# Utilizing Probability Estimates from Machine Learning and Pollen to Understand the Depositional Influences on Branched GDGT in Wetlands, Peatlands, and Lakes

Amy Cromartie[1], Cindy De Jonge[2], Guillemette Ménot[3], Mary Robles[4,5], Lucas Dugerdil[3,5], Odile
Peyron[5], Marta Rodrigo-Gámiz[6], Jon Camuera[7], Maria Jose Ramos-Roman[8], Gonzalo Jiménez-Moreno[6], Claude Colombié[9], Lilit Sahakyan[10], Sébastien Joannin[5]

1 Université Côte d'Azur, CNRS, CEPAM, UMR 7264, 06300, Nice, France

2 Geological Institute, ETH Zürich, 8092, Zurich, Switzerland

3 ENS de Lyon, Université Lyon 1, CNRS, UMR 5276 LGL-TPE, F-69364, Lyon, France

4 Aix-Marseille Univ., CNRS, IRD, INRAE, Coll France, UMR 34 CEREGE, 13545, Aix-en-Provence, France

5 ISEM, Univ Montpellier, CNRS, IRD, Montpellier, France

6 Department of Stratigraphy and Paleontology, University of Granada, 18071, Granada, Spain

7 Unit of Botany, Faculty of Pharmacy, Complutense University of Madrid, Spain

8 Organismal and Evolutionary Biology Research Program, Research Centre for Ecological Change, University of Helsinki, Finland

9 Univ. Lyon, UCBL, ENSL, UJM, CNRS, LGL-TPE, F-69622, Villeurbanne, France

10 Institute of Geological Sciences, National Academy of Sciences of Republic of Armenia, Yerevan, Armenia

*Correspondence to*: Amy Cromartie (aec277@cornell.edu)

**Abstract.** Branched glycerol dialkyl glycerol tetraethers (brGDGTs) serve as critical molecular biomarkers for the quantitative reconstruction of past environments, ambient temperature and pH across various archives. Despite their success, numerous issues persist that limit their application. The distribution of brGDGTs varies significantly based on provenance, resulting in biases in environmental reconstructions that rely on fractional abundances and derived indices, such as the MBT'$_{5ME}$. This issue is especially significant in shallow lakes, wetlands, and peatlands within semi-arid and arid regions, where ecosystems are sensitive to diverse environmental and climatic factors. Recent advancements, such as machine learning techniques, have been developed to identify changes in sources; however, these techniques are insufficient for detecting mixed source environments. The probability estimates derived from five machine learning algorithms are employed here to detect provenance changes in brGDGT downcore records and to identify periods of mixed provenance. A new global modern database (n=2031) was compiled to train, validate, test, and apply these algorithms to two sedimentary

records. Our findings are corroborated by pollen and non-pollen palynomorphs obtained from the identical records. These microfossil proxies are utilized to discuss changes in provenance, hydrology, and ecology that influence the distribution of brGDGTs. Probability estimates derived from Random Forest with a sigmoid calibration are most effective in detecting changes in brGDGT distribution. Minor changes in the relative contributions of brGDGTs provenance can significantly

35 influence the distribution of brGDGTs, especially regarding the MBT'$_{5ME}$ index. This study introduces a novel brGDGT wetland index aimed at monitoring potential biases arising from wetland development.

## 1 Introduction

Branched glycerol dialkyl glycerol tetraethers (brGDGTs), first identified in peat sequences (Weijers et al., 2006), have

40 demonstrated significant potential as a quantitative proxy for paleo-environmental reconstructions. The ubiquity of brGDGTs and their global correlations with temperature and pH, notably across different archive types, positions them as a valuable tool for paleoclimate reconstructions (among others Weijers et al., 2007; Peterse et al., 2012; Loomis et al., 2012; Raberg et al., 2022). Researchers have identified brGDGTs across various depositional environments, such as peat, soils, loess, fossilized bones, and lacustrine, marine, and river sediments (e.g., Weijers et al., 2006; 2007; De Jonge et al., 2014a;

45 Warden et al., 2016; Naafs et al., 2017a,b; Dillon et al., 2018; Baker et al., 2019) at differing geological timescales, indicating their widespread potential as a proxy for reconstructing continental paleoclimate.

Despite its potential, a key challenge in utilizing brGDGT-based reconstructions in continental settings is the temperature independent variability in fractional abundance (FA) distribution across these environments. In the context of lacustrine, wetland, and peat archives, the fractional abundance of brGDGTs produced in the aquatic environments and surrounding soils varies (Tierney and Russell, 2009, Tierney et al., 2010, Zink et al., 2010, Buckles et al., 2014, Loomis et

50 al., 2011, Loomis et al., 2012, Loomis et al., 2014a, Loomis et al., 2014b, Li et al., 2016, Russell et al., 2018; Dang et al., 2018). Potential changes in provenance thus result in distribution differences that may lead to inaccuracies in paleoenvironmental reconstructions. This includes paleotemperature reconstructions based on the widely recognized MBT'$_{5ME}$ index. This index measures the degree of methylation of the 5-methyl brGDGTs, distinguishing it from the 6-

55 methyl brGDGTs to establish calibrations that exhibit a stronger correlation with mean annual air temperature (MAAT) (De Jonge et al., 2014a). The MBT'$_{5ME}$ index has been successfully utilized as grounds for various global temperature calibrations concerning lakes, peats, and soils (e.g., De Jonge et al., 2014a; Hopmans et al., 2016; Naafs et al., 2017a; Dearing Crampton-Flood et al., 2020; Martínez-Sosa et al., 2021; Véquaud et al., 2022).

Furthermore, brGDGTs distributions within these depositional environments may be influenced by distinct

60 environmental characteristics. Soil chemistry, particularly pH, can influence the 5-methyl brGDGTs (De Jonge et al., 2021; 2024). In certain lakes, the 6- methyl brGDGTs exhibit a stronger correlation with mean annual air temperature compared to the 5- methyl brGDGTs, which contrasts with the catchment soils (Dang et al., 2018). In peatlands, the MBT' and MBT'$_{5ME}$ values are higher in dry sites compared to those that are waterlogged (Rao et al., 2022). The potential differential

distributions resulting from depositional environments underscore the influence of changes in provenance or hydrological
65    conditions on brGDGT-based environmental reconstructions.

Factors influencing the distribution of brGDGTs in lakes include lake stratification and redox conditions (Weber et al., 2018), salinity (Wang et al., 2021), conductivity (Tierney et al., 2010; Raberg et al., 2022), dissolved oxygen (Wu et al., 2021), and water depth (Stefanescu et al., 2021), amongst others. In soils, vegetation, and vegetation-mediated factors such as soil temperature (Liang et al., 2019; 2023), soil moisture (Menges et al., 2014; Dang et al., 2016), precipitation (Dugerdil
70    et al., 2021a), and soil chemistry (Dang et al., 2016; De Jonge et al., 2021) all influence distributional changes. BrGDGT distributions in peat may vary in response to flooding, drying of peatlands, and alterations in the water table (Rao et al., 2022; Ofiti et al., 2024).

BrGDGT reconstruction in Quaternary downcore lacustrine records indicates that changes in depositional and mixed provenance significantly affect environmental reconstructions (i.e., Martin et al., 2019; Robles et al., 2022; Ramos-
75    Román et al., 2022; d'Oliveira et al., 2023; Acharya et al., 2023). As climatic or successional changes occur concurrently with temperature variations, isolating the effects of source changes on the MBT'$_{5ME}$ is challenging. Therefore, an alternative method for detecting minor or major provenance shifts is required. This paper presents a strategy for identifying provenance or ecological changes across lacustrine, peat, and soil depositional environments, including mixed contexts, utilizing a new global brGDGT database, machine learning techniques, as well as environmental reconstructions based on pollen and non-
80    pollen palynomorphs.

Two approaches are employed to achieve this objective. First, we use probability estimates derived from machine learning to identify changes in sourcing over time. This study extends the work of Martínez-Sosa et al., (2023), who employed supervised machine learning to identify changes in brGDGT sources using classification models based on modern samples. Rather than employing discrete classification, we utilize the probability estimates from these classification
85    algorithms to analyze the contributions from differential sources at any specific time. This method enhances prior approaches by recognizing environments that integrate brGDGTs from multiple inputs and depositional settings that have not fully transitioned to a new depositional state. The probability estimates are derived from the classification of modern samples (n=2301), categorized into three groups: soil, peat, and lake, utilizing both previously published and new datasets. We implement five algorithms: K-Nearest Neighbor, Support Vector Machines (SVM), Logistic Regression (LR),
90    Classification and Regression Trees (CART), and Random Forest (RF). These are employed using Python and scikit-learn to identify intervals where downcore records are predominantly influences by in-situ lake brGDGTs, mineral soils, and peatlands, as well as combinations of these elements.

Secondly, to ensure accurate identification of provenance changes, comparisons are conducted with published pollen and non-pollen palynomorphs (NPPs) from extensive sediment records and variations in brGDGT distribution (i.e.,
95    Robles et al., 2022; Camuera et al., 2018;2019; Ramos-Román et al., 2018; Rodrigo-Gámiz et al., 2022). The records are situated in the semi-arid mid-latitude zones, where water bodies are subject to temporal variations. Utilizing these proxies allows for an independent comparison of outputs to: i.) confirm machine learning results through the integration of brGDGT-

based reconstructions with pollen and NPPs; ii.) demonstrate how these complementary proxies can aid in identifying potential hydrological, ecological, and depositional shifts that may introduce bias in brGDGT reconstructions. This study

100 demonstrates that alterations in provenance and hydrology can significantly influence the distribution of brGDGTs and, consequently, established indices like MBT'$_{5ME}$, while also offering novel methodologies for identifying changes in paleorecords.

## 2 Materials and Methods

105 ### 2.1 GDGT databases

### 2.1.1 Building a new modern sample database

This study compiles published brGDGT databases for lake (n=591), soil (n = 1197), and peat (n=532) depositional categories (Baxter et al. 2019; Cao et al. 2020; Chen et al., 2021; Dang et al., 2018 ; Dearing Crampton-Flood et al., 2020; De Jonge et al., 2014b; Ding et al., 2015 ; Dugerdil et al., 2021 ; Guo et al. 2020 ; Halffman et al. 2022 ; Jaeschke et al. 2018 ; Kirkels et

110 al. 2020 ; Kou et al., 2022 ; Li et al., 2016 ; Li et al., 2018 ; Martin et al., 2019; 2020 ; Martínez-Sosa et al., 2021 ; Naafs et al., 2017b ; Ning et al., 2019 ; Pérez-Angel et al., 2020 ; Qian et al., 2019 ; Raberg et al., 2021 ; Robles et al., 2021 ; Rao et al., 2022 ; Russell et al., 2018 ; Stefanescu et al., 2021 ; Véquaud et al., 2021a ; 2021b ; Wang et al., 2016 ; 2018 ; 2020a ; 2020b ; 2021 ; Weber et al. 2018 ; Wu et al. 2020 ; Xiao et al., 2015 ; Yang et al., 2015 ; Yao et al., 2020 ; Fig. 1 Map; full data on https://github.com/amycromartie/ProbbrGDGT). Round robin test results show that results from multiple laboratories

115 can be integrated into a single database (De Jonge et al., 2024). Results were included only when chromatography enabled the separate quantification of 5- and 6-methyl brGDGTs (i.e., De Jonge et al., 2014b). The fractional abundances of fifteen distinct brGDGT structural isomers were sourced from the original authors or recalculated from the initial datasets (https://github.com/amycromartie/ProbbrGDGT). We enhanced the training dataset for certain published datasets by obtaining data with greater precision from the original authors, where fractional abundances had been rounded two decimal

120 places. We incorporated the fractional abundances of individual downcore samples from Naafs et al., (2017a) to enhance the sample size of our peat analysis. This facilitated the development of a more robust model for assessing brGDGT distribution across various types. All samples originate from terrestrial environments (Fig. 1). The 7-methyl (Ding et al., 2016) or the ⅚ isomer, also referred to as IIIa" (Weber et al., 2015) were excluded due to limited publication. The original author's description was utilized to categorize the samples into a classification index (i.e., soil, lake, peat). Suspended particulate

125 matter (SPM), moss polsters, marine, and river samples were excluded. Latitude and longitude data were converted to decimal degrees as required.

### 2.1.2. Addition of new samples from Armenia

Thirty new surface samples from the country of Armenia were added to the global dataset to expand the database for semi-

130 arid environments. Nine samples were collected from wetlands at a depth of 0-2 cm, one sample from Lake Sevan at a depth of 2-3 cm, as previously discussed in Robles et al., (2022), along with 20 surface soil samples. For brGDGT extraction, each
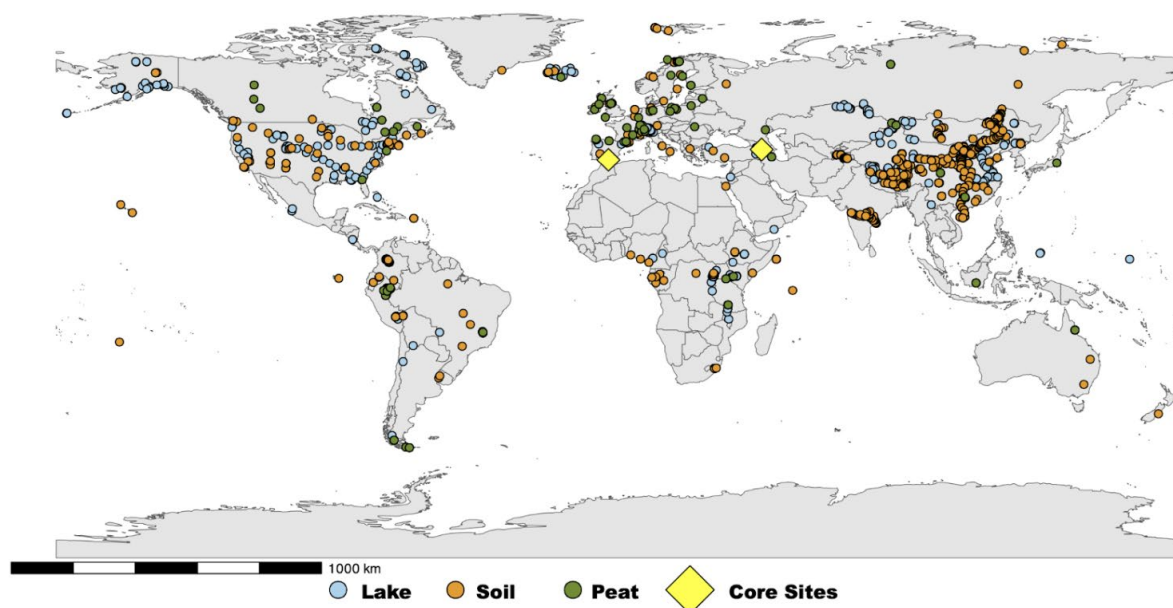
sample was first lyophilized, freeze-dried, and ground. Lipids were extracted from 0.5 to 1g of sample in two rounds with a

MARS 6 CEM microwave, using a 3:1 mixture of dichloromethane (DCM) and methanol (MeOH) (3:1). These samples

were filtered with a silicon SPE cartridge with a mixture of Hexane:DCM (1:1) and then DCM:MeOH (1:1) to separate the

135    apolar and polar fraction, respectively. A standard of $C_{46}$, following Huguet et al. (2006), was added to the total lipid extract

(TLE) prior to separating the fraction. The polar fraction was then analyzed on high-performance liquid chromatography

with atmospheric pressure chemical ionization mass spectrometry (HPLC-APCI-MS, Agilent 1200) at LGL-TPE ENS,

which allows separation of the 5- and 6-methyl GDGTs, following Hopmans et al., (2016). Selective ions monitoring (SIM)

of m/z of 1050, 1048, 1046, 1036, 1032, 1022, 1020, 1018, and 744 was used for the brGDGTs isomers and the $C_{46}$ standard

140    (Hopmans et al., 2016; Davtian et al., 2021; Huguet et al., 2006).

### 2.1.3 Resampling and balancing modern dataset

The distribution of modern brGDGT samples across classification datasets (i.e., soils, lakes, peats) was not uniform, with a

predominance of soil samples and an under-representation of lake and peat samples (Fig. 1). Unbalanced datasets can result

145    in considerable performance issues, such as misclassification of data with limited sample sizes, which may prevent the

learning algorithm from identifying general patterns within the datasets (He and Garcia, 2009). Consequently, a combination

of downsampling and upsampling techniques was utilized for model comparisons (Fig. 2). This involved evaluating each

machine learning model using both the raw and resampled datasets, which incorportaed upsampled synthetic samples. In the

resampled dataset, we initially preformed random downsampling of the soil samples in R to achieve a sample size of 750

150    from the original dataset. The Synthetic Minority Oversampling Technique (SMOTE) function from the R library

smotefamily (Siriseriwan, 2019) was employed to upsample the peat and lake dataset. The SMOTE function is an

oversampling technique that selects a sample from the minority dataset, identifies its nearest neighbor(s), and generates a

new data point between the original pair (Siriseriwan, 2019). SMOTE was utilized to generate 1000 synthetic samples for the

lake and peat datasets derived from the original datasets. Samples were randomly selected from the synthetic dataset to

155    adjust the raw datasets for lake and peat to a total of 750. In the case of peat and lake samples, 219 and 159 synthetic

samples were incorporated into the raw dataset, respectively.



**Figure 1:** Map of modern sample locations used in the compiled database alongside the two sites designated for paleo-reconstructions. Map created with R package ggplot2 (Wickham, 2016)

160

## 2.2 Machine Learning models

### 2.2.1 Building probability and classification machines

In supervised classification problems, machine learning algorithms utilize grouped attributes and features to identify patterns within human-curated datasets (Kalita, 2022). Samples in these datasets are typically assigned a label (class), target value, or
165  dependent variable, which correspond to independent variables and features. The model utilizes this information to understand the relationships between the independent and dependent variables during the training process (Geetha and Sendhilkumar, 2023). The models are subsequently refined and evaluated for accuracy using a subset of the known classification dataset that has not been previously encountered by the model. A distinct validation set is employed to adjust the probability estimates. Numerous classification machine learning models employ probability estimates to determine the
170  appropriate class (Murphy 2012). When calibrated, these probability estimates can provide information that extends beyond merely identifying an individual's category but can also indicating the degree of likeness of an individual belongs to a category (Malley et al., 2012). Most machine learning algorithms, when initially deployed, lack calibration for precise probability predictions. Calibration is essential to ensure that the empirical probability is both valid and accurate (Dawid

6

1985). In the absence of calibration, certain model outputs may push probability estimates toward 0 or 1, necessitating

175 correction through calibration (Niculescu-Mizil and Caruana, 2005). Typically, either Sigmoid ("Platt scaling") or Isotonic regression is employed for calibration on a validation dataset that the model has not previously encountered (ibid). Subsequent to these steps, the model may be utilized to predict a class within a dataset where the classification remains unknown.

We employed Python and scikit-learn (Pedregosa et al., 2011) for the machine learning analysis. Five commonly

180 used supervised machine learning models were evaluated: k-nearest neighbors (KNN), support vector machine (SVM), logistic regression (LR), classification and regression trees (CART), and random forests (RF). The selection of these models is based on their capacity to produce calibrated probability estimates.

Logistic regression (LR) is a parametric model akin to linear regression in its functionality, yet it is more appropriate for classification tasks (i.e., binary outcomes) (Hilbe, 2016). The analysis relies on the likelihood of an event occurring and the

185 alignment of predictor response variables within the probability distribution (Hilbe, 2016). The algorithm is inherently calibrated for precise probability outputs due its foundational reliance on probability.

K-nearest neighbors (KNN), support vector machine (SVM), classification and regression trees (CART), and random forests (RF) are non-parametric models demonstrated to be effective in probability estimation following calibration (i.e., Niculescu-Mizil and Caruana, 2005; Kruppa et al., 2013; Dankowski and Ziegler, 2016, Cearns et al., 2020). K-nearest

190 neighbor (KNN) functions as a "lazy learner" by determining the distance between data points according to the characteristics of the training dataset (Geetha and Sendhilkumar, 2023). The "K" in KNN refers to a small positive integer that determines the number of neighbors taken into account when predicting the category of a data point (ibid). KNN is extensively employed in palaeosciences for paleoclimate regression issues, particularly through the modern analog technique (Simpson, 2007). Supported Vector Machines (SVM) is a model that positions data items in n-dimensional space based on n

195 features (Geetha and Sendhilkumar, 2023). Classification is achieved by identifying a hyperplane in the dimensional space that distinguishes between the classes (ibid).

Classification and Regression Trees (CART) and Random Forests (RF) are tree-based learning algorithms. Trees are formed through three fundamental steps: (1) binary splits are selected, (2) a determination is made regarding the node is terminal or requires further splitting; and (3) a class is assigned to the terminal leaf node (Bell 1999). Random Forest (RF) is

200 founded on the principles of natural variability and randomness inherent in trees, where both the variables and the individual elements exhibit a degree of randomness (Genuer and Poggi, 2020). RF classification problems utilize a committee of decision trees that collectively vote to determine the predicted class (Hastie, 2009). In classification problems, each vote corresponds to a classification in the terminal node of the tree (Malley et al., 2012), with the majority vote determining the final classification outcome. The probability estimates are derived by calculating the fraction of votes from each tree to

205 determine the predicted class probability.

**2.2.2 Verification, tuning, and calibration of models**

The raw and SMOTE datasets were divided into training, testing and validation sets in a 60:20:20 ratio (Fig. 2). The models

210 underwent testing and hyperparameter tuning using a k-fold cross-validation approach, incorporating ten data splits and a

parameter grid with the test dataset. K-fold cross-validation involves partitioning the data into equal-sized subsets, which are

then utilized k times, with k - 1 subsets used for training and one subset reserved for validation (Jung, 2017). The

performance is evaluated by averaging each k iteration. The parameter grid facilitates the iteration over a finite set of values

to identify optimal variables for tuning. After tuning, all models and datasets were retested for accuracy. The distribution

215 was subsequently plotted, and the mean F1 accuracy results were computed (Fig. 3).
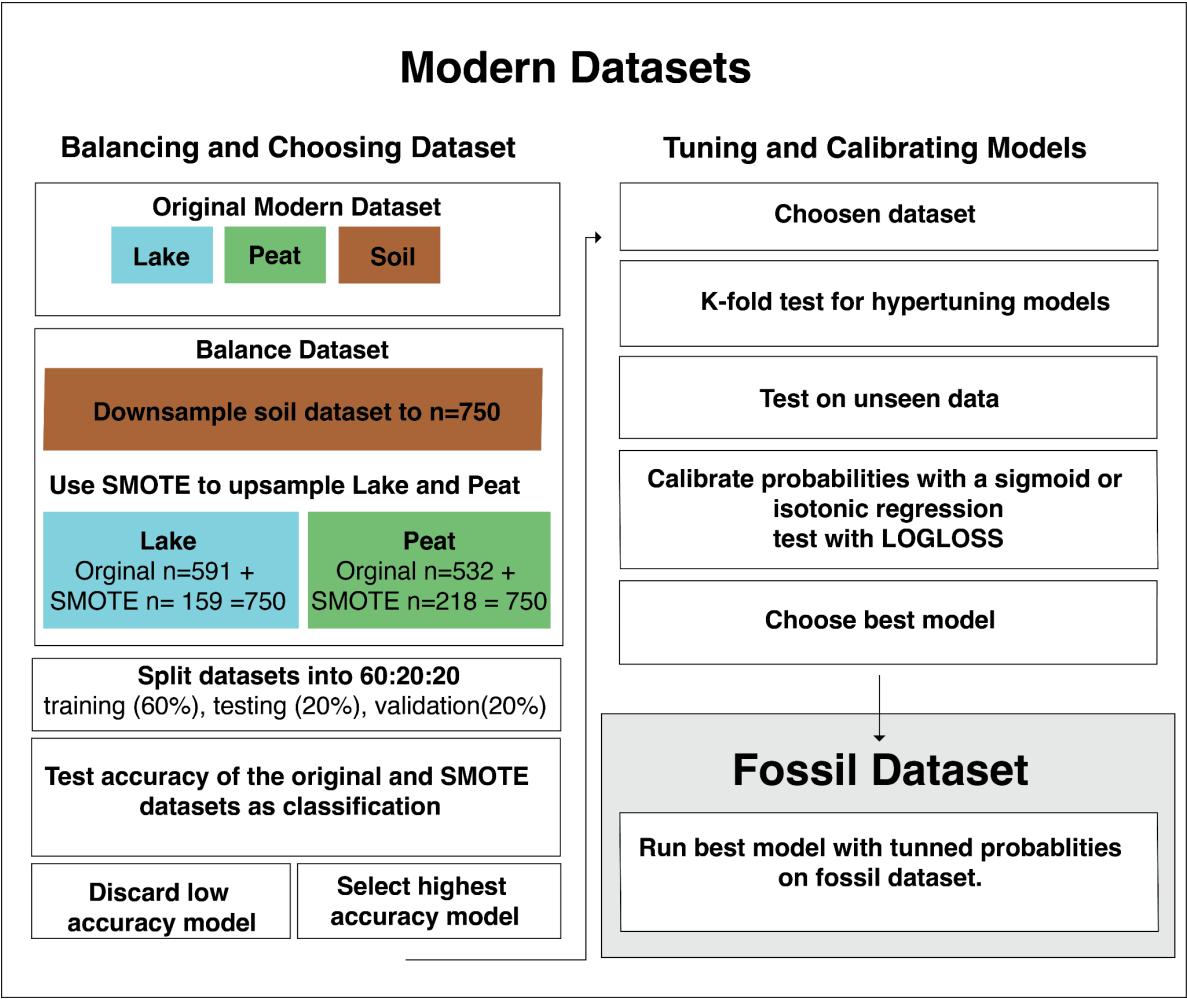


**Figure 2:** Illustration of the methods employed in this study for testing, tuning, and validating the datasets and models.

**2.2.3 Probability estimate calibration and application of classification machines**

220  Instead of solely predicting the class (e.g., soil, peat, lake), we are employing the probability estimate output generated by the classification algorithms as a proxy for environmental change. The probability output enables the estimation of the likelihood that a specific sample belongs to a particular class, thus facilitating the identification of periods of mixed sourcing. Due to the lack of calibration in the default probability estimates of the algorithms employed, we applied Sigmoid and Isotonic regression on the validation dataset to rectify any distortion and assessed the effectiveness using the Log loss

225  function in scikit-learn. Log loss is employed in probability scenarios where the likelihood of an event being true is represented as 1, equally true as 0.5, and false as 0 (Manzali et al., 2017). In Log loss, a greater divergence between the predicted value and the actual value results in a higher log-loss score (Dembla, 2020). A lower score indicates greater accuracy in predictions. Log loss scores were subsequently compared across models to evaluate performance.

230  **2.3 Application of models, downcore pollen, non-pollen palynomorph, brGDGT analysis**

To assess the accuracy of the probability estimates on the downcore record, five machine-learning models were applied to two published brGDGT records that included datasets of pollen and non-pollen palynomorphs (NPPs). Aquatic pollen and NPPs provide critical insights into alterations in lake or wetland ecology (e.g., Cromartie et al., 2020; Robles et al., 2022). We selected two records: one from Armenia in the southern Caucasus (Vanevan peat: 40°12′8.83″N, 45°40′24.03″E, Robles

235  et al., 2022) and another from southern Spain (Padul paleolake: 37°00′39′′N, 3°36′14′′W, Camuera et al., 2018; 2019; Ramos-Román et al., 2018; Rodrigo-Gámiz et al., 2022), both situated at comparable latitudes in Eurasia.

The extraction methods for brGDGTs are detailed in the original articles by Robles et al., (2022) and Rodrigo-Gámiz et al., (2022). We revisited the original chromatograms of Robles et al. (2022) to investigate the presence of the IIIa'' isomer, which had not been published previously to verify the brGDGTs based ML lake probability output. The IIIa'' isomer was

240  reported in Rodrigo-Gámiz et al., 2022.  Robles et al. (2022) provide identification and counting methods for aquatic pollen and NPPs in the Vanevan Peat, while Ramos-Román et al., (2018) and Camuera et al., (2019) address similar methods for the Padul paleolake. Additionally, we re-calculated the reconstructed water-depth based on aquatic pollen and NPPs. The analysis relies on the raw datasets employing the original equations established by Robles et al., (2022) and Camuera et al., (2019). Instead of applying a smoothing technique to the water-depth reconstruction as done by Camuera et al., (2019), we

245  retained the original sample-to-sample curve.  Percentages of the aquatic and NPP taxa were calculated by summing all relevant pollen types for each record and dividing each taxon by the total sum. We calculated and re-calculated key brGDGT-based indices (Table 1) to compare our machine learning results with the brGDGT record as well as the aquatic pollen and NPPs.

| Index | Formulae | Citation |
|-------|----------|----------|
| MBT'$_{5ME}$ | MBT'$_{5ME}$ = ([Ia] + [Ib] + [Ic]) / ([Ia] + [Ib] + [Ic] + [IIa] + [IIb] + [IIc] + [IIIa]) | De Jonge et al. 2014b |

| CBT' | CBT= $^{10}$log([Ic] + [IIa'] + [IIb'] +[IIc'] + [IIIa'] + [IIIb'] + [IIIc']) / ([Ia] + [IIa] +[IIIa]) | De Jonge et al. 2014b |
|---|---|---|

**Table 1:** Table of brGDGT indices employed in this study.

250

## 2.4 Descriptive statistics

The programming languages R (R core team) and Python (Python Software Foundation. Python Language Reference, version 3.7.3. available at http://www.python.org) were utilized alongside ggplot2 (Wickham, 2016) and matplotlib (Hunter, 2007) to visualize the results.

255        Redundancy analysis (RDA) was performed using the vegan package (Oksanen et al., 2019). RDA was employed in two capacities: i.)  To compare the fractional abundances of the brGDGTs in the global modern dataset with the author's descriptive categories (i.e., soil, peat, lake), and ii.) To compare the probability estimates results from the Vanevan and Padul records with the pollen and NPPs. In this analysis, we downsampled the pollen record to align with the brGDGT resolution, selecting samples that were no more than 100 years apart. Bayesian change-point analyses were conducted on the brGDGTs

260   based ML lake probability results using the bcp package (Erdman and Emerson, 2007;2008; Wang and Emerson, 2015) in R to identify significant shifts in depositional environments.
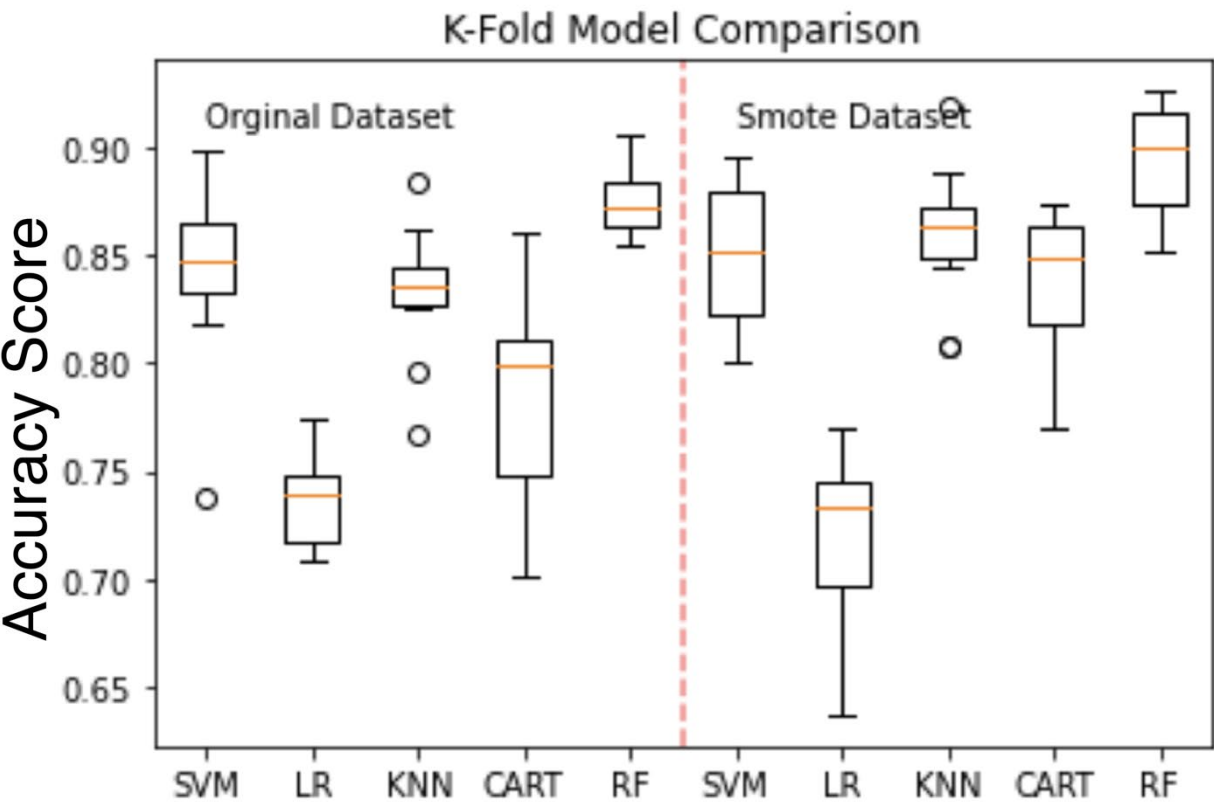

## 3 Results

### 3.1 New modern brGDGT dataset

265   The raw database compiled for this study comprised a total of 2282 samples (591 from lakes, 532 from peats, and 1177 from soils). This addition includes 319 lake samples to the database of Martínez-Sosa et al., (2021), 62 peat samples to Naafs et al., (2017b), and 450 soil samples to the Dearing-Crampton Flood et al. (2019) datasets. Subsequent to the compilation of this dataset, additional datasets have been published (e.g., Raberg et al., 2022; Martínez-Sosa et al., 2023) that are not incorporated in our dataset. Fig. 1 and Fig. S1 illustrate the distribution of the brGDGT datasets.

270

### 3.2. Model accuracy and Log loss

In classification mode, all models demonstrated mean accuracy F1 scores, which measures the predictive accuracy of the models, between 0.72 and 0.90 (Fig. 3). The study compared the performance of various classification models, with the SMOTE dataset showing superior results over the raw unbalanced dataset for SVM, KNN, CART, and RF (Table 2). The raw

275   dataset showed better performance with LR, while the SMOTE dataset improved probability estimates for SVM and RF but decreased for LR, KNN, and CART. The sigmoid calibration improved probabilities for RF, CART, KNN, but decreased for SVM and LR. The isotonic calibration improved probabilities for KNN, CART, but decreased for SVM, LR, and RF over uncalibrated probabilities. The sigmoid function outperformed the isotonic function on both datasets for SVM, LR, and RF.

The RF model with the SMOTE dataset had the highest accuracy and the lowest Log loss score for sigmoid and uncalibrated

280    probabilities. The RF model with the SMOTE dataset and sigmoid calibration was chosen as the best performing model.



**Figure 3:** Comparison of k-fold testing models between datasets, utilizing k-fold cross-validation for classification on both SMOTE and original datasets. The k-fold comparison utilized 10 splits across Supported Vector Machines (SVM), Logistic Regression (LR), K-Nearest Neighbor (KNN), Classification and Regression Tree (CART), and Random Forest (RF).

285

| Model and Database | F1 Mean Accuracy | Standard Deviation | LogLoss uncalibrated | LogLoss Sigmoid | LogLoss Isotonic |
|---|---|---|---|---|---|
| SVM | 0.84 | 0.04 | 0.46 | 0.51 | 0.93 |
| SVM_SMOTE | 0.85 | 0.03 | 0.40 | 0.50 | 0.64 |
| LR | 0.74 | 0.02 | 0.66 | 0.68 | 0.75 |
| LR_SMOTE | 0.72 | 0.04 | 0.69 | 0.70 | 0.73 |

11

| KNN | 0.83 | 0.03 | 1.17 | 0.46 | 0.80 |
| KNN_SMOTE | 0.86 | 0.03 | 2.08 | 0.41 | 0.41 |
| CART | 0.79 | 0.05 | 0.91 | 0.59 | 0.64 |
| CART_SMOTE | 0.84 | 0.03 | 3.95 | 0.55 | 0.55 |
| RF | 0.88 | 0.01 | 0.35 | **0.34** | 0.48 |
| RF_SMOTE | **0.89** | 0.03 | **0.31** | **0.3** | 0.56 |

Table 2. Evaluation of accuracy across various models to determine optimal performance. The mean accuracy results were derived from a k-fold evaluation of the models, focusing on the classification of data categories (i.e., lake, peat, soil) using 10 splits. Log loss was computed for the probability estimates following their calibration using either a sigmoid or isotonic function. For these functions, values approaching 0 signify improved performance.

### 3.3 Downcore analysis

### 3.3.1 Downcore probability estimates and changepoint analysis

**Figure 4:** Downcore probability estimates with changepoint breaks from Random Forests (RF) on the SMOTE dataset with a sigmoid calibration. Results from the Padul (1) and Vanevan (2) records are broken down by lake probabilities (blue curves - a), peat probabilities (green curves -b), and soil probabilities (brown curves - c). Highlighted grey and white boxes indicate changepoint mean breaks identifying phases. Probability estimates from other models can be found in Supplement 1. (Fig. S5 – S8)

300

The Padul record showed mean probabilities for lake, peat, and soil, with lake having the highest probability in 68 out of 93 samples and peat in 25 out of 93 (Fig. 4, column 1). Vanevan's mean probability was .87 for lake, .05 for soil, and .08 for peat, with lake samples having the highest probability in 44 out of 46 samples, peat in 2 out of 46, and no samples having the highest probability in soil (Fig. 4, column 2).
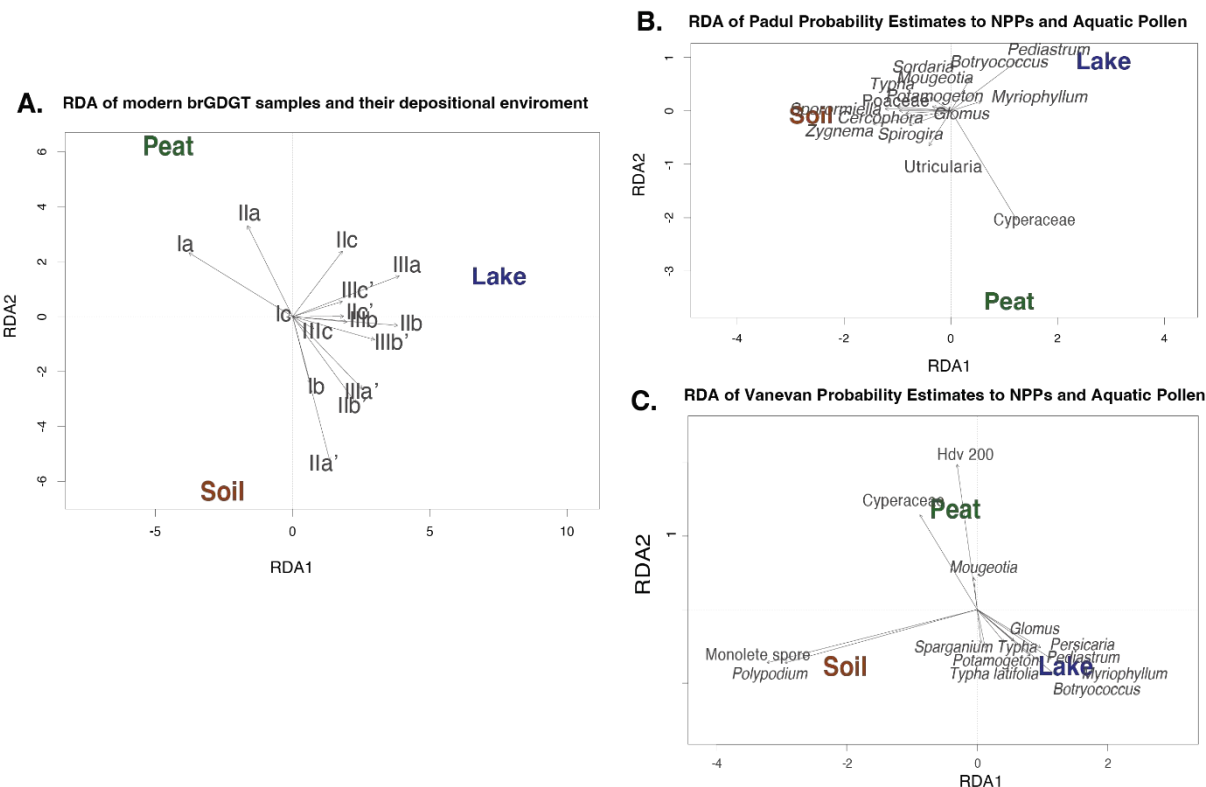
305

### 3.3.2 Probability analysis by changepoint phases

For the Padul record, changepoints were detected at 12837 cal. BP, 20627 cal. BP, and 29617 cal. BP for lake probabilities dividing it into phases 1 - 4 (Fig. 4, column 1) and changepoints in the Vanevan record were identified at 2043 cal. BP, 3577 cal. BP, 4628 cal. BP, 5061 cal. BP, and 8592 cal. BP based lake probabilities (Fig. 4, column 2). Changepoint analysis of

310 the Padul record indicates that brGDGTs based ML lake probabilities peak in phase 3, while in phases 1, 2, and 4, these probabilities vary between soil and peat. BrGDGTs-based ML soil probabilities are the highest in phase 1, while peat

probabilities are consistent, primarily in phases 4 and 2 (Fig. 4, column 2). The probabilities of Vanevan brGDGT-based ML

lake probabilities are elevated across the entire record, peaking during phases 5, 6, and 2, while peat probabilities are

elevated in phases 4 and 1, and soil probabilities exhibit fluctuations with peat and lake, predominantly in phases 4 and 3

315 (Fig. 4, column 2).

## 3.4 RDA analysis on modern and downcore pollen, NPP, and brGDGTs

### 3.4.1 Modern samples



320 **Figure 5:** (A) RDA analysis of the modern fractional abundances of the global brGDGT database with their depositional

environments. (B) Probability estimates results of the depositional environment compared with pollen and NPPs of Padul

record (i.e., Rodrigo-Gámiz et al., 2022; Camuera et al., 2019) and (C) the Vanevan record (i.e., Robles et al., 2022).

The RDA analysis reveals the association of brGDGTs with each depositional unit (i.e., soil, lake, peat) in the global modern

325 database and the downcore probability predictions. Most variance can be explained across RDA-1 (30.3), where peat and soil

sources sit in contrast to lakes in the modern database (Fig. 5a). BrGDGTs Ia and IIa are more clearly associated with peat

and soil depositional environments, while the rest have a stronger association with lake and soil environments (Fig. 5a).

Comparisons between depositional environment probability estimates, aquatic pollen, and NPPs reveal relationships between

14

these variables in the downcore record. Pollen and NPP associations between peat and lake probabilities include Cyperaceae
330 pollen, while algae like *Pediastrum*, *Botryococcus*, and *Myriophyllum* have associations with lake probabilities. For soil,
spores and algae are associated with these depositional environments. Hdv-200 and Cyperaceae have the highest explanatory
power for peat, monolete spores and polypodium for soil, and *Botryococcus*, *Myriophyllum*, and *Pediastrum* for lakes (Fig
5b & 5c).

335 **4 Discussion:**

**4.1 Probability estimates for chosen models and application to downcore records**

**4.1.1 Model accuracy**

The Random Forest model utilizing the SMOTE dataset achieves an F1 score of 89% (Table 2) in classification, which is
lower than the 95% F1 score reported for the BIGMaC model by Martínez-Sosa et al., (2023). The difference in F1 scores is
340 anticipated as a result of differing training datasets and methodologies, including the incorporation of isoGDGTs by
Martínez-Sosa et al., (2023), and their establishment of curated clusters, and the application of classification models. Our
models focus on probability estimate outputs instead of discrete classes, allowing for a nuanced understanding of shifts in
provenance; thus, a lower F1 score is permissible. A score of 89% demonstrated a strong and precise model, despite being
lower than expected. The difference between the new model and the BigMac model is also seen when applied on the
345 downcore records, the BIGMaC model failed to predict soil classification for both cores, whereas our probabilities for soil
were elevated in the Padul record (Supplement 1, Fig. S3 and S4). The inclusion of additional soil samples in our database
enhances soil identification during model training (BIGMaC n= 192, database in this paper n= 750).

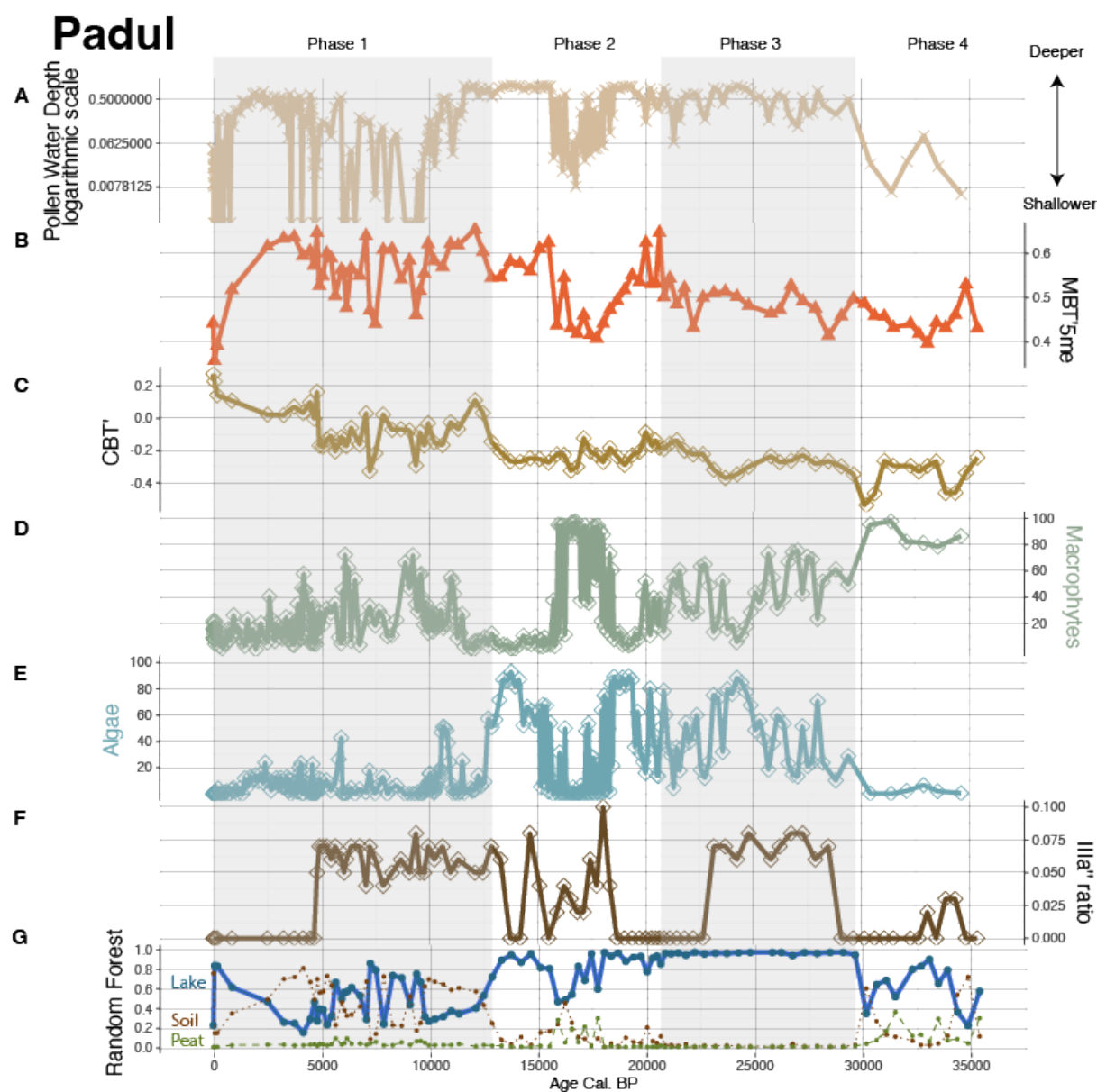**4.1.2 Validating models with pollen, NPPs, and derived water-depth reconstructions.**

350 We evaluated the accuracy of our ML model for detecting provenance change by comparing the probability estimates it
produced with published data on pollen, NPPs, and water depth estimates derived from these proxies (Fig. 5 and 6). The
comparison of GDGT-model variables with depositional environments indicates clear associations. Individual RDA analysis
of both records associates Cyperaceae pollen, prevalent in wetland and peat contexts, with modern brGDGT samples
obtained from peat depositional environments (Fig. 5). Down core records indicate distinct associations between pollen from
355 algae, specifically *Pediastrum* and *Botryococcus*, which are typically associated with open lakes, and the brGDGT based ML
probability estimates associated with lake depositional environment (Fig. 5). Monolete spores in the Vanevan record, along
with *Sordaria* and *Sporormiella* in the Padul record, are related to the brGDGT-trained ML probability estimates of soil
depositional environments. In the Padul record, these spores were likely introduced through human activity (Ramos-Roman
et al., 2018) and may be associated with erosion into the lake. Pollen from semi-aquatic plants, including Cyperaceae and
360 *Typha,* follow similar trends that align with brGDGT-trained ML lake probabilities, as well as increases in brGDGT-trained
ML probabilities for peat and soil across both records (Fig. 6 and 7).

15

The comparison of reconstructed water depth results with the probability estimates from the brGDGT-trained ML model for both cores indicates that the two proxies exhibit similar pattern of increase and decrease. This suggests that our models accurately identify changes in sourcing and hydrology across both records (Fig. 7A). Robles et al., (2022) interpreted the water-depth changes for the Vanevan record based on aquatic pollen and NPPs, identifying a shallow lake from 9700 to 9400 cal. BP, a lake system from 9700 to 5100 cal. BP, a transitional phase from 5100 to 4950 cal. BP, and peatland development from 5100 cal. BP to today. This aligns closely with our changepoint phases, indicating elevated lake probabilities during phases 6 and 5, high peat probabilities during phase 4, and a rise in soil and peat probabilities from our brGDGT-trained ML model over the past 5000 years (Fig. 7E).

**Figure 6:** Comparison of probability estimates from the Padul record with aquatic pollen and NPPs and brGDGT indexes (data from Ramos-Román et al., 2018; Camuera et al., 2019; Rodrigo-Gámiz et al., 2022). (A) Pollen and NPP based water-depth reconstructions (B) MBT'$_{5ME}$ brGDGT index (C) CBT' brGDGT index (D) Selected aquatic plants including Cyperaceae and *Typha*. (E) Selected algae *Pediastrum*, *Botryococcus*, and *Mougeota* as well as aquatic plants Cyperaceae and *Typha*. (F) IIIa'' brGDGT ratio (G) Probability estimates for the lake depositional environments (this study).

17

380 The probability estimates in the Padul record exhibit trends analogous to the Vanevan results, with an alignment with water depth as indicated by pollen and NPPs (Fig. 7B). The estimates derive from the pollen data from Camuera et., (2019), indicating a low water stand in phase 4, a high water stand in phase 3, a fluctuating high to low to high stand in phase 2, and a high, fluctuating to low to high stand in phase 1. The observed trends are reflected in our brGDGT-based ML lake probability estimates. Similar to the Vanevan record, the Padul record predominantly features samples with brGDGT-based

385 ML lake probabilities assigned to lakes. However, there is greater variation among categorical types. This is evident in phases 4 and 3, where peat and soil probabilities are combined with lake probabilities, and in phase 1, where notable fluctuations occur in soil and lake probabilities.

### 4.2 Environmental controls and depositional shifts in downcore brGDGT records

390 **4.2.1 Identifying provenance changes in downcore records (and their impact on the MBT'$_{5ME}$).**
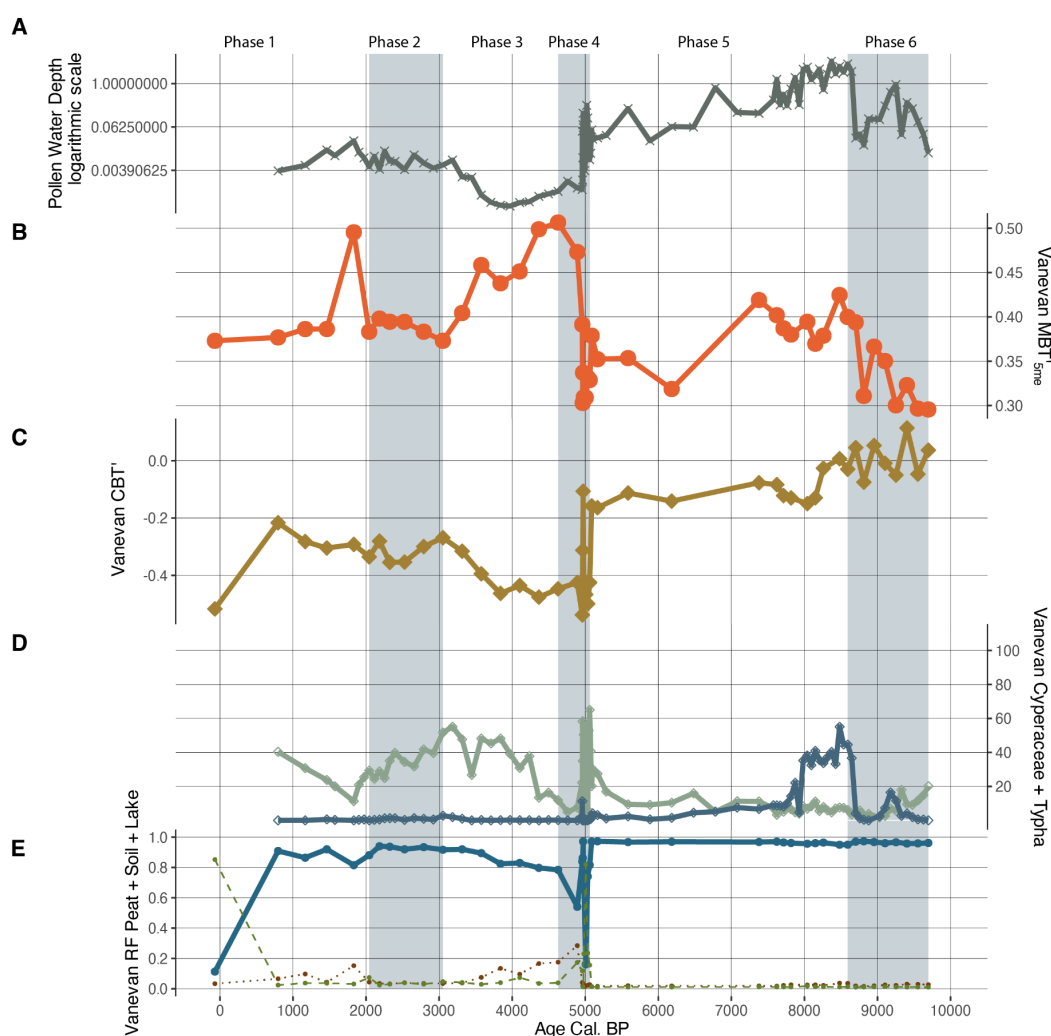
The ML-probability estimates may be interpreted as originating from either a dominant or mixed-sourced sedimentary environment. During periods of high brGDGT-based ML lake probabilities, the findings indicate that brGDGTs are produced in-situ in shallow lakes and wetlands, alongside contributions from other sources. This result is unexpected for the Vanevan record, as the brGDGTs from the last 5000 years exhibit a closer alignment with soil samples when plotted on a

395 ternary diagram (Robles et al., 2022). We compared our results with the classification provided by the BigMAC machine learning model from Martínez-Sosa et al., (2023), which classified most samples as a lake depositional environment, similar to our results, while the samples at 5000 cal. BP for Vanevan were categorized as soil and peat rather than soil and lake. In the Padul record between 35000 – 30000 cal. BP and from 10000 cal. BP – to the present differences between models include samples classified as soil (our model) rather than peat or lake (Supplement 1, Fig. S3). A potential bias in both

400 models may arise from samples in the global database categorized as lakes, but with substantial contributions of brGDGTs from soil or peat depositional environments.

Secondly, our findings indicate that the identification of in-situ produced lacustrine brGDGTs from lake depositional environments using a model provides more nuance than the quantification of the IIIa'' brGDGT isomers. Rodrigo-Gámiz et al., (2022) identified IIIa'' in the Padul record, which is attributed to in-situ brGDGT lake production.

405 Here, the ratio of brGDGTs IIIa'' aligns with the brGDGT-based ML lake probability estimates for this record (Figure 6). The IIIa'' isomer is completely absent from the Vanevan record; however, the brGDGT-based ML lake probability estimates approach 100%. This supports earlier studies indicating that the IIIa'' isomer is not universally found in all lake systems (i.e., Weber et al., 2015; Ding et al., 2018). However, the lack of discussion regarding the absence of IIIa'' isomer in both modern and downcore records is notable and warrants attention in future investigations.

410 The probability and RDA results underscore the need for multivariate methods in the analyses of depositional environments. RDA analysis of the global brGDGT database reveals a distinct separation along RDA-2 between 5- and 6-methyl pentamethylated brGDGTs in various modern depositional environments (Fig. 5a). This is observed between IIa and

IIa', associated with peat and soil depositional environments, respectively (Fig. 6). Martínez-Sosa et al., (2023) identified IIa' as the most significant isomer for provenance classification using their random forest model, a finding that aligns with

415 our models. Furthermore, brGDGT Ia exhibits a stronger association with peat depositional environments compared to other tetramethylated brGDGTs linked to lake environments, highlighting the need to advance beyond ternary diagrams for provenance identification.



**Figure 7:** Comparison between the probability estimates on the Vanevan record and the aquatic pollen and NPPs and

420 brGDGT indexes (data from Robles et al. 2022). (A) Pollen and NPP based water-depth reconstructions (B) MBT'$_{5ME}$ brGDGT index (C) CBT' brGDGT index (D) Select algae and aquatic plants for the Vanevan record algae is *Pediastrum*,

*Botryococcus*, and *Mougeota* and aquatic plants are Cyperaceae and *Typha*. (E) Probability estimates for the lake depositional environments on both records.

425    The results of our models indicate that variations in brGDGT provenance, even within mixed sedimentary environments, significantly influence widely used indexes like the MBT'$_{5ME}$ (Fig. 6 and 7). Pollen and NPPs provide independent confirmation of the impact these changes have on the MBT'$_{5ME}$, particularly when analyzed alongside data from our new brGDGTs global database. In this database, soil (0.56) and peat (0.58) exhibit higher mean MBT'$_{5ME}$ values compared to the lower values observed in lakes (0.39) (Supplement 1, Fig. S2). However, it must be noted that analytical

430    differences between laboratories (i.e., De Jonge et al., 2024) suggests the accuracy of these values may fluctuate, but the overall trends remain. Both the Vanevan and Padul records indicate that MBT'$_{5ME}$ values are elevated during periods characterized by high brGDGT-based ML soil and peat probabilities, while they are reduced during periods of high lake probabilities. The pollen-based water depth reconstructions, serving as an independent proxy, exhibit trends analogous to both the MBT'$_{5ME}$, and the brGDGT-based ML probability estimates. These changes are documented in additional proxies

435    from these records, including XRF and sediment analysis (e.g., Robles et al., 2022; Camerua et al., 2018), highlighting the necessity of identifying the appropriate depositional contexts.

Our findings indicate that even minor changes in provenance can affect the MBT'$_{5ME}$ and CBT' indexes. Where increased brGDGT-based ML soil probabilities occur in the Vanevan and Padul records, they do not reach the threshold indicative of a complete depositional environment shift, instead indicating mixed provenance (Fig. 6 and 7).  The Vanevan

440    record indicates that the large shifts in the MBT'$_{5ME}$ and CBT' occur with increased inputs of soil and peat brGDGTs during phases 3 and 4. In the Padul record, variations in CBT' correlate with increases in brGDGT-based soil ML probabilities, particularly during phase 2.

The co-occurrence of aquatic pollen, NPPs, and MBT'$_{5ME}$ variations indicates that provenance, rather than temperature, drives these changes. Increases in MBT'$_{5ME}$ during phases 3 and 4 of the Vanevan record correspond to shifts in

445    aquatic pollen, indicating a transition from lake to peatland driven by a local catchment fire event (Leroyer et al., 2016; Robles et al., 2022). The observed changes are inconsistent with regional climate reconstructions (i.e., Joannin et al., 2014; Cromartie et al., 2020), confirming that provenance change is the primary driver of this alteration.

### 4.2.2 Environmental drivers of provenance changes

450    The impact of provenance changes on MBT'$_{5ME}$ highlights various factors that can alter the distribution of brGDGTs over time, making it a crucial aspect for environmental reconstructions. The environmental changes that cause a change in GDGT provenance, will also affect the environmental chemistry. While large pH changes have the potential to impact MBT'$_{5ME}$ values in soils, muted pH changes in soils, and the impact on GDGTs produced in lakes, are less well constrained. The introduction of soil brGDGTs into a lake, even in small amounts, can alter the MBT'$_{5ME}$ distribution, and also potentially
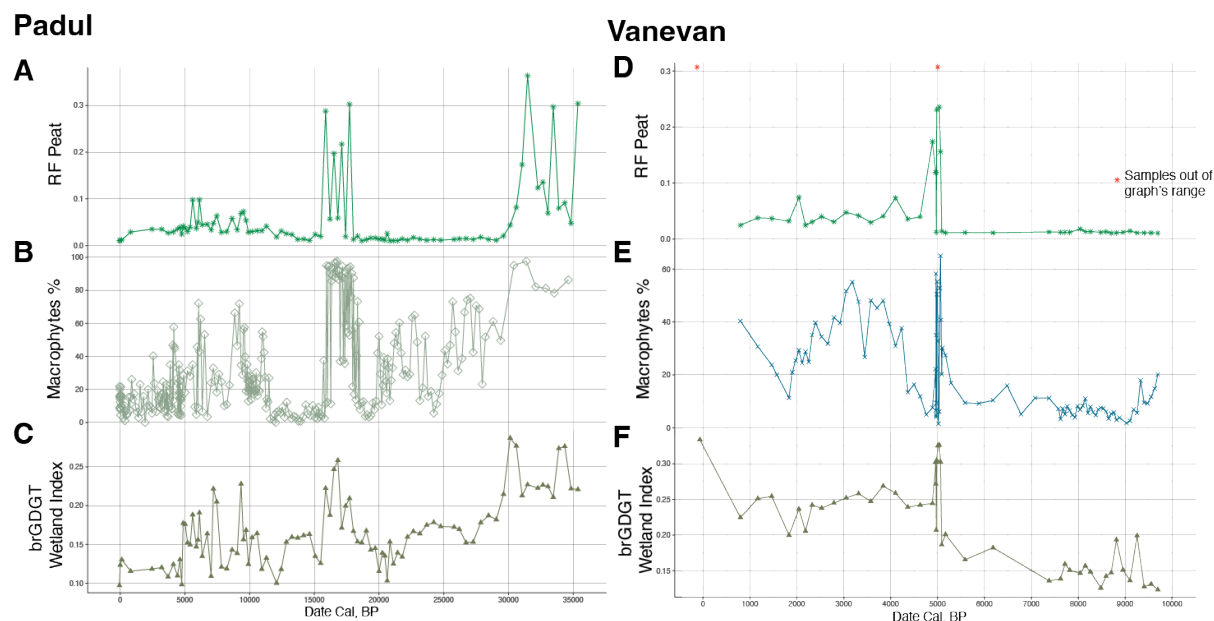
455    introduce pH-related changes.

Rodrigo-Gámiz et al., (2022) identified a relationship in the Padul record between increases in reconstructed pH and MAAT variability within the upper 116 cm, approximately correlating to the last 5000 years. They associated this with a dried ephemeral lake and suggested potential bias in the MBT'$_{5ME}$ reconstruction. In this section of the Padul record, high brGDGT-based ML soil probabilities align with increased CBT' and MBT'$_{5ME}$ values, indicating the contribution of soil-

460 derived brGDGTs, and potentially pH on this variation (Fig. 7). Soil provenance changes appear to exert a greater influence on the MBT'$_{5ME}$ compared to peat and are more prevalent in the record; however, brGDGT-based ML probability estimates for these depositional environments overlap in certain sections of both records.

Our results also highlight the impact of both sudden and gradual depositional changes on the distribution of brGDGTs, driven by hydrological and ecological shifts. Hydrological changes, including variations in water depth in lakes

465 (Stefanescu et al., 2021) and alterations in water-table levels in peat (Ofiti et al., 2024), have been demonstrated to affect brGDGT distribution. The water-depth equations for the Padul and Vanevan records incorporate Cyperaceae pollen at one end of the equation. Cyperaceae is typically associated with the development of wetland ecosystems and the process of lake shallowing. A correlation is observed between MBT'$_{5ME}$ and Cyperaceae for the Padul record (- 0.52, p-value: 0.000003672). There is a correlation between MBT'$_{5ME}$ and the water depth reconstruction for the Vanevan record (0.40, p-value: 0.006).

470 Alterations in hydrology influence shifts in ecological communities, which may or may not be driven by climate. Our results indicate that wetland development, resulting from ecological shifts such as the introduction of aquatic plants and/or lake shallowing can influence the MBT'$_{5ME}$ and the distribution of brGDGTs (Fig. 6 & 7).

### 4.2.4 Use of brGDGT indices to trace environmental changes

475 The shifts in the CBT' which align with the increased soil probabilities in our brGDGT-based ML model suggests that the CBT' can possibly be utilized as a screening tool to detect soil influences. In addition, a new wetland index was developed (Equation 1) to monitor wetland development over time, based on the strong correlation between Cyperaceae pollen and brGDGT [IIa], the association with peat depositional environments in the modern database, and higher median [IIa] values for peat in the global database distribution (Supplement 1 Fig. S1). BrGDGT [Ia] is linked to peatland depositional

480 environments and exhibits a positive correlation with Cyperaceae in the Padul record; however, its inclusion in the equation resulted in a diminished correlation between Cyperaceae and the pollen-based water depth. Consequently, it has been excluded from the equation. In peatlands, incorporating the Ia into the equation may facilitate the monitoring of change.

Equation 1: Wetland Index: ([IIa]) / ([Ia] + [Ib] + [Ic] + [IIa] + [IIa'] +

485 $\quad\quad\quad\quad\quad$ [IIb] + [IIb'] + [IIc] + [IIc'] +

$\quad\quad\quad\quad\quad$ [IIIa] + [IIIb] + [IIIc] +

$\quad\quad\quad\quad\quad$ [IIIa'] + [IIIb'] + [IIIc'])

**Figure 8:** Comparison between the Random Forest peat probabilities (A and D), pollen from Cyperaceae and *Typha* (B and
490   E), and the brGDGT wetland index (Equation 1, C and F) on the Padul and Vanevan records.


The wetland index results align closely with the pollen percentages of Cyperaceae (Padul: 0.52, p-value: 0.000003, Vanevan:
0.70 p-value: 0.0000001) and with *Typha* (Padul: 0.50, p-value: 0.000007, Vanevan 0.65, p-value) for both records (Fig. 8).
In testing the index on the modern sample database, an adequate division between depositional environments is observed,
495   with median values of 0.135, 0.167, and 0.337 for soil, lake, and peat, respectively (Supplement 1, Fig. S2). Furthermore,
while our emphasis was on the frequently employed MBT'$_{5ME}$, the correlation between shifts and individual brGDGTs
underscores the potential impact these modifications may have on any brGDGT index.


### 4.3 Limitations and Future Directions

500   Our findings indicate that multivariate methods, such as machine learning, are needed for analyzing brGDGT distributions.
To effectively utilize these tools, a standardized collection of datasets across research groups is essential, along with an
increased datasets from a variety of environments. A notable limitation of our study was our reliance on the published
sample name. The identification of depositional environments was successful with these names; however, variations within a
depositional environment, for example a shallow versus deep lakes, remain inadequately represented. To effectively utilize
505   these tools, it is essential to collect additional information, including water depth, salinity, pH, redox conditions, and others,
in a standardized manner across research teams.

This study demonstrates the necessity of multi-proxy approaches to comprehend the influence of ecological,
hydrological, and depositional changes on brGDGT-based reconstructions. Numerous studies are presently employing

pollen-based climate reconstructions in conjunction with brGDGT reconstructions (e.g., Watson et al., 2018; Martin et al.,

510 2019; Dugerdil et al., 2021b; Robles et al., 2022:2023; Stefanescu et al., 2021). The findings indicate that aquatic pollen and NPPs provide valuable insights for understanding biases introduced by alterations in depositional environments and provenance.

Many downcore studies are derived from smaller lakes, wetlands, and peatlands (e.g., Martin et al., 2019; Dugerdil et al., 2021b; Robles et al., 2022;2023; Ramos-Román et al., 2022; Acharya et al., 2023; Barhoumi et al., 2024). This study

515 emphasizes the importance of comprehending changes in depositional environments over geological time and advocates for studies in smaller lakes, wetlands, and peatlands. There is a particular necessity for improved classification of wetland environments within the brGDGT literature.


## 5 Conclusion

520 This study demonstrated that multivariate methods enhance the understanding of how provenance changes affect brGDGT distributions and the MBT'$_{5ME}$. A new database of modern samples (n=2301 samples) has been utilized to apply probability estimates from five machine learning algorithms to downcore sediments, facilitating the identification of changes in brGDGT provenance across lake, soil, and peat depositional environments. Utilizing calibrated probability estimates enhances the identification of the provenance of brGDGTs, including those originating from mixed sources.

525 The results indicate that alterations in provenance, depositional environments, and hydrology, particularly the transition from open lakes to wetlands and variations in water depth, can substantially influence the brGDGT signal. The introduction of soil-derived brGDGTs, even in minimal quantities, significantly influences the brGDGT distribution and the MBT'$_{5ME}$. We developed a wetland index, utilizing the fractional abundance of IIa, to analyze shifts between these systems.

This study confirms that independent proxies, including aquatic pollen and non-pollen palynomorphs, can

530 effectively quantify hydrological and ecological changes, thereby influencing the gradual depositional alterations that may affect brGDGT distribution. Our models can accurately and independently identify changes in provenance and are applicable to existing datasets. We suggest that complementary environmental proxies, including fossil pollen, non-pollen palynomorphs, XRF, diatoms, and testate amoebae, among others, are essential for confirming changes in provenance in brGDGT environmental reconstructions.

535

**Competing interests:** The authors declare no competing interests

**Code availability:** Code and data for this project will be publicly available on https://github.com/amycromartie/ProbbrGDGT

**Data availability:** Additionally, database will be uploaded to Pangea

550

**References:**

Acharya, Sudip, Roland Zech, Paul Strobel, Marcel Bliedtner, Maximilian Prochnow, and Cindy De Jonge. 2023.
"Environmental Controls on the Distribution of GDGT Molecules in Lake Höglwörth, Southern Germany." *Organic*
560 *Geochemistry* 186. https://doi.org/10.1016/j.orggeochem.2023.104689.

Baker, Andy, Alison J Blyth, Catherine N Jex, James A Mcdonald, Martijn Woltering, and Stuart J Khan. 2019.
"Glycerol Dialkyl Glycerol Tetraethers (GDGT) Distributions from Soil to Cave: Refining the Speleothem
Paleothermometer." *Organic Geochemistry* 136: 103890.
https://doi.org/https://doi.org/10.1016/j.orggeochem.2019.06.011.

565 Barhoumi, Chéima, Guillemette Ménot, Sébastien Joannin, Adam A Ali, Salomé Ansanay-Alex, Yulia Golubeva,
Dmitry Subetto, Alexander Kryshen, Igor Drobyshev, and Odile Peyron. 2023. "Temperature and Fire Controls on
Vegetation Dynamics in Northern Ural (Russia) Boreal Forests during the Holocene Based on BrGDGT and Pollen
Data." *Quaternary Science Reviews* 305: 108014.

Bell, John F. 1999. "Tree-Based Methods." In *Machine Learning Methods for Ecological Applications*, edited by
570 Alan H Fielding, 89–105. Boston, MA: Springer US. https://doi.org/10.1007/978-1-4615-5289-5_3.

Baxter, A. J., Hopmans, E. C., Russell, J. M., & Sinninghe Damsté, J. S. (2019). Bacterial GMGTs in East African lake
sediments: Their potential as palaeotemperature indicators. *Geochimica et Cosmochimica Acta*, *259*, 155–169.
https://doi.org/https://doi.org/10.1016/j.gca.2019.05.039

575

Buckles, L K, J W H Weijers, X.-M. Tran, S Waldron, and J S Sinninghe Damsté. 2014. "Provenance of Tetraether Membrane Lipids in a Large Temperate Lake (Loch Lomond, UK): Implications for Glycerol Dialkyl Glycerol Tetraether (GDGT)-Based Palaeothermometry." *Biogeosciences* 11 (19): 5539–63. https://doi.org/10.5194/bg-11-5539-2014.

Camuera, Jon, Gonzalo Jiménez-Moreno, María J Ramos-Román, Antonio García-Alix, Jaime L Toney, R Scott Anderson, Francisco Jiménez-Espejo, et al. 2018. "Orbital-Scale Environmental and Climatic Changes Recorded in a New ~200,000-Year-Long Multiproxy Sedimentary Record from Padul, Southern Iberian Peninsula." *Quaternary Science Reviews* 198: 91–114. https://doi.org/https://doi.org/10.1016/j.quascirev.2018.08.014.

Camuera, Jon, Gonzalo Jiménez-Moreno, María J. Ramos-Román, Antonio García-Alix, Jaime L. Toney, R. Scott Anderson, Francisco Jiménez-Espejo, et al. 2019. "Vegetation and Climate Changes during the Last Two Glacial-Interglacial Cycles in the Western Mediterranean: A New Long Pollen Record from Padul (Southern Iberian Peninsula)." *Quaternary Science Reviews* 205. https://doi.org/10.1016/j.quascirev.2018.12.013.

Cao, J., Rao, Z., Shi, F., & Jia, G. (2020). Ice formation on lake surfaces in winter causes warm-season bias of lacustrine brGDGT temperature estimates. Biogeosciences, 17(9), 2521–2536. https://doi.org/10.5194/bg-17-2521-2020

Cearns, M., Hahn, T., Clark, S., & Baune, B. T. (2020). Machine learning probability calibration for high-risk clinical decision-making. In *Australian \& New Zealand Journal of Psychiatry* (Vol. 54, Issue 2, pp. 123–126). SAGE Publications Sage UK: London, England.

Chen, C., Bai, Y., Fang, X., Zhuang, G., Khodzhiev, A., Bai, X., & Murodov, A. (2021). Evaluating the potential of soil bacterial tetraether proxies in westerlies dominating western Pamirs, Tajikistan and implications for paleoenvironmental reconstructions. Chemical Geology, 559, 119908. https://doi.org/https://doi.org/10.1016/j.chemgeo.2020.119908

Cromartie, Amy, Claire Blanchet, Chéïma Barhoumi, Erwan Messager, Odile Peyron, Vincent Ollivier, Pierre Sabatier, et al. 2020. "The Vegetation, Climate, and Fire History of a Mountain Steppe: A Holocene Reconstruction from the South Caucasus, Shenkani, Armenia." *Quaternary Science Reviews* 246: 106485.

Cearns, Micah, Tim Hahn, Scott Clark, and Bernhard T Baune. 2020. "Machine Learning Probability Calibration for High-Risk Clinical Decision-Making." *Australian \& New Zealand Journal of Psychiatry*. SAGE Publications Sage UK: London, England.

d'Oliveira, L, L Dugerdil, G Ménot, A Evin, S D Muller, S Ansanay-Alex, J Azuara, et al. 2023. "Reconstructing 15,000 Years of Southern France Temperatures from Coupled Pollen and Molecular (Branched Glycerol Dialkyl Glycerol Tetraether) Markers (Canroute, Massif Central)." *Climate of the Past* 19 (11): 2127–56. https://doi.org/10.5194/cp-19-2127-2023.

Dankowski, Theresa, and Andreas Ziegler. 2016. "Calibrating Random Forests for Probability Estimation." *Statistics in Medicine* 35 (22): 3949–60.

610         Dang, Xinyue, Weihua Ding, Huan Yang, Richard D. Pancost, B. David A. Naafs, Jiantao Xue, Xiao Lin, Jiayi Lu, and
Shucheng Xie. 2018. "Different Temperature Dependence of the Bacterial BrGDGT Isomers in 35 Chinese Lake
Sediments Compared to That in Soils." *Organic Geochemistry* 119. https://doi.org/10.1016/j.orggeochem.2018.02.008.
Dang, Xinyue, Huan Yang, B David A Naafs, Richard D Pancost, and Shucheng Xie. 2016. "Evidence of Moisture
Control on the Methylation of Branched Glycerol Dialkyl Glycerol Tetraethers in Semi-Arid and Arid
615         Soils."*Geochimica et Cosmochimica Acta* 189: 24–36. https://doi.org/https://doi.org/10.1016/j.gca.2016.06.004.
Davtian, Nina, Edouard Bard, Guillemette Ménot, and Yoann Fagault. 2018. "The Importance of Mass Accuracy in
Selected Ion Monitoring Analysis of Branched and Isoprenoid Tetraethers." *Organic Geochemistry* 118.
https://doi.org/10.1016/j.orggeochem.2018.01.007.
De Jonge, C., E. E. Kuramae, D. Radujković, J. T. Weedon, I. A. Janssens, and F. Peterse. 2021. "The Influence of
620         Soil Chemistry on Branched Tetraether Lipids in Mid- and High Latitude Soils: Implications for BrGDGT- Based
Paleothermometry." *Geochimica et Cosmochimica Acta* 310. https://doi.org/10.1016/j.gca.2021.06.037.
De Jonge, Cindy, Alina Stadnitskaia, Ellen C. Hopmans, Georgy Cherkashov, Andrey Fedotov, and Jaap S. Sinninghe
Damsté. 2014a. "In Situ Produced Branched Glycerol Dialkyl Glycerol Tetraethers in Suspended Particulate Matter
from the Yenisei River, Eastern Siberia." *Geochimica et Cosmochimica Acta* 125.
625         https://doi.org/10.1016/j.gca.2013.10.031.
De Jonge, Cindy, Ellen C Hopmans, Claudia I Zell, Jung-Hyun Kim, Stefan Schouten, and Jaap S Sinninghe Damsté.
2014b. "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., E. E. Kuramae, D. Radujković, J. T. Weedon, I. A. Janssens, and F. Peterse. 2021. "The Influence of Soil
630         Chemistry on Branched Tetraether Lipids in Mid- and High Latitude Soils: Implications for BrGDGT- Based
Paleothermometry." *Geochimica et Cosmochimica Acta* 310. https://doi.org/10.1016/j.gca.2021.06.037.
De Jonge, Cindy, Jingjing Guo, Petter Hällberg, Marco Griepentrog, Hamdi Rifai, Andreas Richter, Edson Ramirez, et
al. 2024. "The Impact of Soil Chemistry, Moisture and Temperature on Branched and Isoprenoid GDGTs in Soils: A
Study Using Six Globally Distributed Elevation Transects." *Organic Geochemistry* 187.
635         https://doi.org/10.1016/j.orggeochem.2023.104706.
De Jonge, Cindy, Francien Peterse, Klaas G J Nierop, Thomas M Blattmann, Marcelo Alexandre, Salome Ansanay-
Alex, Thomas Austin, et al. 2024. "Interlaboratory Comparison of Branched GDGT Temperature and PH Proxies
Using Soils and Lipid Extracts." *Geochemistry, Geophysics, Geosystems* 25 (7): e2024GC011583.
Dearing Crampton-Flood, Emily, Jessica E Tierney, Francien Peterse, Frédérique M S A Kirkels, and Jaap S Sinninghe
640         Damsté. 2020. "BayMBT: A Bayesian Calibration Model for Branched Glycerol Dialkyl Glycerol Tetraethers in Soils
and Peats." *Geochimica et Cosmochimica Acta* 268: 142–59.
Dembla, Gaurav. 2020. "Intuition behind Log-Loss Score." *Towards Data Science*.

Dillon, James T., Sam Lash, Jiaju Zhao, Kevin P. Smith, Peter van Dommelen, Andrew K. Scherer, and Yongsong
Huang. 2018. "Bacterial Tetraether Lipids in Ancient Bones Record Past Climate Conditions at the Time of Disposal."
645     *Journal of Archaeological Science* 96. https://doi.org/10.1016/j.jas.2018.05.009.

Ding, Su, Valérie F Schwab, Nico Ueberschaar, Vanessa-Nina Roth, Markus Lange, Yunping Xu, Gerd Gleixner, and
Georg Pohnert. 2016. "Identification of Novel 7-Methyl and Cyclopentanyl Branched Glycerol Dialkyl Glycerol
Tetraethers in Lake Sediments." *Organic Geochemistry* 102: 52–58.

Ding, S., Xu, Y., Wang, Y., He, Y., Hou, J., Chen, L., & He, J.-S. (2015). Distribution of branched glycerol dialkyl
650     glycerol tetraethers in surface soils of the Qinghai–Tibetan Plateau: implications of brGDGTs-based proxies in cold
and dry regions. *Biogeosciences*, *12*(11), 3141–3151. https://doi.org/10.5194/bg-12-3141-2015

Dillon, James T., Sam Lash, Jiaju Zhao, Kevin P. Smith, Peter van Dommelen, Andrew K. Scherer, and Yongsong
Huang. 2018. "Bacterial Tetraether Lipids in Ancient Bones Record Past Climate Conditions at the Time of Disposal."
*Journal of Archaeological Science* 96. https://doi.org/10.1016/j.jas.2018.05.009.

655     Dugerdil, L, S Joannin, O Peyron, I Jouffroy-Bapicot, B Vannière, B Boldgiv, J Unkelbach, H Behling, and G Ménot.
2021a. "Climate Reconstructions Based on GDGT and Pollen Surface Datasets from Mongolia and Baikal Area:
Calibrations and Applicability to Extremely Cold--Dry Environments over the Late Holocene." *Climate of the Past* 17
(3): 1199–1226. https://doi.org/10.5194/cp-17-1199-2021.

Dugerdil, Lucas, Guillemette Ménot, Odile Peyron, Isabelle Jouffroy-Bapicot, Salomé Ansanay-Alex, Ingrid
660     Antheaume, Hermann Behling, et al. 2021b. "Late Holocene Mongolian Climate and Environment Reconstructions
from BrGDGTs, NPPs and Pollen Transfer Functions for Lake Ayrag: Paleoclimate Implications for Arid Central
Asia." *Quaternary Science Reviews* 273. https://doi.org/10.1016/j.quascirev.2021.107235.

Genuer, Robin, Jean-Michel Poggi, Robin Genuer, and Jean-Michel Poggi. 2020. *Random Forests*. Springer.

Geetha, T.V., Sendhilkumar, S. (2023). Machine Learning: Concepts, Techniques and Applications (1st ed.). Chapman
665     and Hall/CRC. https://doi.org/10.1201/9781003290100

Guo, J., Glendell, M., Meersmans, J., Kirkels, F., Middelburg, J. J., & Peterse, F. (2020). 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

Hastie, Trevor, Robert Tibshirani, Jerome Friedman, Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2009.
670     "Random Forests." *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 587–604.

Halffman, R., Lembrechts, J., Radujković, D., de Gruyter, J., Nijs, I., & de Jonge, C. (2022). Soil chemistry,
temperature and bacterial community composition drive brGDGT distributions along a subarctic elevation gradient.
Organic Geochemistry, 163, 104346. https://doi.org/https://doi.org/10.1016/j.orggeochem.2021.104346

Hilbe, Joseph M. 2016. *Practical Guide to Logistic Regression*. crc Press.

675     Hopmans, Ellen C, Stefan Schouten, and Jaap S Sinninghe Damsté. 2016. "The Effect of Improved Chromatography
on GDGT-Based Palaeoproxies." *Organic Geochemistry* 93: 1–6.

Hunter, J D. 2007. "Matplotlib: A 2D Graphics Environment." *Computing in Science \& Engineering* 9 (3): 90–95. https://doi.org/10.1109/MCSE.2007.55.

Huguet, Carme, Ellen C Hopmans, Wilma Febo-Ayala, David H Thompson, Jaap S Sinninghe Damsté, and Stefan Schouten. 2006. "An Improved Method to Determine the Absolute Abundance of Glycerol Dibiphytanyl Glycerol Tetraether Lipids." *Organic Geochemistry* 37 (9): 1036–41. https://doi.org/https://doi.org/10.1016/j.orggeochem.2006.05.008.

Jaeschke, A., Rethemeyer, J., Lappé, M., Schouten, S., Boeckx, P., & Schefuß, E. (2018). Influence of land use on distribution of soil n-alkane δD and brGDGTs along an altitudinal transect in Ethiopia: Implications for (paleo)environmental studies. Organic Geochemistry, 124, 77–87. https://doi.org/https://doi.org/10.1016/j.orggeochem.2018.06.006

Joannin, Sébastien, Adam A. Ali, Vincent Ollivier, Paul Roiron, Odile Peyron, Samy Chevaux, Samuel Nahapetyan, Petros Tozalakyan, Arkadi Karakhanyan, and Christine Chataigner. 2014. "Vegetation, Fire and Climate History of the Lesser Caucasus: A New Holocene Record from Zarishat Fen (Armenia)." *Journal of Quaternary Science* 29 (17): 70–82. https://doi.org/10.1002/jqs.2679.

Jung, Yoonsuh. 2018. "Multiple Predicting K-Fold Cross-Validation for Model Selection." *Journal of Nonparametric Statistics* 30 (1): 197–215. https://doi.org/10.1080/10485252.2017.1404598.

Kalita, Jugal. 2022. *Machine Learning: Theory and Practice*. Chapman and Hall/CRC.

Kirkels, F., Peterse, F., Ponton, C., Feakins, S. J., & West, A. J. (2016). Soil Organic Carbon Transport in Headwater Tributaries of the Amazon River Traced by Branched GDGTs. AGU Fall Meeting Abstracts, 2016, B13C--0593.

Kou, Q., Zhu, L., Ju, J., Wang, J., Xu, T., Li, C., & Ma, Q. (2022). Influence of salinity on glycerol dialkyl glycerol tetraether-based indicators in Tibetan Plateau lakes: Implications for paleotemperature and paleosalinity reconstructions. *Palaeogeography, Palaeoclimatology, Palaeoecology*, *601*, 111127. https://doi.org/https://doi.org/10.1016/j.palaeo.2022.111127

Kruppa, Jochen, Alexandra Schwarz, Gerhard Arminger, and Andreas Ziegler. 2013. "Consumer Credit Risk: Individual Probability Estimates Using Machine Learning." *Expert Systems with Applications* 40 (13): 5125–31.

Leroyer, Chantal, Sébastien Joannin, David Aoustin, Adam A Ali, Odile Peyron, Vincent Ollivier, Petros Tozalakyan, Arkady Karakhanyan, and Fany Jude. 2016. "Mid Holocene Vegetation Reconstruction from Vanevan Peat (South-Eastern Shore of Lake Sevan, Armenia)." *Quaternary International* 395 (22): 5–18. https://hal-univ-rennes1.archives-ouvertes.fr/hal-01201997/document.

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/https://doi.org/10.1016/j.orggeochem.2016.06.007

Li, Y., Zhao, S., Pei, H., Qian, S., Zang, J., Dang, X., & Yang, H. (2018). Distribution of glycerol dialkyl glycerol tetraethers in surface soils along an altitudinal transect at cold and humid Mountain Changbai: Implications for the

reconstruction of paleoaltimetry and paleoclimate. Science China Earth Sciences, 61(7), 925–939. https://doi.org/10.1007/s11430-017-9168-9

Liang, Jie, Nora Richter, Haichao Xie, Boyang Zhao, Guicai Si, Jian Wang, Juzhi Hou, Gengxin Zhang, and James M. Russell. 2023. "Branched Glycerol Dialkyl Glycerol Tetraether (BrGDGT) Distributions Influenced by Bacterial

715      Community Composition in Various Vegetation Soils on the Tibetan Plateau." *Palaeogeography, Palaeoclimatology, Palaeoecology* 611. https://doi.org/10.1016/j.palaeo.2022.111358.

Liang, Jie, James M. Russell, Haichao Xie, Rachel L. Lupien, Guicai Si, Jian Wang, Juzhi Hou, and Gengxin Zhang. 2019. "Vegetation Effects on Temperature Calibrations of Branched Glycerol Dialkyl Glycerol Tetraether (BrGDGTs) in Soils." *Organic Geochemistry* 127. https://doi.org/10.1016/j.orggeochem.2018.10.010.

720      Loomis, Shannon E., James M. Russell, and Jaap S. Sinninghe Damsté. 2011. "Distributions of Branched GDGTs in Soils and Lake Sediments from Western Uganda: Implications for a Lacustrine Paleothermometer." *Organic Geochemistry* 42 (7). https://doi.org/10.1016/j.orggeochem.2011.06.004.

Loomis, Shannon E, James M Russell, Bethany Ladd, F Alayne Street-Perrott, and Jaap S Sinninghe Damsté. 2012. "Calibration and Application of the Branched GDGT Temperature Proxy on East African Lake Sediments." *Earth and*
725      *Planetary Science Letters* 357–358: 277–88. https://doi.org/https://doi.org/10.1016/j.epsl.2012.09.031.

Loomis, Shannon E, James M Russell, Hilde Eggermont, Dirk Verschuren, and Jaap S Sinninghe Damsté. 2014a. "Effects of Temperature, PH and Nutrient Concentration on Branched GDGT Distributions in East African Lakes: Implications for Paleoenvironmental Reconstruction." *Organic Geochemistry* 66: 25–37. https://doi.org/https://doi.org/10.1016/j.orggeochem.2013.10.012.

730      Loomis, Shannon E, James M Russell, Ana M Heureux, William J D'Andrea, and Jaap S Sinninghe Damsté. 2014b. "Seasonal Variability of Branched Glycerol Dialkyl Glycerol Tetraethers (BrGDGTs) in a Temperate Lake System." *Geochimica et Cosmochimica Acta* 144: 173–87. https://doi.org/https://doi.org/10.1016/j.gca.2014.08.027.

Li, Jingjing, Richard D Pancost, B David A Naafs, Huan Yang, Cheng Zhao, and Shucheng Xie. 2016. "Distribution of Glycerol Dialkyl Glycerol Tetraether (GDGT) Lipids in a Hypersaline Lake System." *Organic Geochemistry* 99:
735      113–24. https://doi.org/https://doi.org/10.1016/j.orggeochem.2016.06.007.

Malley, J. D., J. Kruppa, A. Dasgupta, K. G. Malley, and A. Ziegler. 2012. "Probability Machines: Consistent Probability Estimation Using Nonparametric Learning Machines." *Methods of Information in Medicine* 51 (1). https://doi.org/10.3414/ME00-01-0052.

Manzali, Youness, Mohamed Chahhou, and Mohammed El Mohajir. 2017. "Impure Decision Trees for Auc and Log
740      Loss Optimization." In *2017 International Conference on Wireless Technologies, Embedded and Intelligent Systems (WITS)*, 1–6.

Martin, Céline, Guillemette Ménot, Nicolas Thouveny, Nina Davtian, Valérie Andrieu-Ponel, Maurice Reille, and Edouard Bard. 2019. "Impact of Human Activities and Vegetation Changes on the Tetraether Sources in Lake St Front

(Massif Central, France).” *Organic Geochemistry* 135: 38–52.

745      https://doi.org/https://doi.org/10.1016/j.orggeochem.2019.06.005.

Martin, C., Ménot, G., Thouveny, N., Peyron, O., Andrieu-Ponel, V., Montade, V., Davtian, N., Reille, M., & Bard, E. (2020). Early Holocene Thermal Maximum recorded by branched tetraethers and pollen in Western Europe (Massif Central, France). *Quaternary Science Reviews*, *228*, 106109. https://doi.org/https://doi.org/10.1016/j.quascirev.2019.106109

750      Martínez-Sosa, Pablo, Jessica E. Tierney, Ioana C. Stefanescu, Emily Dearing Crampton-Flood, Bryan N. Shuman, and Cody Routson. 2021. “A Global Bayesian Temperature Calibration for Lacustrine BrGDGTs.” *Geochimica et Cosmochimica Acta* 305. https://doi.org/10.1016/j.gca.2021.04.038.

Martínez-Sosa, Pablo, Jessica E. Tierney, Lina C. Pérez-Angel, Ioana C. Stefanescu, Jingjing Guo, Frédérique Kirkels, Julio Sepúlveda, Francien Peterse, Bryan N. Shuman, and Alberto V. Reyes. 2023. “Development and Application of

755      the Branched and Isoprenoid GDGT Machine Learning Classification Algorithm (BIGMaC) for Paleoenvironmental Reconstruction.” *Paleoceanography and Paleoclimatology* 38 (7). https://doi.org/10.1029/2023PA004611.

Menges, J, C Huguet, J M Alcañiz, S Fietz, D Sachse, and A Rosell-Melé. 2014. “Influence of Water Availability in the Distributions of Branched Glycerol Dialkyl Glycerol Tetraether in Soils of the Iberian Peninsula.” *Biogeosciences* 11 (10): 2571–81. https://doi.org/10.5194/bg-11-2571-2014.

760      Murphy, Kevin P. 2012. *Machine Learning: A Probabilistic Perspective*. MIT press.

Naafs, B D A, A V Gallego-Sala, G N Inglis, and R D Pancost. 2017a. “Refining the Global Branched Glycerol Dialkyl Glycerol Tetraether (BrGDGT) Soil Temperature Calibration.” *Organic Geochemistry* 106: 48–56.

Naafs, B. D.A., G. N. Inglis, Y. Zheng, M. J. Amesbury, H. Biester, R. Bindler, J. Blewett, et al. 2017b. “Introducing Global Peat-Specific Temperature and PH Calibrations Based on BrGDGT Bacterial Lipids.” *Geochimica et*

765      *Cosmochimica Acta* 208. https://doi.org/10.1016/j.gca.2017.01.038.

Niculescu-Mizil, Alexandru, and Rich Caruana. 2005. “Predicting Good Probabilities with Supervised Learning.” In *Proceedings of the 22nd International Conference on Machine Learning*, 625–32.

Ning, D., Zhang, E., Shulmeister, J., Chang, J., Sun, W., & Ni, Z. (2019). Holocene mean annual air temperature (MAAT) reconstruction based on branched glycerol dialkyl glycerol tetraethers from Lake Ximenglongtan,

770      southwestern China. *Organic Geochemistry*, *133*, 65–76. https://doi.org/https://doi.org/10.1016/j.orggeochem.2019.05.003

Ofiti, Nicholas O.E., Arnaud Huguet, Paul J. Hanson, and Guido L.B. Wiesenberg. 2024. “Peatland Warming Influences the Abundance and Distribution of Branched Tetraether Lipids: Implications for Temperature Reconstruction.” *Science of the Total Environment* 924. https://doi.org/10.1016/j.scitotenv.2024.171666.

775      Oksanen, Jari, F Guillaume Blanchet, Michael Friendly, Roeland Kindt, Pierre Legendre, Dan McGlinn, Peter R Minchin, et al. 2019. “Vegan: Community Ecology Package.” https://cran.r-project.org/package=vegan.

Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, et al. 2011. "Scikit-Learn: Machine Learning in Python." *Journal of Machine Learning Research* 12 (Oct): 2825–30.

780 Pérez-Angel, L. C., Sepúlveda, J., Molnar, P., Montes, C., Rajagopalan, B., Snell, K., Gonzalez-Arango, C., & Dildar, N. (2020). 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.

Peterse, Francien, Jaap van der Meer, Stefan Schouten, Johan W H Weijers, Noah Fierer, Robert B Jackson, Jung-Hyun Kim, and Jaap S Sinninghe Damsté. 2012. "Revised Calibration of the MBT--CBT Paleotemperature Proxy

785 Based on Branched Tetraether Membrane Lipids in Surface Soils." *Geochimica et Cosmochimica Acta* 96: 215–29.

Raberg, Jonathan H., Edgart Flores, Sarah E. Crump, Greg de Wet, Nadia Dildar, Gifford H. Miller, Áslaug Geirsdóttir, and Julio Sepúlveda. 2022. "Intact Polar BrGDGTs in Arctic Lake Catchments: Implications for Lipid Sources and Paleoclimate Applications." *Journal of Geophysical Research: Biogeosciences* 127 (10). https://doi.org/10.1029/2022JG006969.

790 Qian, S., Yang, H., Dong, C., Wang, Y., Wu, J., Pei, H., Dang, X., Lu, J., Zhao, S., & Xie, S. (2019). Rapid response of fossil tetraether lipids in lake sediments to seasonal environmental variables in a shallow lake in central China : Implications for the use of tetraether-based proxies. *Organic Geochemistry*, *128*, 108–121. https://doi.org/https://doi.org/10.1016/j.orggeochem.2018.12.007

Raberg, J. H., Harning, D. J., Crump, S. E., de Wet, G., Blumm, A., Kopf, S., Geirsdóttir, Á., Miller, G. H., &

795 Sepúlveda, J. (2021). Revised fractional abundances and warm-season temperatures substantially improve brGDGT calibrations in lake sediments. *Biogeosciences*, *18*(12). https://doi.org/10.5194/bg-18-3579-2021

Raberg, Jonathan H., Gifford H. Miller, Áslaug Geirsdóttir, and Julio Sepúlveda. 2022. "Near-Universal Trends in BrGDGT Lipid Distributions in Nature." *Science Advances* 8 (20). https://doi.org/10.1126/sciadv.abm7625.

Ramos-Román, María J., Cindy De Jonge, Eniko Magyari, Daniel Veres, Liisa Ilvonen, Anne Lise Develle, and Heikki

800 Seppä. 2022. "Lipid Biomarker (BrGDGT)- and Pollen-Based Reconstruction of Temperature Change during the Middle to Late Holocene Transition in the Carpathians." *Global and Planetary Change* 215. https://doi.org/10.1016/j.gloplacha.2022.103859.

Ramos-Román, María J., Gonzalo Jiménez-Moreno, Jon Camuera, Antonio García-Alix, R. Scott Anderson, Francisco J. Jiménez-Espejo, and José S. Carrión. 2018. "Holocene Climate Aridification Trend and Human Impact Interrupted

805 by Millennial- and Centennial-Scale Climate Fluctuations from a New Sedimentary Record from Padul (Sierra Nevada, Southern Iberian Peninsula)." *Climate of the Past* 14 (1). https://doi.org/10.5194/cp-14-117-2018.

Robles, Mary, Odile Peyron, Elisabetta Brugiapaglia, Guillemette Ménot, Lucas Dugerdil, Vincent Ollivier, Salomé Ansanay-Alex, et al. 2022. "Impact of Climate Changes on Vegetation and Human Societies during the Holocene in the South Caucasus (Vanevan, Armenia): A Multiproxy Approach Including Pollen, NPPs and BrGDGTs."

810 *Quaternary Science Reviews* 277 (February). https://doi.org/10.1016/j.quascirev.2021.107297.

Robles, Mary, Odile Peyron, Guillemette Ménot, Elisabetta Brugiapaglia, Sabine Wulf, Oona Appelt, Marion Blache, et al. 2023. "Climate Changes during the Late Glacial in Southern Europe: New Insights Based on Pollen and BrGDGTs of Lake Matese in Italy." *Climate of the Past* 19 (2): 493–515.

Rodrigo-Gámiz, Marta, Antonio García-Alix, Gonzalo Jiménez-Moreno, María J. Ramos-Román, Jon Camuera, Jaime
L. Toney, Dirk Sachse, R. Scott Anderson, and Jaap S. Sinninghe Damsté. 2022. "Paleoclimate Reconstruction of the Last 36 Kyr Based on Branched Glycerol Dialkyl Glycerol Tetraethers in the Padul Palaeolake Record (Sierra Nevada, Southern Iberian Peninsula)." *Quaternary Science Reviews* 281. https://doi.org/10.1016/j.quascirev.2022.107434.

Rao, Zhiguo, Haichun Guo, Shikai Wei, Jiantao Cao, and Guodong Jia. 2022. "Influence of Water Conditions on Peat BrGDGTs: A Modern Investigation and Its Paleoclimatic Implications." *Chemical Geology* 606.
https://doi.org/10.1016/j.chemgeo.2022.120993.

Simpson, Gavin L. 2007. "Analogue Methods in Palaeoecology: Using the Analogue Package." *Journal of Statistical Software* 22: 1–29.

Russell, James M, Ellen C Hopmans, Shannon E Loomis, Jie Liang, and Jaap S Sinninghe Damsté. 2018.
"Distributions of 5- and 6-Methyl Branched Glycerol Dialkyl Glycerol Tetraethers (BrGDGTs) in East African Lake
Sediment: Effects of Temperature, PH, and New Lacustrine Paleotemperature Calibrations." *Organic Geochemistry* 117: 56–69. https://doi.org/https://doi.org/10.1016/j.orggeochem.2017.12.003.

Siriseriwan, W. 2019. "A Collection of Oversampling Techniques for Class Imbalance Problem Based on SMOTE."

Stefanescu, Ioana C., Bryan N. Shuman, and Jessica E. Tierney. 2021. "Temperature and Water Depth Effects on BrGDGT Distributions in Sub-Alpine Lakes of Mid-Latitude North America." *Organic Geochemistry* 152.
https://doi.org/10.1016/j.orggeochem.2020.104174.

Team, R Core. 2020. "R Core Team R: A Language and Environment for Statistical Computing." *Foundation for Statistical Computing*.

Tierney, Jessica E., and James M. Russell. 2009. "Distributions of Branched GDGTs in a Tropical Lake System: Implications for Lacustrine Application of the MBT/CBT Paleoproxy." *Organic Geochemistry* 40 (9).
https://doi.org/10.1016/j.orggeochem.2009.04.014.

Tierney, J E, J M Russell, H Eggermont, E C Hopmans, D Verschuren, and J S Sinninghe Damsté. 2010.
"Environmental Controls on Branched Tetraether Lipid Distributions in Tropical East African Lake Sediments."
*Geochimica et Cosmochimica Acta* 74 (17): 4902–18. https://doi.org/https://doi.org/10.1016/j.gca.2010.06.002.

Véquaud, Pierre, Alexandre Thibault, Sylvie Derenne, Christelle Anquetil, Sylvie Collin, Sergio Contreras, Andrew T
Nottingham, Pierre Sabatier, Josef P Werne, and Arnaud Huguet. 2022. "FROG: A Global Machine-Learning Temperature Calibration for Branched GDGTs in Soils and Peats." *Geochimica et Cosmochimica Acta* 318: 468–94.
https://doi.org/https://doi.org/10.1016/j.gca.2021.12.007.

Véquaud, P., Derenne, S., Thibault, A., Anquetil, C., Bonanomi, G., Collin, S., Contreras, S., Nottingham, A. T., Sabatier, P., Salinas, N., Scott, W. P., Werne, J. P., & Huguet, A. (2021a.). Development of global temperature and pH

845     calibrations based on bacterial 3-hydroxy fatty acids in soils. *Biogeosciences*, *18*(12), 3937–3959.

https://doi.org/10.5194/bg-18-3937-2021

Véquaud, P., Derenne, S., Anquetil, C., Collin, S., Poulenard, J., Sabatier, P., & Huguet, A. (2021b). Influence of

environmental parameters on the distribution of bacterial lipids in soils from the French Alps: Implications for paleo-

reconstructions. *Organic Geochemistry*, *153*, 104194.

850     https://doi.org/https://doi.org/10.1016/j.orggeochem.2021.104194

Wang, H., Liu, W., & Lu, H. (2016). Appraisal of branched glycerol dialkyl glycerol tetraether-based indices for North

China. *Organic Geochemistry*, *98*, 118–130. https://doi.org/https://doi.org/10.1016/j.orggeochem.2016.05.013

Wang, H., Liu, W., & Lu, H. (2016). Appraisal of branched glycerol dialkyl glycerol tetraether-based indices for North

China. *Organic Geochemistry*, *98*, 118–130. https://doi.org/https://doi.org/10.1016/j.orggeochem.2016.05.013

855     Wang, M., Zong, Y., Zheng, Z., Man, M., Hu, J., & Tian, L. (2018). Utility of brGDGTs as temperature and

precipitation proxies in subtropical China. *Scientific Reports*, *8*(1), 194. https://doi.org/10.1038/s41598-017-17964-0

Wang, H., An, Z., Lu, H., Zhao, Z., & Liu, W. (2020a). Calibrating bacterial tetraether distributions towards in situ soil

temperature and application to a loess-paleosol sequence. *Quaternary Science Reviews*, *231*, 106172.

https://doi.org/https://doi.org/10.1016/j.quascirev.2020.106172

860     Wang, M., Yang, H., Zheng, Z., & Tian, L. (2020b). Altitudinal climatic index changes in subtropical China indicated

from branched glycerol dialkyl glycerol tetraethers proxies. *Chemical Geology*, *541*, 119579.

https://doi.org/https://doi.org/10.1016/j.chemgeo.2020.119579

Wang, Huanye, Weiguo Liu, Yuxin He, Aifeng Zhou, Hui Zhao, Hu Liu, Yunning Cao, et al. 2021. "Salinity-

Controlled Isomerization of Lacustrine BrGDGTs Impacts the Associated MBT5ME' Terrestrial Temperature Index."

865     *Geochimica et Cosmochimica Acta* 305. https://doi.org/10.1016/j.gca.2021.05.004.

Watson, Benjamin I., John W. Williams, James M. Russell, Stephen T. Jackson, Linda Shane, and Thomas V. Lowell.

2018. "Temperature Variations in the Southern Great Lakes during the Last Deglaciation: Comparison between Pollen

and GDGT Proxies." *Quaternary Science Reviews* 182. https://doi.org/10.1016/j.quascirev.2017.12.011.

Warden, L, J.-H. Kim, C Zell, G.-J. Vis, H de Stigter, J Bonnin, and J S Sinninghe Damsté. 2016. "Examining the

870     Provenance of Branched GDGTs in the Tagus River Drainage Basin and Its Outflow into the Atlantic Ocean over the

Holocene to Determine Their Usefulness for Paleoclimate Applications." *Biogeosciences* 13 (20): 5719–38.

https://doi.org/10.5194/bg-13-5719-2016.

Weber, Yuki, Jaap S.Sinninghe Damsté, Jakob Zopfi, Cindy De Jonge, Adrian Gilli, Carsten J. Schubert, Fabio Lepori,

Moritz F. Lehmann, and Helge Niemann. 2018. "Redox-Dependent Niche Differentiation Provides Evidence for

875     Multiple Bacterial Sources of Glycerol Tetraether Lipids in Lakes." *Proceedings of the National Academy of Sciences

of the United States of America* 115 (43). https://doi.org/10.1073/pnas.1805186115.

Weber, Yuki, Cindy De Jonge, W. Irene C. Rijpstra, Ellen C. Hopmans, Alina Stadnitskaia, Carsten J. Schubert,

Moritz F. Lehmann, Jaap S. Sinninghe Damsté, and Helge Niemann. 2015. "Identification and Carbon Isotope

Composition of a Novel Branched GDGT Isomer in Lake Sediments: Evidence for Lacustrine Branched GDGT

880    Production." *Geochimica et Cosmochimica Acta* 154. https://doi.org/10.1016/j.gca.2015.01.032.

Weijers, Johan W H, Stefan Schouten, Ellen C Hopmans, Jan A J Geenevasen, Olivier R P David, Joanna M Coleman, Rich D Pancost, and Jaap S Sinninghe Damsté. 2006. "Membrane Lipids of Mesophilic Anaerobic Bacteria Thriving in Peats Have Typical Archaeal Traits." *Environmental Microbiology* 8 (4): 648–57.

Weijers, Johan W H, Stefan Schouten, Jurgen C van den Donker, Ellen C Hopmans, and Jaap S Sinninghe Damsté.

885    2007. "Environmental Controls on Bacterial Tetraether Membrane Lipid Distribution in Soils." *Geochimica et Cosmochimica Acta* 71: 703–13.

Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. https://ggplot2.tidyverse.org.

Wu, D., Cao, J., Jia, G., Guo, H., Shi, F., Zhang, X., & Rao, Z. (2020). Peat brGDGTs-based Holocene temperature

890    history of the Altai Mountains in arid Central Asia. *Palaeogeography, Palaeoclimatology, Palaeoecology*, *538*, 109464. https://doi.org/https://doi.org/10.1016/j.palaeo.2019.109464

Wu, Jie, Huan Yang, Richard D. Pancost, B. David A. Naafs, Shi Qian, Xinyue Dang, Huiling Sun, et al. 2021. "Variations in Dissolved O2 in a Chinese Lake Drive Changes in Microbial Communities and Impact Sedimentary GDGT Distributions." *Chemical Geology* 579. https://doi.org/10.1016/j.chemgeo.2021.120348.

895    Xiao, W., Xu, Y., Ding, S., Wang, Y., Zhang, X., Yang, H., Wang, G., & Hou, J. (2015). Global calibration of a novel, branched GDGT-based soil pH proxy. *Organic Geochemistry*, *89–90*, 56–60. https://doi.org/https://doi.org/10.1016/j.orggeochem.2015.10.005

Yang, H., Lü, X., Ding, W., Lei, Y., Dang, X., & Xie, S. (2015). The 6-methyl branched tetraethers significantly affect the performance of the methylation index (MBT′) in soils from an altitudinal transect at Mount Shennongjia. *Organic*

900    *Geochemistry*, *82*, 42–53. https://doi.org/https://doi.org/10.1016/j.orggeochem.2015.02.003

Yao, Y., Zhao, J., Vachula, R. S., Werne, J. P., Wu, J., Song, X., & Huang, Y. (2020). Correlation between the ratio of 5-methyl hexamethylated to pentamethylated branched GDGTs (HP5) and water depth reflects redox variations in stratified lakes. *Organic Geochemistry*, *147*, 104076. https://doi.org/https://doi.org/10.1016/j.orggeochem.2020.104076

905    Zink, Klaus-G., Marcus J Vandergoes, Kai Mangelsdorf, Ann C Dieffenbacher-Krall, and Lorenz Schwark. 2010. "Application of Bacterial Glycerol Dialkyl Glycerol Tetraethers (GDGTs) to Develop Modern and Past Temperature Estimates from New Zealand Lakes." *Organic Geochemistry* 41 (9): 1060–66. https://doi.org/https://doi.org/10.1016/j.orggeochem.2010.03.004.