 31(5), 463–466. [https://doi.org/10.1130/0091-7613\(2003\)031<0463:heftme>2.0.co;2](https://doi.org/10.1130/0091-7613(2003)031<0463:heftme>2.0.co;2)
- Kim, J.-H., Van der Meer, J., Schouten, S., Helmke, P., Willmott, V., Sangiorgi, F., et al. (2010). New indices and calibrations derived from the distribution of crenarchaeal isoprenoid tetraether lipids: Implications for past sea surface temperature reconstructions. *Geochimica et Cosmochimica Acta*, 74(16), 4639–4654. <https://doi.org/10.1016/j.gca.2010.05.027>
- Kirkels, F. M., Ponton, C., Galy, V., West, A. J., Feakins, S. J., & Peterse, F. (2020). From Andes to Amazon: Assessing branched tetraether lipids as tracers for soil organic carbon in the Madre de Dios River system. *Journal of Geophysical Research: Biogeosciences*, 125(1), e2019JG005270. <https://doi.org/10.1029/2019jg005270>
- Kirkels, F. M., Usman, M. O., & Peterse, F. (2022). Distinct sources of bacterial branched GMGTs in the Godavari River basin (India) and Bay of Bengal sediments. *Organic Geochemistry*, 167, 104405. <https://doi.org/10.1016/j.orggeochem.2022.104405>
- Kirkels, F. M., Zwart, H. M., Usman, M. O., Hou, S., Ponton, C., Giosan, L., et al. (2022). From soil to sea: Sources and transport of organic carbon traced by tetraether lipids in the monsoonal Godavari River, India. *Biogeosciences*, 19(17), 3979–4010. <https://doi.org/10.5194/bg-19-3979-2022>
- Kuhn, M., & Wickham, H. (2020). Tidymodels: A collection of packages for modeling and machine learning using tidyverse principles. [Computer software manual]. Retrieved from <https://www.tidymodels.org>
- Kumar, D. M., Woltering, M., Hopmans, E. C., Damsté, J. S. S., Schouten, S., & Werne, J. P. (2019). The vertical distribution of Thaumarchaeota in the water column of Lake Malawi inferred from core and intact polar tetraether lipids. *Organic Geochemistry*, 132, 37–49. <https://doi.org/10.1016/j.orggeochem.2019.03.004>
- Langworthy, T. A. (1977). Long-chain diglycerol tetraethers from *Thermoplasma acidophilum*. *Biochimica et Biophysica Acta (BBA) - Lipids and Lipid Metabolism*, 487(1), 37–50. [https://doi.org/10.1016/0005-2760\(77\)90042-x](https://doi.org/10.1016/0005-2760(77)90042-x)
- Lauretano, V., Kennedy-Asser, A. T., Korasidis, V. A., Wallace, M. W., Valdes, P. J., Lunt, D. J., et al. (2021). Eocene to oligocene terrestrial southern hemisphere cooling caused by declining  $pCO_2$ . *Nature Geoscience*, 14(9), 659–664. <https://doi.org/10.1038/s41561-021-00788-z>
- Li, J., Pancost, R. D., Naafs, B. D. A., Yang, H., Zhao, C., & Xie, S. (2016). Distribution of glycerol dialkyl glycerol tetraether (GDGT) lipids in a hypersaline lake system. *Organic Geochemistry*, 99, 113–124. <https://doi.org/10.1016/j.orggeochem.2016.06.007>
- Liu, X., Lipp, J. S., & Hinrichs, K.-U. (2011). Distribution of intact and core GDGTs in marine sediments. *Organic Geochemistry*, 42(4), 368–375. <https://doi.org/10.1016/j.orggeochem.2011.02.003>
- Locarnini, M., Mishonov, A., Baranova, O., Boyer, T., Zweng, M., Garcia, H., et al. (2018). World Ocean Atlas 2018, volume 1: Temperature. Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M., Hornik, K., Studer, M., et al. (2019). “Finding groups in data”: Cluster analysis extended Rousseeuw et al, (Vol. 2). R package version.
- Martínez-Sosa, P., Tierney, J., Pérez-Angel, L., Stefanescu, I. C., Guo, J., Kirkels, F., et al. (2023). BIGMaC GDGT algorithm [Software]. Zenodo. (V. 1.0). <https://doi.org/10.5281/zenodo.7513557>
- Martínez-Sosa, P., Tierney, J., Pérez-Angel, L., Stefanescu, I. C., Guo, J., Kirkels, F., et al. (2023a). Giraffe kimberlite pipe core GDGT fractional abundance [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.7540094>
- Martínez-Sosa, P., Tierney, J., Pérez-Angel, L. C., Stefanescu, I. C., Guo, J., Kirkels, F., et al. (2023b). Environmental data and fractional abundance of ISO and branched GDGT data used to train the BIGMaC algorithm [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.7522415>
- Martínez-Sosa, P., Tierney, J. E., & Meredith, L. K. (2020). Controlled lacustrine microcosms show a brGDGT response to environmental perturbations. *Organic Geochemistry*, 104041.

- Martínez-Sosa, P., Tierney, J. E., Stefanescu, I. C., Crampton-Flood, E. D., Shuman, B. N., & Routson, C. (2021). A global Bayesian temperature calibration for lacustrine brGDGTs. *Geochimica et Cosmochimica Acta*, 305, 87–105. <https://doi.org/10.1016/j.gca.2021.04.038>
- Menze, B. H., Kelm, B. M., Masuch, R., Himmelreich, U., Bachert, P., Petrich, W., & Hamprecht, F. A. (2009). A comparison of random forest and its GINI importance with standard chemometric methods for the feature selection and classification of spectral data. *BMC Bioinformatics*, 10, 1–16. <https://doi.org/10.1186/1471-2105-10-213>
- Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., & Leisch, F. (2020). e1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071), TU Wien, 2018, R package version 1.7-0.
- Naafs, B., Rohrsen, M., Inglis, G. N., Lähteenoja, O., Feakins, S. J., Collinson, M. E., et al. (2018). High temperatures in the terrestrial mid-latitudes during the early Palaeogene. *Nature Geoscience*, 11(10), 766–771. <https://doi.org/10.1038/s41561-018-0199-0>
- Naafs, B. D. A. (2017). Global biomarker (GDGT) database for peatlands [Dataset]. PANGAEA. <https://doi.org/10.1594/PANGAEA.883765>
- Naafs, B. D. A., Inglis, G. N., Zheng, Y., Amesbury, M., Biester, H., Bindler, R., et al. (2017). Introducing global peat-specific temperature and pH calibrations based on brGDGT bacterial lipids. *Geochimica et Cosmochimica Acta*, 208, 285–301. <https://doi.org/10.1016/j.gca.2017.01.038>
- Naeher, S., Peterse, F., Smittenberg, R. H., Niemann, H., Zigah, P. K., & Schubert, C. J. (2014). Sources of glycerol dialkyl glycerol tetraethers (GDGTs) in catchment soils, water column and sediments of Lake Rotsee (Switzerland)—Implications for the application of GDGT-based proxies for lakes. *Organic Geochemistry*, 66, 164–173. <https://doi.org/10.1016/j.orggeochem.2013.10.017>
- Parmar, A., Kataria, R., & Patel, V. (2019). A review on random forest: An ensemble classifier. In *International Conference on Intelligent Data Communication Technologies and Internet of Things (ICICI) 2018* (pp. 758–763).
- Peaple, M. D., Tierney, J. E., McGee, D., Lowenstein, T. K., Bhattacharya, T., & Feakins, S. J. (2021). Identifying plant wax inputs in lake sediments using machine learning. *Organic Geochemistry*, 156, 104222. <https://doi.org/10.1016/j.orggeochem.2021.104222>
- Pérez-Angel, L. C., Sepúlveda, J., Molnar, P., Montes, C., Rajagopalan, B., Snell, K., et al. (2020a). In situ temperature and brGDGTs measurements in soils from the Tropical Andes of Colombia and a tropical soil brGDGT compilation dataset [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.3939270>
- Pérez-Angel, L. C., Sepúlveda, J., Molnar, P., Montes, C., Rajagopalan, B., Snell, K., et al. (2020b). Soil and air temperature calibrations using branched GDGTs for the Tropical Andes of Colombia: Toward a pan-tropical calibration. *Geochemistry, Geophysics, Geosystems*, 21(8), e2020GC008941. <https://doi.org/10.1029/2020gc008941>
- R Core Team. (2022). R: A language and environment for statistical computing. [Computer software manual]. Retrieved from <https://www.R-project.org/>
- Raberg, J. H., Miller, G. H., Geirsdóttir, Á., & Sepúlveda, J. (2022). Near-universal trends in brGDGT lipid distributions in nature. *Science Advances*, 8(20), eabm7625. <https://doi.org/10.1126/sciadv.abm7625>
- Rattanasriampaipong, R., Zhang, Y. G., Pearson, A., Hedlund, B. P., & Zhang, S. (2022). Archaeal lipids trace ecology and evolution of marine ammonia-oxidizing archaea. *Proceedings of the National Academy of Sciences*, 119(31), e2123193119. <https://doi.org/10.1073/pnas.2123193119>
- Revelle, W., & Revelle, M. W. (2015). Package ‘psych’. The comprehensive R archive network, (pp. 337–338).
- 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. <https://doi.org/10.1016/j.orggeochem.2012.09.006>
- Schouten, S., Hopmans, E. C., Schefuß, E., & Damsté, 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. [https://doi.org/10.1016/s0012-821x\(02\)00979-2](https://doi.org/10.1016/s0012-821x(02)00979-2)
- 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. <https://doi.org/10.1016/j.orggeochem.2018.06.005>
- Sinninghe Damsté, J. S., Rijpstra, W. I. C., Hopmans, E. C., Weijers, J. W., Foesel, 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. <https://doi.org/10.1128/aem.00466-11>
- 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. <https://doi.org/10.1194/jlr.m200148-jlr200>
- 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. <https://doi.org/10.1016/j.quascirev.2021.107352>
- 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. <https://doi.org/10.1017/jpa.2022.101>
- 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. <https://doi.org/10.1016/j.gloplacha.2013.06.011>
- 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. *Geochimica et Cosmochimica Acta*, 74(17), 4902–4918. <https://doi.org/10.1016/j.gca.2010.06.002>
- 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. <https://doi.org/10.1029/2017gc007401>
- Véquaud, P., Thibault, A., Derenne, S., Anquetil, C., Collin, S., Contreras, S., et al. (2022). FROG: A global machine-learning temperature calibration for branched GDGTs in soils and peats. *Geochimica et Cosmochimica Acta*, 318, 468–494. <https://doi.org/10.1016/j.gca.2021.12.007>
- Weijers, J. W., Schouten, S., Hopmans, E. C., Geenevasen, J. A., David, O. R., Coleman, J. M., et al. (2006). Membrane lipids of mesophilic anaerobic bacteria thriving in peats have typical archaeal traits. *Environmental Microbiology*, 8(4), 648–657. <https://doi.org/10.1111/j.1462-2920.2005.00941.x>
- 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. <https://doi.org/10.1016/j.gca.2006.10.003>
- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., et al. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/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 years. *Earth and Planetary Science Letters*, 516, 66–76. <https://doi.org/10.1016/j.epsl.2019.03.038>

- Wolfe, A. P., Reyes, A. V., Royer, D. L., Greenwood, D. R., Doria, G., Gagen, M. H., et al. (2017). Middle Eocene CO<sub>2</sub> and climate reconstructed from the sediment fill of a subarctic kimberlite maar. *Geology*, *45*(7), 619–622. <https://doi.org/10.1130/g39002.1>
- Wright, M. N., Wager, S., Probst, P., & Wright, M. M. N. (2019). Package 'ranger' version 0.11, 2.
- Yavitt, J., Harms, K., Garcia, M., Wright, S., He, F., & Mirabello, M. (2009). Spatial heterogeneity of soil chemical properties in a lowland tropical moist forest, Panama. *Soil Research*, *47*(7), 674–687. <https://doi.org/10.1071/sr08258>
- Zell, C., Kim, J.-H., Moreira-Turcq, P., Abril, G., Hopmans, E. C., Bonnet, M.-P., et al. (2013). Disentangling the origins of branched tetraether lipids and crenarchaeol in the lower Amazon River: Implications for GDGT-based proxies. *Limnology & Oceanography*, *58*(1), 343–353. <https://doi.org/10.4319/lo.2013.58.1.0343>
- Zhang, Y. G., Zhang, C. L., Liu, X.-L., Li, L., Hinrichs, K.-U., & Noakes, J. E. (2011). Methane index: A tetraether archaeal lipid biomarker indicator for detecting the instability of marine gas hydrates. *Earth and Planetary Science Letters*, *307*(3–4), 525–534. <https://doi.org/10.1016/j.epsl.2011.05.031>
- Zhao, B., Castañeda, I. S., Salacup, J. M., Thomas, E. K., Daniels, W. C., Schneider, T., et al. (2022). Prolonged drying trend coincident with the demise of Norse settlement in southern Greenland. *Science Advances*, *8*(12), eabm4346. <https://doi.org/10.1126/sciadv.abm4346>
- Zheng, Y., Heng, P., Conte, M. H., Vachula, R. S., & Huang, Y. (2019). Systematic chemotaxonomic profiling and novel paleotemperature indices based on alkenones and alkenoates: Potential for disentangling mixed species input. *Organic Geochemistry*, *128*, 26–41. <https://doi.org/10.1016/j.orggeochem.2018.12.008>
- Zheng, Y., Liu, H., Yang, H., Wang, H., Zhao, W., Zhang, Z., et al. (2022). Decoupled Asian monsoon intensity and precipitation during glacial-interglacial transitions on the Chinese Loess Plateau. *Nature Communications*, *13*(1), 5397. <https://doi.org/10.1038/s41467-022-33105-2>
- Zheng, Y., Pancost, R. D., Liu, X., Wang, Z., Naafs, B., Xie, X., et al. (2017). Atmospheric connections with the north Atlantic enhanced the deglacial warming in northeast China. *Geology*, *45*(11), 1031–1034. <https://doi.org/10.1130/g39401.1>