Climate
of the Past
Discussions

Open Access

EGU

1 **OPTiMAL: A new machine learning approach for GDGT-based palaeothermometry**

2

3 Yvette L. Eley,[1,] William Thomson[2], Sarah E. Greene[1], Ilya Mandel[3,4], Kirsty Edgar[1], James A. Bendle[1],

4 and Tom Dunkley Jones[1]

5

6 **Affiliations:**

7 [1]School of Geography, Earth and Environmental Sciences, University of Birmingham, Edgbaston, B15
8 2TT, UK

9 [2]School of Mathematics, University of Birmingham, Edgbaston, B15 2TT, UK

10 [3]Institute of Gravitational Wave Astronomy, School of Physics and Astronomy, University of
11 Birmingham, Edgbaston, B15 2TT, UK

12 [4]Monash Centre for Astrophysics, School of Physics and Astronomy, Monash University, Clayton,
13 Victoria 3800, Australia

14

15

16 **Abstract**

17

18 In the modern oceans, the relative abundances of Glycerol dialkyl glycerol tetraether (GDGTs) compounds

19 produced by marine archaeal communities show a significant dependence on the local sea surface

20 temperature at the site of formation. When preserved in ancient marine sediments, the measured abundances

21 of these fossil lipid biomarkers thus have the potential to provide a geological record of long-term variability

22 in planetary surface temperatures. Several empirical calibrations have been made between observed GDGT

23 relative abundances in late Holocene core top sediments and modern upper ocean temperatures. These

24 calibrations form the basis of the widely used $TEX_{86}$ palaeothermometer. There are, however, two

25 outstanding problems with this approach, first the appropriate assignment of uncertainty to estimates of

26 ancient sea surface temperatures based on the relationship of the ancient GDGT assemblage to the modern

27 calibration data set; and second, the problem of making temperature estimates beyond the range of the

28 modern empirical calibrations (>30 ºC). Here we apply modern machine-learning tools, including Gaussian

29 Process Emulators and forward modelling, to develop a new mathematical approach we call OPTiMAL

30 (**O**ptimised **P**alaeothermometry from **T**etraethers v**i**a **MA**chine **L**earning) to improve temperature

31 estimation and the representation of uncertainty based on the relationship between ancient GDGT

32 assemblage data and the structure of the modern calibration data set. We reduce the root mean square

33 uncertainty on temperature predictions (validated using the modern data set) from ~± 6 ºC using $TEX_{86}$

34 based estimators to ± 3.6 ºC using Gaussian Process estimators for temperatures below 30 ºC. We also

35 provide a new but simple quantitative measure of the distance between an ancient GDGT assemblage and

Climate
of the Past
Discussions

36   the nearest neighbour within the modern calibration dataset, as a test for significant non-analogue

37   behaviour. Finally, we advocate against the use of temperature estimates beyond the range of the modern

38   empirical calibration dataset, given the absence – to date - of a robust predictive biological model or

39   extensive and reproducible mesocosm experimental data in this elevated temperature range.

40

41   **1. Introduction**

42

43   Glycerol dialkyl glycerol tetraethers (GDGTs) are membrane lipids consisting of isoprenoid carbon

44   skeletons ether-bound to glycerol (Schouten et al., 2013). In marine systems they are primarily produced

45   by ammonia oxidising marine Thamarchaeota (Schouten et al., 2013), although some bacterial production

46   may also be important, especially in sub-freezing ecosystems (Siliakus et al., 2017). In modern marine core

47   top sediments, the relative abundance of GDGT compounds with more ring structures increases with the

48   mean annual sea surface temperature (SST) of the overlying waters (Schouten et al., 2002). This trend is

49   most likely driven by the need for increased cell membrane stability and rigidity at higher temperatures

50   (Sinninghe Damsté et al., 2002). On this basis, the $TEX_{86}$ (tetraether index of tetraethers containing 86

51   carbon atoms) ratio was derived to provide an index to represent the extent of cyclisation (Eq. 1; where

52   GDGT-x represents the fractional abundance of GDGT-x by LCMS peak area, and cren' is the peak area

53   of the regioisomer of crenarchaeol) (Schouten et al., 2002) and was shown to be positively correlated with

54   mean annual SSTs:

55

56   $TEX_{86} = (GDGT\text{-}2 + GDGT\text{-}3 + cren')/ (GDGT\text{-}1 + GDGT\text{-}2 + GDGT\text{-}3 + cren')$  (Eq. 1)

57

58   Early applications of $TEX_{86}$ to reconstruct ancient SSTs were promising, especially in providing

59   temperature estimates in environments where standard carbonate-based proxies are hampered by poor

60   preservation ( Schouten et al., 2003; Herfort et al., 2006; Schouten et al., 2007; Huguet et al., 2006; Sluijs

61   et al., 2006; Brinkhuis et al., 2006; Pearson et al., 2007; Slujis et al., 2009). The $TEX_{86}$ approach also

62   extended beyond the range of the widely used alkenone-based $U^{k'}_{37}$ thermometer, in both temperature space,

63   where $U^{k'}_{37}$ saturates at ~28ºC (Brassell, 2014; Zhaung et al., 2017), and back into the early Cenozoic (Bijl

64   et al., 2009; Hollis et al., 2009; Bijl et al., 2013; Inglis et al., 2015) and Mesozoic (Schouten et al., 2002;

65   Jenkyns et al., 2012; O'Brien et al., 2017) where haptophyte-derived alkenones are typically absent from

66   marine sediments (Brassell, 2014). Initially, $TEX_{86}$ was converted to SSTs using an equation derived by

67   Schouten et al. (2002) (Eq. 2):

68

69   $SST \ ºC = 66.7 * TEX_{86} - 18.7$ (Eq. 2)

Climate
of the Past
Discussions

70

71  However as the number and range of applications of $TEX_{86}$ palaeothermometry grew, concerns arose about

72  proxy behaviour at both the high (Liu et al., 2009) and low (Kim et al., 2008) temperature ends of the

73  modern calibration. In response to these criticisms, a new expanded modern core top dataset (Kim et al.,

74  2010) was used to generate two new forms of the GDGT proxy – $TEX_{86}^L$ (Eq. 3), an exponential function

75  that does not include the crenarchaeol regio-isomer and was recommended for use across the entire

76  temperature range of the new core top data (-3 to 30 ºC, particularly when SSTs are lower than 15 ºC), and

77  $TEX_{86}^H$ (Eq. 4), also exponential, and recommended for use when SSTs exceeded 15 ºC (Kim et al. 2010).

78  $TEX_{86}^L$ also excludes GDGT abundance data from the high-temperature regimes of the Red Sea, with the

79  rationale that high-salinity conditions are responsible for somewhat anomalous GDGT compositions in this

80  region (Kim et al. 2010).

81

82  $$TEX_{86}^L = log\left(\frac{[GDGT2]}{[GDGT1]+[GDGT2]+[GDGT3]}\right) \qquad \text{Eq. 3}$$

83

84

85  $$TEX_{86}^H = log\left(\frac{[GDGT2]+[GDGT3]+[Cren']}{[GDGT1]+[GDGT2]+[GDGT3]+[Cren']}\right) \text{Eq. 4}$$

86

87  Despite the recommendations of Kim et al. (2010), both $TEX_{86}^H$ and $TEX_{86}^L$ were widely used and tested

88  across a range of temperatures and palaeoenvironments, including comparisons against other

89  palaeotemperature proxy systems (Hertzberg et al., 2016; Zhang et al., 2014; Seki et al., 2014; Douglas et

90  al., 2014; Linnert et al., 2014; Tyler; Hollis et al. 2012; Dunkley Jones et al. 2013; Lunt 2012). The rationale

91  was that both $TEX_{86}^L$ and $TEX_{86}^H$ were calibrated across a full temperature range, with the exception of the

92  inclusion or exclusion of Red Sea core-top data. The difference in model fit between the two proxy

93  formulations to the calibration dataset was also minor (Kim et al. 2010). In certain environments, however,

94  $TEX_{86}^L$ was subject to significant variability in derived temperatures that were not apparent in $TEX_{86}^H$

95  (Taylor et al. 2013). This was mostly due to changing GDGT2 to GDGT3 ratios, which strongly influence

96  $TEX_{86}^L$, and may be related to the GDGT productivity environment and deep-water lipid production, (Taylor

97  et al., 2013). As a result, $TEX_{86}^L$ is no longer regarded as an appropriate tool for palaeotemperature

98  reconstructions, except in limited Polar conditions (Kim et al., 2010; Tierney, 2012).

99

100  Three fundamental issues have always troubled the $TEX_{86}$ proxy. The first is a concern about undetected

101  non-analogue palaeo-GDGT assemblages, for which the modern calibration data set is inadequate to

102  provide a robust temperature estimation. Although various screening protocols, with independent indices

103    and thresholds, have been proposed to test for an excessive influence of terrestrial lipids (BIT index;

104    Hopmans et al. 2004), within sediment methanogenesis (Methane Index, 'MI'; Zhang et al., 2011) and non-

105    thermal effects such as nutrient levels and archaeal community structure to impact the weighted average of

106    cyclopentane moieties (Ring Index, 'RI;' Zhang et al., 2016), these do not provide a fundamental measure

107    of the proximity between GDGT abundance distributions in the modern, and ancient GDGT abundance

108    distributions recorded in sediment samples. The fundamental question remains – are measured ancient

109    assemblages of GDGT compounds anything like the modern assemblages, from which palaeotemperatures

110    are being estimated? Understanding this question is not helped by the use of indices – $TEX_{86}$ itself, or BIT

111    and MI – that collapse the dimensionality of GDGT abundance relationships onto a single axis of variation.

112

113    Second, from the earliest applications of the $TEX_{86}$ proxy to deep-time warm climate states (Schouten et

114    al., 2003) it was recognized that reconstructed temperatures beyond the range of the modern calibration

115    (>30 ºC), were highly sensitive to model choice within the modern calibration range. Thus, Schouten et al.

116    (2003) restricted their calibration data for deep-time temperature estimates to core-top data in the modern

117    with mean annual SSTs over 20 ºC. However, this problem of model choice, and its impact on temperature

118    estimation beyond the modern calibration range, persists (Hollis et al. in review), with current arguments

119    focused on whether there is an exponential (e.g. Cramwinckel et al., 2018) or linear (Tierney & Tingley,

120    2015) dependency of $TEX_{86}$ on SSTs, and the effect of these models on temperature estimates over 30 ºC.

121

122    Culture and mesocosm studies are sometimes cited in support of extrapolations beyond the modern

123    calibration range when reconstructing ancient SSTs (Kim et al., 2010, Hollis et al., in review). However,

124    close examination of these culture studies reveals a significant variation in the patterns of archaeal GDGT

125    production in response to increasing growth temperature (e.g., Elling et al., 2015). At present, there are a

126    limited number of pure Thaumarchaeotal strains that can be cultured in the laboratory (Qin et al., 2014,

127    2015). Of the existing studies on these cultures, several focus on non-thermal environmental or

128    physiological variables, such as oxygen availability (Qin et al., 2015), that may also influence $TEX_{86}$. Early

129    mesocosm studies indicated that a $TEX_{86}$  to temperature relationship was maintained up to ~35-40 ºC, but

130    with differing linear slopes from modern core top calibrations (Schouten et al., 2007; Wuchter et al., 2004).

131    The production of the crenarchaeol regional isomer (cren') in these elevated temperature mesocosms is

132    typically lower than that found in sediments (Pearson & Ingalls, 2013; Schouten et al., 2007), indicating

133    that the GDGT production environment in the mesocosms was in some fundamental way not analogous to

134    the time-averaged GDGT assemblages recovered from core-top samples. A more recent study of three pure

135    archaeal cultures (Elling et al. 2015) found that $TEX_{86}$ ratios for two of the three cultures showed positive

136    relationships with growth temperature (but with different slopes and intercepts), while the third, an isolate

**Climate of the Past** Open Access

**Discussions**

EGU

137    from the surface waters of the South Atlantic, showed no relationship between $TEX_{86}$ and growth

138    temperature. The available evidence is that there is not a uniform response in GDGT-production to growth

139    temperature across distinct strains of archaea in culture (Elling et al., 2015). More fundamentally, in natural

140    systems, it is likely that aggregated GDGT abundance variations in response to growth temperatures result

141    from changing compositions of archaeal populations as well as the physiological response of individual

142    strains to growth temperature (Elling et al. 2015). For instance, a multiproxy study of Mediterranean

143    Pliocene-Pleistocene sapropels indicates that specific distributions of archaeal lipids might be reflective of

144    temporal changes in thaumarchaeael communities rather than temperature alone (Polik et al., 2018). Indeed,

145    the potential influence of community switching on GDGT composition can be seen in mesocosm studies,

146    with different species preferentially thriving at different growth temperatures (e.g., Schouten et al., 2007).

147    To use the responses of single, selected archaeal strains in culture to validate a particular model of

148    community-level responses to growth temperature is clearly problematic even in the modern system (Elling

149    et al., 2015). For deep time applications it is even more difficult, where there is no independent constraint

150    on the archaeal strains dominating production or their evolution through time (Elling et al. 2015). What is

151    notable, however, is that the Ring Index (RI) - calculated using all commonly measured GDGTs (Zhang et

152    al., 2016) – has a more robust relationship with culture temperature between archaeal strains than $TEX_{86}$,

153    indicating a potential loss of information within the $TEX_{86}$ index (Elling et al. 2015).

154

155    Finally, traditional uses of the $TEX_{86}$ proxy poorly represent the true uncertainty of palaeotemperature

156    estimates, as they include no assessment of non-analogue behavior relative to the modern core-top data.

157    Instead, uncertainty is typically based on the residuals on the modern calibration, with no reference to the

158    relationship between GDGT distributions of an ancient sample and the modern calibration data. An

159    improved Bayesian uncertainty model "BAYSPAR" is now in widespread use for SST estimation, which

160    uses sub-sampling approaches to improve temperature estimation and model uncertainty (Tierney and

161    Tingley, 2015). The Bayesian approach, however, still does not appropriately model uncertainty based on

162    the structure of fossil GDGT abundances relative to modern data, and is insensitive to detecting wildly non-

163    analogue behaviour in ancient GDGT distributions, as it still functions on one-dimensional $TEX_{86}$ index

164    values.

165

166    All empirical calibrations of GDGT-based proxies assume that mean annual SST is the master control on

167    GDGT assemblages both today and in the past. Mean annual SST, however, is strongly correlated with

168    many other environmental variables (e.g., seasonality, pH, mixed layer depth, and productivity). In the

169    modern calibration dataset, mean annual SST shows the strongest correlation with $TEX_{86}$ index (Schouten

170    et al., 2002), but this does not preclude an important (but undetectable) influence of these other

Climate
of the Past
Discussions

171   environmental variables. The use of empirical GDGT calibrations to infer ancient sea surface temperatures
172   thus implicitly assumes that the relationships between mean annual SST and all other GDGT-influencing
173   variables are invariant through time. This assumption is inescapable until, and unless, a more complete
174   biological mechanistic model of GDGT production emerges.

175

176   Here, we return to the primary modern core-top GDGT assemblage data (Tierney and Tingley, 2015), and
177   systematically explore the relationships between the modern GDGT distributions and surface ocean
178   temperatures using powerful mathematical tools. These tools can investigate correlations without prior
179   assumptions on the best form of relationship or *a priori* selection of GDGT compounds to be used. This
180   analysis is then extended through the exploration of the relationships between the modern core top GDGT
181   distributions and two compilations of ancient GDGT datasets, one from the Eocene (Inglis et al. 2015) and
182   one from the Cretaceous (O'Brien et al. 2017). We explore simple metrics to answer the fundamental
183   question – are modern core-top GDGT distributions good analogues for ancient distributions? We propose
184   the first robust methodology to answer this question, and so screen for significantly non-analogue palaeo-
185   assemblages. From this, we go on to derive a new machine learning approach 'OPTiMAL' (**O**ptimised
186   **P**alaeothermometry from **T**etraethers **v**i**a** **MA**chine **L**earning) for reconstructing SSTs from GDGT
187   datasets, which outperforms previous GDGT palaeothermometers and includes robust error estimates that,
188   for the first time, accounts for model uncertainty.

189

**2. Models for GDGT-based Temperature Reconstruction**

191

192   Our new analyses use the modern core-top data compilation, and satellite-derived estimates of SSTs, of
193   Tierney and Tingley (2015) as well as compilations of Eocene (Inglis et al. 2015) and Cretaceous (O'Brien
194   et al. 2017) GDGT assemblages. Within these fossil assemblages, only data points with full characterisation
195   of individual GDGT relative abundances were used. We also note that, in the first instance, all available
196   fossil assemblage data were interrogated, although later comparisons between BAYSPAR and our new
197   temperature predictor excludes fossil data that was regarded as unreliable based on standard pre-screening
198   indices, as noted within the original compilations (Inglis et al. 2015; O'Brien et al. 2017). All data used in
199   this study are tabulated in the supplementary information.

200

201   In order to enable meaningful comparison between new and existing temperature predictors, we use the
202   following consistent procedure for evaluating all predictors throughout this paper.  We divide the modern
203   core-top data set of 854 data points into 85 validation data points (chosen randomly) and 769 calibration
204   points.  We calibrate the predictor on the calibration points, and then judge its performance on the validation

205     points using the square root of the average of the square of the difference between the prediction at each

206     validation point, and the true (i.e. measured) temperature value:

207

208

$$\delta T = \sqrt{\frac{1}{N_v - 1} \sum_{k=1}^{N_v} (\hat{T}(x_k) - T(x_k))^2}$$

209                                                                                     (Eq. 5)

210

211     where the sum is taken over each of $N_v = 85$ validation points, $T$ is the known measured temperature (which

212     we refer to as the true temperature) and $\hat{T}$ is the predicted temperature.  For conciseness, we refer to $\delta T$ as

213     the predictor standard error.  It is useful to compare the accuracy of the predictor to the standard deviation

214     of all temperatures in the data set $\sigma T$, which corresponds to using the mean temperature as the predictor in

215     Equation 1; for the modern data set, $\sigma T = 10.0$ °C.  The so-called coefficient of determination $R^2$, given by

216

217

$$R^2 \equiv 1 - \left(\frac{\delta T}{\sigma T}\right)^2$$

218                                                                                     (Eq. 6)

219

220     provides a measure of the fraction of the fluctuation in the temperature explained by the predictor.  To

221     facilitate performance comparisons between different methods of predicting temperature, we use the same

222     subset of validation points for all analyses. To avoid sensitivity to the choice of validation points, we repeat

223     the calibration-validation procedure for 10 random choices from the validation dataset.

224

225     *2.1 Nearest neighbours*

226

227     We begin with an agnostic approach to using some combination of the six observables - GDGT-0, GDGT-

228     1, GDGT-2, GDGT-3, crenarchaeol and the crenarchaeol regio-isomer (cren'), which we will jointly refer

229     to as GDGTs - to predict sea surface temperatures. Whatever functional form the predictor might take, it

230     can only provide accurate temperature predictions if nearby points in the six-dimensional observable space

231     - i.e. the distribution of all of the six commonly reported GDGTs - can be translated to nearby points in

232     temperature space. Conversely, if nearby points in the observable space correspond to vastly different

233     temperatures, then no predictor, regardless of which combination of GDGTs are used, will be able to

234     provide a useful temperature estimate. In other words, the structuring of GDGT distributions within multi-

235     dimensional space, must have some correspondence to the temperatures of formation (or rather the mean

236     annual SSTs used for standard calibrations).

237

238     We therefore consider the prediction offered by the temperature at the nearest point in the GDGT parameter

239     space. Of course, nearness depends on the choice of the distance metric. For example, it may be that sea

240     surface temperatures are very sensitive to one observable, so even a small change in that observable

241     corresponds to a significant distance, and rather insensitive to another, meaning that even with a large

242     difference in the nominal value of that observable the distance is insignificant. In the first instance, we use

243     a very simple Euclidian distance estimate $D_{x,y}$ where the distance along each observable is normalised by

244     the total spread in that observable across the entire data set. This normalisation ensures that a dimensionless

245     distance estimate can be produced even when observables have very different dynamical ranges, or even

246     different units. Thus, the normalised distance $D$ between parameter data points $x$ and $y$ is

247

248 
$$D_{x,y}^2 \equiv \sum_{i=0}^{6} \frac{GDGT_i(x) - (GDGT_i(y))^2}{var(GDGT_i)}$$

249     (Eq. 7)

250

251     We show the distribution of nearest distances of points in the modern data set, excluding the sample itself,
252     in (Fig. 1).

253

254     The nearest-sample temperature predictor is $\hat{T}_{nearest}$ (x) = $T(y)$ where $y$ is the nearest point to $x$ over the

255     calibration data set, i.e., one that minimises $D_{x,y}$. Fig. 2 shows the scatter in the predicted temperature when

256     using the temperature of the nearest data point to make the prediction. Overall, the failure of the nearest-

257     neighbour predictor to provide accurate temperature estimates even when the normalised distance to the

258     nearest point is small, $D_{x,y} \leq 0.5$, casts doubt on the possibility of designing an accurate predictor for

259     temperature based on GDGT observations. This is most likely due to additional environmental controls on

260     GDGT abundance distributions in natural systems, in particular the water depth (Zhang and Liu, 2018),

261     nutrient availability (Hurley et al., 2018; Polik et al., 2018; Park et al., 2018), seasonality, growth rate

262     (Elling et al., 2014) and ecosystem composition (Polik et al., 2018), that obscure a predominant relationship

263     to mean annual SSTs.

264

265     On the other hand, the standard error for the nearest-neighbour temperature predictor is $\delta T_{nearest}$ = 4.5 ºC.

266     This is less than half of the standard deviation $\sigma T$ in the temperature values across the modern data set.

267     Thus, the temperatures corresponding to nearby points in GDGT observable space also cluster in

268     temperature space. Consequently, there is hope that we can make some useful, if imperfect, temperature

269    predictions. The value of $\delta T_{nearest}$ will also serve as a useful benchmark in this design: while we may hope

270    to do better by, say, suitably averaging over multiple nearby calibration points rather than adopting the

271    temperature at one nearest point as a predictor, any method that performs worse than the nearest-neighbour

272    predictor is clearly suboptimal.

273

274    *2.2 TEX$_{86}$ and Bayesian applications*

275

276    The TEX$_{86}$ index reduces the six-dimensional observable GDGT space to a single number. While this has

277    the advantage of convenience for manipulation and the derivation of simple analytic formulae for

278    predictors, as illustrated below, this approach has one critical disadvantage: it wastes significant information

279    embedded in the hard-earned GDGT distribution data. Fig. 3 illustrates both the advantage and

280    disadvantage of TEX$_{86}$. On the one hand, there is a clear correlation between TEX$_{86}$ and temperature (top

281    panel of Fig. 3), with a correlation coefficient of 0.81 corresponding to an overwhelming statistical

282    significance of $10^{-198}$. On the other hand, very similar TEX$_{86}$ values can correspond to very different

283    temperatures. We can apply the nearest-neighbour temperature prediction approach to the TEX$_{86}$ value

284    alone rather than the full GDGT parameter space; this predictor yields a large standard error of $\delta T_{nearestTEX86}$

285    = 8.0 ºC (bottom panel of Fig. 3). While smaller than $\sigma T$, this is significantly larger than $\delta T_{nearest}$ (Fig. 2),

286    consistent with the loss of information in TEX$_{86}$. We therefore do not expect other predictors based on

287    TEX$_{86}$ to perform as well as those based on the full available data set.

288

289    Indeed, this is what we find when we consider predictors of the form $\hat{T}_{1/TEX} = a + b/TEX_{86}$ and $\hat{T}_{TEXH} = c$

290    $+ d \log_{TEX86}$ (Liu et al., 2009; Kim et al., 2010), i.e., the established relationships between GDGT

291    distributions and SST. We fit the free parameters $a, b, c,$ and $d$ by minimising the sum of squares of the

292    residuals over the calibration data sets. We find that $\delta T_{1/TEX} = 6.1$ ºC (note that this is slightly better than

293    using the fixed values of $a$ and $b$ from (Kim et al., 2010), which yield $\delta T_{1/TEX} = 6.2$ ºC). We note that the

294    corresponding $R^2$ value associated with these TEX$_{86}$ based predictors is 0.64, which is lower than the $R^2$

295    values in Kim et al. (2010). We attribute this to the fact that we are using a larger dataset based on Tierney

296    and Tingley (2015), including data from the Red Sea (Kim et al. 2010).

297

298    Tierney and Tingley (2014) proposed a more sophisticated approach to obtaining the transfer function from

299    TEX$_{86}$ to temperature, continuing to use simple linear regression, but with the addition of Gaussian

300    processes to model spatial variability in the temperature-TEX$_{86}$ relationship and working with a forward

301    model which is subsequently inverted to produce temperature predictions. This forward model

302    'BAYSPAR' is capable of generating an infinite number of calibration curves relating TEX$_{86}$ to sea surface

303    temperatures (Tierney and Tingley, 2014). In order to derive a calibration for a specific dataset, the user

304    edits a range of parameters which vary depending on whether the dataset in question is from the relatively

305    recent past or deep time (Tierney and Tingley, 2014). For deep time applications, the authors propose a

306    modern analogue-type approach, in which they search the modern data for 20º x 20º grid boxes containing

307    `nearby' $TEX_{86}$ measurements and subsequently apply linear regression models calibrated on the analogous

308    samples for making predictions.

309

310    However, along with the simpler $TEX_{86}$-based models described above, this approach still suffers from the

311    arbitrary reduction of a six-dimensional data set to a single number. Therefore, it is not surprising that even

312    the simplest nearest-neighbour predictor (such as the one described above) that makes use of the full six-

313    dimensional dataset outperforms single-dimensional forward modelling approaches. Additionally,

314    uncertainty estimates do not account for the fact that $TEX_{86}$ is, fundamentally, an empirical proxy, and so

315    its validity outside the range of the modern calibration is not guaranteed. This is a fundamental issue for

316    attempts to reconstruct surface temperatures during Greenhouse climate states, when tropical and sub-

317    tropical SSTs were likely hotter than those observed in the modern oceans.

318

319    *2.3 Machine learning Approaches – Random Forests*

320

321    There are a number of options to improve on nearest-neighbour predictions using machine learning

322    techniques such as artificial neural networks and random forests. These flexible, non-parametric models

323    would ideally be based on the underlying processes driving the GDGT response to temperature, but since

324    these processes remain unconstrained at present, we choose to deploy models which can reasonably reflect

325    predictive uncertainty and will be sufficiently adaptable in future (as new information regarding controls

326    on GDGTs emerge). These machine learning approaches are all based on the idea of training a predictor by

327    fitting a set of coefficients in a sufficiently complex multi-layer model in order to minimise residuals on

328    the calibration data set. As an example of the power of this approach, we train a random forest of decision

329    trees with 100 learning cycles using a least-squares boosting to fit the regression ensemble. Figure 4 shows

330    the prediction accuracy for this random forest implementation. This machine learning predictor yields $\delta T$

331    = 4.1 ºC degrees, outperforming the naive nearest-neighbour predictor by effectively applying a suitable

332    weighted average over multiple near neighbours. This corresponds to a very respectable $R^2 = 0.83$, meaning

333    that 83% of the variation in the observed temperature is successfully explained by our GDGT-based model.

334

335    *2.4 Gaussian Process Regression*

336

337  One downside of the random forest predictor is the difficulty of accurately estimating the uncertainty on
338  the prediction (Mentch and Hooker, 2016), although this is possible with, e.g., a bootstrapping approach
339  (Coulston et al., 2016). Fortunately, Gaussian process (GP) regression provides a robust alternative. For
340  full details on GP regression refer to Williams and Rasmussen (20060 and Rasmussen and Nickisch (2010).
341  Loosely, the objective here is to search among a large space of smoothly varying functions of GDGT
342  compositions for those functions which adequately describe temperature variability. This, essentially, is a
343  way of combining information from all calibration data points, not just the nearest neighbours, assigning
344  different weights to different calibration points depending on their utility in predicting the temperature at
345  the input of interest. The trained Gaussian process learns the best choice of weights to fit the data. Typically,
346  the GP will give greater weight to closer points, but, as we discuss below, it will learn the appropriate
347  distance metric on the multi-dimensional GDGT input space.

348

349  The weighting coefficients learned by the GP emulator represent a covariance matrix on the GDGT
350  parameter space. We can use this as a distance metric to provide meaningfully normalised distances
351  between points, removing the arbitrariness from the nearest neighbour distance ($D_{x,y}$) definition used
352  earlier. If the temperature is insensitive to a particular GDGT input coordinate (i.e., the value of that input
353  has a minimal effect on the temperature) then points within GDGT space that have large differences in
354  absolute input values in that coordinate are still near. We find that Cren has very limited predictive power,
355  and so points with large Cren differences are close in term of the normalised distance. Conversely, if the
356  temperature is sensitive to small changes in a particular GDGT variant, then points with relatively nearby
357  absolute input values in that coordinate are still distant. We find that most GDGT parameters other than
358  Cren are comparably useful in predicting temperature, with GDGT-0 and GDGT-3 marginally the most
359  informative.

360

361  We use a Gaussian process model with a squared exponential kernel with automatic relevance
362  determination (ARD) to allow for a separate length scale for each GDGT predictor. We fit the GP
363  parameters with an optimiser based on quasi-Newton approximation to the Hessian. Prediction accuracy is
364  shown in Figure 5, and we find that $\delta T$ = 3.72 ºC, which is a substantial improvement over the existing
365  indices, at least on the modern data. As mentioned, the GP framework provides a natural quantification of
366  predictive uncertainty, which includes uncertainty about the learned function. This is in contrast to, for
367  example, the TEX$_{86}$ proxy, whereby the uncertainty associated with the selection of the particular functional
368  form used for predictions is ignored. While Tierney & Tingley (2014) also use Gaussian processes to model
369  uncertainty, they model spatial variability in the TEX$_{86}$-temperature relationship with a Gaussian process

Climate
of the Past
Discussions

370    prior. While this is a valuable approach to understand regional effects in the $TEX_{86}$-temperature

371    relationship, it does not deal with the `non-analogue' situations we are concerned with in this paper.

372

373    *2.5 Data Structure*

374

375    The random forest (Section 2.3) and GPR approaches (Section 2.4) are agnostic about any underlying bio-

376    physical model that might impart the observed temperature-dependence on GDGT relative abundances

377    produced by archaea. They are essentially optimized interpolation tools for mapping correlations between

378    temperature and GDGT abundances within the range of the modern calibration data set; they can make no

379    sensible inference about the behavior of this relationship outside of the range of this training data. To move

380    from interpolation within, to extrapolation beyond, the modern calibration requires an understanding of,

381    and model for, the temperature-dependence of GDGT production. To explore these relationships and the

382    extent to which the ancient and modern data reside in a coherent relationship within GDGT space, we

383    employed two forms of dimensionality reduction to enable visualisation of the data in two or three

384    dimensions. The fundamental point is that if temperature is the dominant control, all of the data should lie

385    approximately on a one-dimensional curve in GDGT space, and the arclength along this curve should

386    correspond to temperature; we will revisit this point below.

387

388    We first employed a version of principal component analysis (PCA) tailored to compositional data

389    (Aitcheson, 1982, 1983; Aitcheston and Greenacre, 2002; Filzmoser et al., 2009a; Filzmoser et al., 2009b;

390    Filzmoser et al., 2012). Taking into account the compositional nature of the data is important because the

391    sum-to-one constraint induces correlations between variables which are not accounted for by classical PCA.

392    Furthermore, apparently nonlinear structure in Euclidean space often corresponds to linearity in the simplex

393    (i.e. the restricted space in which all elements sum to one) (Egozcue et al., 2003). Figure 6 shows the

394    modern, Eocene and Cretaceous data projected onto the first two principal components. Aside from the

395    obvious outlying cluster of Cretaceous data, characterised by GDGT-3 fractions above 0.6, the bulk of the

396    data occupy a two-dimensional point cloud with a small amount of curvature. The large majority of the

397    Cretaceous data has more positive PC1 values relative to the modern data.

398

399    We also explored the data using diffusion maps (Coifman et al., 2005; Haghverdi et al., 2015), a nonlinear

400    dimensionality reduction tool designed to extract the dominant modes of variability in the data. Such

401    diffusion maps have been successfully used to infer latent variables that can explain patterns of gene

402    expression. In the case of biological organisms, this latent variable is commonly developmental age (called

403    pseudo-time) (Haghverdi et al., 2016). In our case, the assumption would be that this latent variable

Climate
of the Past
Discussions

404   corresponds to temperature. Inspection of the eigenvalues of the diffusion map transition matrix suggests

405   that four diffusion components are adequate to represent the data; we plot the second, third and fourth of

406   these components in Figure 7 for the modern and ancient data. The separate clusters marked `A' are the

407   outlying Cretaceous points with high GDGT-3 values. The bulk of the modern data lies on the branch

408   marked `B', while the bulk of the Cretaceous data lies on the branch marked `C'. Notably, the majority of

409   the modern points lying on branch C are from the Red Sea, which suggests that the Red Sea data is essential

410   for understanding ancient climates (particularly Cretaceous climates).

411

412   The relationship between the first diffusion component and $TEX_{86}$ for all data is shown in Figure 8. There

413   is a clear correlation, despite the presence of some outlying Cretaceous points, some of which are not shown

414   because they lie so far outside the majority data range within this projection. This suggests that $TEX_{86}$ is,

415   in one sense, a natural one-dimensional representation of the data. We also plot the first diffusion

416   component for the modern data as a function of temperature (Figure 9). We see a similar pattern emerging

417   to that displayed by $TEX_{86}$ - there is little sensitivity to temperature below 15 ºC, and between ~20 and 25

418   ºC. An interesting avenue for future research might be to explore the temperature-GDGT system from a

419   dynamical systems perspective, i.e. use simple mechanistic mathematical models to explore the

420   temperature-dependence of steady-state GDGT distributions. It may be that such models suggest that only

421   a few steady-states exist, and that temperature is a bifurcation parameter, i.e. it controls the switch between

422   the steady states. Note also the downward slope in the residual pattern in Figure 4 between 0 and 15-17

423   degrees celsius, and again at higher temperatures. This pattern is consistent with predictions that are biased

424   towards the centre of each `cluster', i.e. a system which is not very sensitive to temperature, but can

425   distinguish between high and low temperatures reasonably well. This observation also links to recent culture

426   studies (Elling et al., 2015) and Pliocene-Pleistocene sapropel data (Polik et al., 2018), which support the

427   existence of discrete populations with unique GDGT-temperature relationships and that temporal changes

428   in population over time can drive changes in $TEX_{86}$.

429

430   *2.6 Forward Modelling*

431

432   Based on the analysis of the combined modern and ancient data structure outlined above, there appears to

433   be some consistency to underlying trends in the overall variance of GDGT relative abundances. These

434   trends provide some hope that models of this variance, and its relationship to sea surface temperature, within

435   the modern dataset could be developed to predict ancient SSTs. $TEX_{86}$ and BAYSPAR are such models,

436   but they are limited by, first, the reduction of six-dimensional GDGT space to a one-dimensional index;

437   and second, by an *ad hoc* model choice – linear, exponential – that does not account for uncertainty in

438 model fit to the modern calibration data, and the resultant uncertainty in the estimation of ancient SSTs

439 relating to model choice. To overcome these issues, we develop a forward model based on a multi-output

440 Gaussian Process (Alvarez et al., 2012), which models GDGT compositions as functions of temperature,

441 accounting for correlations between GDGT measurements. This model is then inverted to obtain

442 temperatures which are compatible with a measured GDGT composition. In simple terms, we posit that a

443 measured GDGT composition is generated by some unknown function of temperature and corrupted by

444 noise, which may be due to measurement error or some unmodelled particularity of the environment in

445 which the sample was generated. We proceed by defining a large (in this case infinite) set of functions of

446 temperature to explore and compare them to the available data, throwing away those functions which do

447 not adequately fit the data. This means, of course, that the behaviour of the functions we accept is allowed

448 to vary more widely outside the range of the modern data than within it. With no mechanistic underpinning,

449 choosing only one function (such as the inverse of $TEX_{86}$) based on how well it fits the modern data grossly

450 underestimates our uncertainty about temperature where no modern analogue is available.

451

452 The forward modelling approach is similar to that of Haslett et al. (2006), who argue that it is preferable to

453 model measured compositions as functions of climate, before probabilistically inverting the model to infer

454 plausible climates given a composition. The cost of modelling the data in this more natural way is the loss

455 of degrees of freedom -- we are now attempting to fit a one-dimensional line through a multidimensional

456 point cloud rather than fit a multidimensional surface to the GDGT data, which means that the predictive

457 power of the model suffers, at least on the modern data. The existing BAYSPAR calibration also specifies

458 the model in the forward direction, but ignores model uncertainty. As with all GP models, the choice of

459 kernel has a substantial impact on predictions (and their associated uncertainty) outside the range of the

460 modern data, where predictions revert to the prior implied by the kernel. Given that we have no mechanistic

461 model for the data generating process, we recommend the use of kernels which do not impose strong prior

462 assumptions on the form of the GDGT-temperature relationship (e.g. kernels with a linear component) and

463 thus reasonably represent model uncertainty outside the range of the modern data. We choose a zero-mean

464 Matern 3/2 kernel for the applications below. Note, however, that since we are working in ilr-transformed

465 coordinates, this corresponds to a prior assumption of uniform compositions at all temperatures, i.e. all

466 components are equally abundant.

467

468 The residuals for the forward model are shown in Figure 10. The clear pattern in the residuals does not

469 necessarily indicate model misspecification, since no explicit noise model is specified for temperatures.

470 Predictive distributions are to be interpreted in the Bayesian sense, in that they represent a 'degree of belief'

471 in temperatures given the model and the modern data. The residual pattern is similar to that of the random

472  forest (Figure 4) with two clear downward slopes, suggesting again that the data are clustered into

473  temperatures above and below 16-17 degrees celsius, and that predictions tend towards temperatures at the

474  centres of these clusters.

475

476  An advantage of the forward modelling approach is that the inversion can incorporate substantive prior

477  information about temperatures for individual data points. In particular, other proxy systems can be used to

478  elicit prior distributions over temperatures to constrain GDGT-based predictions, particularly when

479  attempting to reconstruct ancient climates with no modern analogue in GDGT-space. We emphasise that

480  outside the range of the modern data, the utility of the models is almost solely due to the prior information

481  included in the reconstruction. At present, the only priors being used in the forward model prescribe a

482  reasonable upper limit and lower limit on temperatures (see Supplementary Information). The only way to

483  improve these reconstructions will be for future iterations to incorporate prior information from other

484  proxies. It is worth noting that the predictive uncertainty, while reasonably well-described by the standard

485  deviation in cases where ancient data lie quite close to the modern data in GDGT space, can be highly

486  multimodal (Fig. 11). This is the case when estimates are significantly outside of the modern calibration

487  dataset, such as low latitude data in the Cretaceous, or where there is considerable scatter in the modern

488  calibration data, for example in the low temperature range (<5 °C).

489

490  **3. Non-analogue behavior and Extrapolation**

491

492  In principle, the predictors described above can be applied directly to ancient data, such as data from the

493  Eocene or Cretaceous (Inglis et al., 2015; O'Brien et al., 2017).  In practice, one should be careful with

494  using models outside their domain of applicability.  The machine learning tools described above, which are

495  ultimately based on the analysis of nearby calibration data in GDGT space, are fundamentally designed for

496  *interpolation*.  To the extent that ancient data occupy a very different region in GDGT space, *extrapolation*

497  is required, which the models do not adequately account for. The divergence between modern calibration

498  data and ancient data is evident from Fig. 12, which shows histograms of minimum normalised distances

499  between 'high quality' Eocene/Cretaceous data points (those that passed the screening tests applied by

500  O'Brien et al., 2017 and Inglis et al., 2015) and the nearest point in the full modern data set. We strongly

501  recommend the use of the nearest neighbor distance metric ($D_{nearest}$) as a screening method to determine

502  whether the modern core top GDGT assemblage data is an appropriate basis for ancient SST estimation on

503  a case-by-case basis. Note that this distance measure is weighted by the scale length of the relevant

504  parameter as estimated by the Gaussian process emulator in order to quantify the relative position of ancient

505  GDGT assemblages to the modern core-top data.  By using the GP-estimated covariance as the distance

506  metric, we account for the sensitivity of different GDGT components to temperature. Our inference is that

507  samples with $D_{nearest}$ >0.5, *regardless of the calibration model or approach applied*, are unlikely to generate

508  temperature estimates that are much better than informed guesswork. In these instances, in both our GPR

509  and Fwd models, the constraints provided by the modern calibration data set are so weak that estimates of

510  temperature have large uncertainty bands that are dictated by model priors; i.e. are unconstrained by the

511  calibration data (e.g., Figure 13 and Figure 14). This uncertainty is not apparent from estimates generated

512  by BAYSPAR or $TEX_{86}^{H}$ models, although the underlying and fundamental lack of constraints are the same.

513  While 93% of validation data points in the modern data have $D_{nearest}$ <0.5, this is the case for only 33% of

514  Eocene samples and 3% for Cretaceous samples.

515

516  Where ancient GDGT distributions lie far from the modern calibration data set ($D_{nearest}$ >0.5), we argue that

517  there is no suitable set of modern analogue GDGT distributions from which to infer growth temperatures

518  for this ancient GDGT distribution.  Both the GPR and Fwd models revert to imposed priors once the

519  distance from the modern calibration dataset increases.  We propose that this is more rigorous and justified

520  model behavior than extrapolation of $TEX_{86}$ or BAYSPAR predictors to non-analogue samples far from

521  the modern calibration data.  As a result, the predictive models can only be applied to a subset of the Eocene

522  and Cretaceous data. We also note that there are two broad, non-mutually-exclusive categories of samples

523  that lie far from the modern calibration dataset ($D_{nearest}$ >0.5), the first are samples that seem to lie 'beyond'

524  the temperature-GDGT calibration relationship, likely with (unconstrained) GDGT formation temperatures

525  higher than the modern core-top calibrations; the second are samples with anomalous GDGT distributions

526  lying on the margins of, or far away from the main GDGT clustering in 6-dimensional space (see outliers

527  in Fig. 8).

528

529  Given the (current) limit on natural mean annual surface ocean temperatures of ~30 ºC, extending the

530  GDGT-temperature calibration might be possible through, 1) integration of full GDGT abundance

531  distributions produced in high temperature culture, mesocosm or artificially warmed sea surface conditions

532  into the models; followed by, 2) validation through robust inter-comparisons of any new GDGT

533  palaeothermometer for high temperatures conditions with other temperature proxies from past warm

534  climate states. As discussed in the introduction, the first approach is limited by the ability of culture or

535  mesocosm experiments to accurately represent the true diversity and growth environments and dynamics

536  of natural microbial populations. Such studies clearly indicate a more complex, community-scale control

537  on changing GDGT relative abundances to growth temperatures (e.g., Elling et al., 2015). Community-

538  scale temperature dependency can be modelled relatively well with analyses of natural production preserved

539  in core-top sediments, especially with more sophisticated model fitting, including the GPR and Fwd model

540    presented here. Above ~30ºC, however, the behavior of even single strains of archaea are not well-
541    constrained by culture experiments, and the natural community-level responses above this temperature are,
542    so far, completely unknown. Significantly more culture and mesocosm data at these high temperatures,
543    spanning a range of microbial diversity and growth conditions, could provide some of these constraints in
544    future. Until such data exist, we see no robust justification for any particular extrapolation of modern core-
545    top calibration data sets into the unknown above 30 ºC, although the coherent patterns apparent across
546    GDGT space, between modern, Eocene and Cretaceous data (Figures 7), does provide some grounds for
547    hope that the extension of GDGT palaeothermometry beyond 30ºC might be possible in future.

548

549    **4. OPTiMAL and $D_{nearest}$: A more robust method for GDGT-based paleothermometry**

550

551    A more robust framework for GDGT-based palaeothermometry, could be achieved with a flexible
552    predictive model that uses the full range of six GDGT relative abundances, and has transparent and robust
553    estimates of the prediction uncertainty. In this context, the Gaussian Process Regression model (GPR;
554    Section 2.4) outperforms the Forward model (Fwd; Section 2.6) within the modern calibration dataset and
555    we recommend standard use of the GPR model, henceforth called OPTiMAL, over the Fwd model. Model
556    code for the calculation of $D_{nearest}$ values and OPTiMAL SST estimates (Matlab script) and the Fwd Model
557    SST    estimates    (R    script)    are    archived    in    the    GITHUB    repository,
558    https://github.com/carbonatefan/OPTiMAL.

559

560    To investigate the behaviour of the new OPTiMAL model, we compare temperature predictions including
561    uncertaintities for the Eocene and Cretaceous datasets, made by OPTiMAL and the BAYSPAR
562    methodology of Tierney and Tingley (2014) (Figures 13 and 14). The OPTiMAL model systematically
563    estimates slightly cooler temperatures than BAYSPAR, with the biggest offsets below ~15 ºC (Figure 13).
564    Fossil GDGT assemblages that fail the $D_{nearest}$ test are shown in grey, which clearly illustrate the regression
565    to the mean in the OPTiMAL model, whereas BAYSPAR continues to make SST predictions up to and
566    exceeding 40 ºC for these "non-analogue" samples. A comparison of error estimation between OPTiMAL
567    and BAYSPAR is shown in Figure 14.  For most of the predictive range below the $D_{nearest}$ cut-off of 0.5,
568    OPTiMAL has smaller errors than BAYSPAR, especially in the lower temperature range. As $D_{nearest}$
569    increases, i.e. as the fossil GDGT assemblage moves further from the constraints of the modern calibration
570    dataset, the error on OPTiMAL increases, until it reaches the standard deviation of the modern calibration
571    dataset (i.e., is completely unconstrained). In other words, OPTiMAL generates maximum likelihood SSTs
572    with robust confidence intervals, which appropriately reflect the relative position of an ancient sample used
573    for SST estimation and the structure of the modern calibration data set. Where there are strong constraints

574    from near analogues in the modern data, uncertainties will be small, where there are weak constraints,

575    uncertainty increases. In contrast, BAYSPAR, because it is fundamentally based on a *parametric* linear

576    model and therefore does not account for model uncertainty, assigns similar uncertainty intervals as to the

577    rest of the data, despite there being no way of reasonably testing whether the linear model is an appropriate

578    description of the data far from the modern dataset.

579

580    **5. Conclusions**

581

582    Although the fundamental issue of non-analogue behaviour is a key problem for GDGT-temperature

583    estimation, it has an undue impact on the community's general confidence in this method. In part, this is

584    because these issues have not been clearly stated and circumscribed - rather they have been allowed to erode

585    confidence in the entire GDGT-based methodology through the inappropriate use of GDGT-based

586    palaeothermometry far outside the modern constraints on the behavior of this system. The use of GDGT

587    abundances to estimate temperatures in clearly non-analogue conditions is, at present, difficult to justify on

588    the basis of the available calibration constraints or a good understanding of underlying biophysical models.

589    We hope that this study prompts further investigations that will improve these constraints for the use of

590    GDGTs in deep-time paleoclimate studies, where they clearly have substantial potential as temperature

591    proxies. Temperature estimates based on fossil GDGT assemblages that are within range of, or similar to,

592    modern GDGT calibration data, do, however, rest on a strong, underlying temperature-dependence

593    observed in the empirical data The failure to have an effective means of separating the "good from the bad",

594    either leads to false confidence and inappropriate inferences in non-analogue conditions, or a false

595    pessimism when the community's trust in the overall method is eroded by the clear influence of

596    methodological choice on SST estimates outside of the modern calibration range.

597

598    In this study, we apply modern machine-learning tools, including Gaussian Process Emulators and forward

599    modelling, to improve temperature estimation and the representation of uncertainty in GDGT-based SST

600    reconstructions. Using our new nearest neighbour test, we demonstrate that >60% of Eocene, and >90% of

601    Cretaceous, fossil GDGT distribution patterns differ so significantly from modern that it is inappropriate to

602    interpret them using modern empirical calibrations of any formulation. For data that does show sufficient

603    similarity to modern, we present OPTiMAL, a new multi-dimensional Gaussian Process Regression tool

604    which uses all six GDGTs (GDGT-0, -1, -2, -3, Cren and Cren') to generate an SST estimate with associated

605    uncertainty. The key advantages of the OPTiMAL approach are: 1) that these uncertainty estimates are

606    intrinsically linked to the strength of the relationship between the fossil GDGT distributions and the modern

607    calibration data set, and 2) by considering all GDGT compounds in a multi-dimensional regression model

Climate
of the Past
Discussions

Open Access

EGU

608   it avoids the dimensionality reduction and loss of information that takes place when calibrating single

609   parameters (TEX$_{86}$) to temperature. The methods presented above make very few assumptions about the

610   data. We argue that such methods are appropriate with the current absence of any reasonable mechanistic

611   model for the data generating process, in that they reflect model uncertainty in a natural way.

612

613

620

621

622

623   **Figure Captions:**

624

625

626   **Figure 1**. A histogram of the normalised distance to the nearest neighbour in GDGT space ($D_{x,yt}$) for all

627   samples in the modern calibration dataset of Tierney and Tingley (2015).

628

629   **Figure 2.** The error of the nearest-neighbour temperature ($D_{x,y}$) predictor, for modern core-top data, as a

630   function of the distance to the nearest calibration sample.

631

632   **Figure 3.** Top: The temperature of the modern data set as a function of the TEX$_{86}$ value, showing a clear

633   correlation between the two, but also significant scatter.  Bottom: the error of the predictor based on the

634   nearest TEX$_{86}$ calibration point.

635

636   **Figure 4.** The error of a random forest predictor as a function of the true temperature.

637

638   **Figure 5**. The error of the GPR (Gaussian Process regression) predictor as a function of the true

639   temperature.

640

641   **Figure 6.** Modern and ancient data projected onto the first two compositional principal components. Black:

642   Modern; Blue: Eocene (Inglis et al., 2015); Red: Cretaceous (O'Brien et al., 2017).

643

644    **Figure 7.** Diffusion map projection of the modern and ancient data. Black: Modern; Blue: Eocene (Inglis

645    et al., 2015); Red: Cretaceous (O'Brien et al., 2017). separate clusters marked `A' are the outlying

646    Cretaceous points with high GDGT-3 values. Branch 'B' is dominated by modern data points; branch 'C'

647    by Cretaceous data.

648

649    **Figure 8.** The first diffusion component as a function of $TEX_{86}$ . Some outlying points have been excluded

650    from the plot for the purposes of visualisation. Black: Modern; Blue: Eocene (Inglis et al., 2015); Red:

651    Cretaceous (O'Brien et al., 2017).

652

653    **Figure 9.** The first diffusion component as a function of temperature (modern data only).

654

655    **Figure 10.** Temperature residuals for the forward model.

656

657    **Figure 11.** The posterior distributions over temperature from the forward model for selected examples of

658    high and low temperature, Eocene and Cretaceous, data points. The Gaussian error envelope from the GPR

659    model is shown for comparison.

660

661    **Figure 12.** A histogram of normalised distances to the nearest sample in the modern data set for Eocene

662    and Cretaceous data, excluding samples that had been screened out in previous compilations using BIT, MI

663    and RI following the approach of (Inglis et al., 2015; O'Brien et al., 2017).

664

665    **Figure 13.** Comparison of temperature estimates for the BAYSPAR and the OPTiMAL GPR model, greyed

666    out data fails the $D_{nearest}$ test (>0.5), and the colour scaling reflects $D_{nearest}$ values for those datapoints that

667    pass

668

669    **Figure 14.** Inter-comparison of temperature estimates (top) and errors (bottom) for the Eocene and

670    Cretaceous data calculated using BAYSPAR and OPTiMAL. Greyed out data fails the $D_{nearest}$ test (>0.5),

671    and the colour scaling reflects $D_{nearest}$ values for those datapoints that pass. The black dashed line shows the

672    $D_{nearest}$ threshold (>0.5).

673

674

675

676

Climate
of the Past
Discussions

## References:

Aitchison, J.: The Statistical Analysis of Compositional Data. J. R. Stat. Soc. Series B Stat. Methodol. 44, 139–160, 1982.

Aitchison, J.: Principal component analysis of compositional data. Biometrika 70, 57–65, 1983.

Aitchison, J., Greenacre, M.: Biplots of compositional data. J. R. Stat. Soc. Ser. C Appl. Stat. 51, 375–392, 2002.

Álvarez, M.A., Rosasco, L., Lawrence, N.D.: Kernels for Vector-Valued Functions: A Review. Foundations and Trends® in Machine Learning 4, 195–266, 2012.

Bijl, P. K., S. Schouten, A. Sluijs, G.-J. Reichart, J. C. Zachos, and H. Brinkhuis.: Early Palaeogene temperature evolution of the southwest Pacific Ocean. Nature, 461, 776–779, 2009.

Bijl, P. K., Bendle, J.A.P., Bohaty, S.M., Pross, J., Schouten, S., Tauxe, L., Stickley, C., McKay, R.M., Röhl, U., Olney, M., Slujis, A., Escutia, C., Brinkhuis, H. and Expedition 318 Scientists.: Eocene cooling linked to early flow across the Tasmanian Gateway. Proc. Natl. Acad. Sci. U.S.A., 110, 9645–9650, 2013.

Brassell, S. C.: Climatic influences on the Paleogene evolution of alkenones, Paleoceanography, 29, 255-272, doi:10.1002/2013PA002576, 2014.

Brinkhuis, H., Schouten, S., Collinson, M. E., Sluijs, A., Damsté, J. S. S., Dickens, G. R., Huber, M., Cronin, T. M., Onodera, J., Takahashi, K., Bujak, J. P., Stein, R., van der Burgh, J., Eldrett, J. S., Harding, I. C., Lotter, A. F., Sangiorgi, F., Cittert, H. v. K.-v., de Leeuw, J. W., Matthiessen, J., Backman, J., Moran, K., and the Expedition, Scientists.: Episodic fresh surface waters in the Eocene Arctic Ocean, Nature, 441, 606 – 609, 2006.

Coifman, R.R., Lafon, S., Lee, A.B., Maggioni, M., Nadler, B., Warner, F., Zucker, S.W.: Geometric diffusions as a tool for harmonic analysis and structure definition of data: diffusion maps. Proc. Natl. Acad. Sci. U. S. A. 102, 7426–7431, 2005.

Coulston, J.W., Blinn, C.E., Thomas, V.A.: Approximating prediction uncertainty for random forest regression models.  Photogrammetric Engineering & Remote Sensing, Volume 82, 189-197, https://doi.org/10.14358/PERS.82.3.189, 2016.

Cramwinckel, M. J., Huber, M., Kocken, I. J., Agnini, C., Bijl, P. K., Bohaty, S. M., Frieling, J., Goldner, A., Hilgen, F. J., Kip, E. L., Peterse, F., van der Ploeg, R., Röhl, U., Schouten, S., and Sluijs, A.: Synchronous tropical and polar temperature evolution in the Eocene, Nature, 559, 382-386, 2018.

Douglas, P. M. J., Affek, H. P., Ivany, L. C., Houben, A. J. P., Sijp, W. P., Sluijs, A., Schouten, S., Pagani, M.: Pronounced zonal heterogeneity in Eocene southern high-latitude sea surface temperatures. Proceedings of the National Academy of Sciences, 111, 6582–6587, 2014.

Dunkley Jones, T., Lunt, D. J., Schmidt, D. N., Ridgwell, A., Sluijs, A., Valdes, P. J., and Maslin, M.: Climate model and proxy data constraints on ocean warming across the Paleocene–Eocene Thermal Maximum, Earth-Science Reviews, 125, 123-145, 2013.

714  Egozcue, J.J., Pawlowsky-Glahn, V., Mateu-Figueras, G., Barceló-Vidal, C.: Isometric Logratio
715       Transformations for Compositional Data Analysis. Math. Geol. 35, 279–300, 2003.

716  Elling, F. J., Könneke, M., Lipp, J. S., Becker, K. W., Gagen, E. J., and Hinrichs, K.-U.: Effects of growth
717       phase on the membrane lipid composition of the thaumarchaeon Nitrosopumilus maritimus and
718       their implications for archaeal lipid 20 distributions in the marine environment, Geochimica et
719       Cosmochimica Acta, 141, 579-597, 2014.

720  Elling, F. J., Könneke, M., Mußmann, M., Greve, A., and Hinrichs, K.-U.: Influence of temperature, pH,
721       and salinity on membrane lipid composition and $TEX_{86}$ of marine planktonic thaumarchaeal
722       isolates, Geochimica et Cosmochimica Acta, 171, 238-255, 2015.

723  Filzmoser, P., Hron, K., Reimann, C.: Principal component analysis for compositional data with outliers.
724       Environmetrics 20, 621–632, 2009a.

725  Filzmoser, P., Hron, K., Reimann, C., Garrett, R.: Robust factor analysis for compositional data. Comput.
726       Geosci. 35, 1854–1861, 2009b.

727  Filzmoser, P., Hron, K., Reimann, C.: Interpretation of multivariate outliers for compositional data.
728       Comput. Geosci. 39, 77–85, 2012.

729  Haghverdi, L., Buettner, F., Theis, F.J: Diffusion maps for high-dimensional single-cell analysis of
730       differentiation data. Bioinformatics 31, 2989–2998, 2015.

731  Haghverdi, L., Büttner, M., Wolf, F.A., Buettner, F., Theis, F.J.: Diffusion pseudotime robustly
732       reconstructs lineage branching. Nat. Methods 13, 845–848, 2016.

733  Haslett, J., Whiley, M., Bhattacharya, S., Salter-Townshend, M., Wilson, S.P., Allen, J.R.M., Huntley, B.,
734       Mitchell, F.J.G.: Bayesian palaeoclimate reconstruction. J Royal Statistical Soc A 169, 395–438;,
735       2006.

736  Herfort, L., Schouten, S., Boon, J. P., and Sinninghe Damsté, J. S.: Application of the $TEX_{86}$ temperature
737       proxy to the southern North Sea, Organic Geochemistry, 37, 1715-1726, 2006.

738  Hertzberg, J. E., Schmidt, M. W., Bianchi, T. S., Smith, R. K., Shields, M. R., & Marcantonio, F.:
739       Comparison of eastern tropical Pacific $TEX_{86}$ and Globigerinoides ruber Mg/Ca derived sea surface
740       temperatures: Insights from the Holocene and Last Glacial Maximum. Earth and Planetary Science
741       Letters, 434, 320–332, 2016.

742  Hollis, C. J., Taylor, K. W. R., Handley, L., Pancost, R. D., Huber, M., Creech, J. B., Hines, B. R., Crouch,
743       E. M., Morgans, 25 H. E. G., Crampton, J. S., Gibbs, S., Pearson, P. N., and Zachos, J. C.: Early
744       Paleogene temperature history of the Southwest Pacific Ocean: Reconciling proxies and models,
745       Earth and Planetary Science Letters, 349–350, 53-66, 2012.

746  Hollis, C. J., Dunkley Jones, T., Anagnostou, E., Bijl, P. K., Cramwinckel, M. J., Cui, Y., Dickens, G. R.,
747       Edgar, K. M., Eley, Y., Evans, D., Foster, G. L., Frieling, J., Inglis, G. N., Kennedy, E. M.,
748       Kozdon, R., Lauretano, V., Lear, C. H., Littler, K., Meckler, N., Naafs, B. D. A., Pälike, H.,
749       Pancost, R. D., Pearson, P., Royer, D. L., Salzmann, U., Schubert, B., Seebeck, H., Sluijs, A.,
750       Speijer, R., Stassen, P., Tierney, J., Tripati, A., Wade, B., Westerhold, T., Witkowski, C., Zachos,
751       J. C., Zhang, Y. G., Huber, M., and Lunt, D. J.: The DeepMIP contribution to PMIP4:
752       methodologies for selection, compilation and analysis of latest Paleocene and early Eocene climate

753    proxy data, incorporating version 0.1 of the DeepMIP database, Geosci. Model Dev. Discuss.,
754    https://doi.org/10.5194/gmd-2018-309, in review, 2019.

755    Hollis, C. J., Handley, L., Crouch, E. M., Morgans, H. E., Baker, J. A., Creech, J., Collins, K. S., Gibbs, S.
756    J., Huber, M., Schouten, S.: Tropical sea temperatures in the high-latitude South Pacific during the
757    Eocene. Geology, 37, 99–102, 2009.

758    Hopmans, E. C., Weijers, J. W. H., Schefuss, E., Herfort, L., Sinninghe Damsté, J. S., Schouten, S.: A
759    novel proxy for terrestrial organic matter in sediments based on branched and isoprenoid tetraether
760    lipids, Earth and Planetary Science Letters, 224, 107-116, 2004.

761    Huguet C, Kim J-H, Sinninghe Damste´ J.S., Schouten S: Reconstruction of sea surface temperature
762    variations in the Arabian Sea over the last 23 kyr using organic proxies ($TEX_{86}$ and UK 0 37 .
763    Paleoceanography 21(3): PA3003, 2006.

764    Hurley, S. J., Elling, F. J., Könneke, M., Buchwald, C., Wankel, S. D., Santoro, A. E., Lipp. J.S., Hinrichs,
765    K., Pearson, A.: Influence of ammonia oxidation rate on thaumarchaeal lipid composition and the
766    $TEX_{86}$ temperature proxy. Proceedings of the National Academy of Sciences, 113, 7762–7767,
767    2016.

768    Inglis, G. N., Farnsworth, A., Lunt, D., Foster, G. L., Hollis, C. J., Pagani, M., Jardine, P. E., Pearson, P.
769    N., Markwick, P., Galsworthy, A. M. J., Raynham, L., Taylor, K. W. R., and Pancost, R. D.:
770    Descent toward the Icehouse: Eocene sea surface cooling inferred from GDGT distributions,
771    Paleoceanography, 30, 1000-1020, 2015.

772    Jenkyns H.C., Schouten-Huibers L., Schouten S., Damsté J.S.S.: Warm Middle Jurassic-Early Cretaceous
773    high-latitude sea-surface temperatures from the Southern Ocean. Clim Past 8 (1):215–226, 2012.

774    Kim, J.-H., Schouten, S., Hopmans, E. C., Donner, B., and Sinninghe Damsté, J. S.: Global sediment core-
775    top calibration of the $TEX_{86}$ paleothermometer in the ocean, Geochimica et Cosmochimica Acta,
776    72, 1154-1173, 2008.

777    Kim, J.-H., van der Meer, J., Schouten, S., Helmke, P., Willmott, V., Sangiorgi, F., Koç, N., Hopmans, E.
778    C., and Sinninghe Damsté, J. S.: New indices and calibrations derived from the distribution of
779    crenarchaeal isoprenoid tetraether lipids: Implications for past sea surface temperature
780    reconstructions, Geochimica et Cosmochimica Acta, 74, 4639-4654, 2010.

781    Linnert, C., Robinson, S. A., Lees, J. A., Bown, P. R., Perez-Rodriguez, I., Petrizzo, M. R., Falzoni, F.,
782    Littler, K., Antonio Arz, J., Russell, E. E. : Evidence for global cooling in the Late Cretaceous.
783    Nature Communications, 5, 1–7, 2014.

784    Lunt, D. J., Dunkley Jones, T., Heinemann, M., Huber, M., LeGrande, A., Winguth, A., Loptson, C.,
785    Marotzke, J., Tindall, J., 15 Valdes, P., Winguth, C.: A model-data comparison for a multi-model
786    ensemble of early Eocene atmosphere-ocean simulations: EoMIP, Clim. Past Discuss., 8, 1229-
787    1273, 2012.

788    Mentch, L., Hooker, G.: Quantifying Uncertainty in Random Forests via Confidence Intervals and
789    Hypothesis Tests. Journal of Machine Learning Research, 17, 1-41, 2016.

790    O'Brien, C. L., Robinson, S. A., Pancost, R. D., Sinninghe Damsté, J. S., Schouten, S., Lunt, D. J., Alsenz,
791    H., Bornemann, 20 A., Bottini, C., Brassell, S. C., Farnsworth, A., Forster, A., Huber, B. T., Inglis,

792    G. N., Jenkyns, H. C., Linnert, C., Littler, K., Markwick, P., McAnena, A., Mutterlose, J., Naafs, B.
793        D. A., Püttmann, W., Sluijs, A., van Helmond, N. A. G. M., Vellekoop, J., Wagner, T., Wrobel, N.
794        E.: Cretaceous sea-surface temperature evolution: Constraints from $TEX_{86}$ and planktonic
795        foraminiferal oxygen isotopes, Earth-Science Reviews, 172, 224-247, 2017.

796    Park, E., Hefter, J., Fischer, G., Mollenhauer, G.: $TEX_{86}$ in sinking particles in three eastern Atlantic
797        upwelling regimes. Organic Geochemistry, 124, 151–163, 2018.

798    Pearson, P. N., van Dongen, B. E., Nicholas, C. J., Pancost, R. D., Schouten, S., Singano, J. M., Wade, B.
799        S.: Stable warm tropical climate through the Eocene Epoch. Geology, 35, 211-214, 2007.

800    Polik, C. A., Elling, F. J., Pearson, A.: Impacts of Paleoecology on the $TEX_{86}$ Sea Surface Temperature
801        Proxy in the Pliocene-Pleistocene Mediterranean Sea. Paleoceanography and Paleoclimatology, 33,
802        1472–1489, 2018.

803    Qin, W., Amin, S. A., Martens-Habbena, W., Walker, C. B., Urakawa, H., Devol, A. H., Ingalls, A.E.,
804        Moffett, J.W., Ambrust, E.V., Stahl, D. A.: Marine ammonia-oxidizing archaeal isolates display
805        obligate mixotrophy and wide ecotypic variation. Proceedings of the National Academy of
806        Sciences of the United States of America, 111, 12504–12509, 2014.

807    Qin, W., Carlson, L. T., Armbrust, E. V., Stahl, D. A., Devol, A. H., Moffett, J. W., Ingalls, A. E.:
808        Confounding effects of oxygen and temperature on the TEX 86 signature of marine
809        Thaumarchaeota. Proceedings of the National Academy of Sciences, 112, 10979–10984, 2015.

810    Rasmussen, C.E., Nickisch, H.: Gaussian Processes for Machine Learning (GPML) Toolbox. J. Mach.
811        Learn. Res. 11, 3011–3015, 2010.

812    Sangiorgi, F., van Soelen Els, E., Spofforth David, J. A., Pälike, H., Stickley Catherine, E., St. John, K.,
813        Koç, N., Schouten, S., Sinninghe Damsté Jaap, S., Brinkhuis, H.: Cyclicity in the middle Eocene
814        central Arctic Ocean sediment record: Orbital forcing and environmental response,
815        Paleoceanography, 23, 10.1029/2007PA001487, 2008.

816    Schouten, E., Hopmans, E.C., Forster, A., Van Breugel, Y., Kuypers, M.M.M., Sinninghe Damsté, J.S.:
817        Extremely high seasurface temperatures at low latitudes during the middle Cretaceous as revealed
818        by archaeal membrane lipids. Geology, 31, 1069–1072, 2003.

819    Schouten, S., Forster, A., Panoto, F. E., and Sinninghe Damsté, J. S.: Towards calibration of the $TEX_{86}$
820        palaeothermometer for 20 tropical sea surface temperatures in ancient greenhouse worlds, Organic
821        Geochemistry, 38, 1537-1546, 2007.

822    Schouten, S., Hopmans, E. C., Sinninghe Damsté, J. S.: The organic geochemistry of glycerol dialkyl
823        glycerol tetraether lipids: A review, Organic Geochemistry, 54, 19-61, 2013.

824    Schouten, S., Hopmans, E. C., Schefuß, E., Sinninghe Damsté, J. S.: Distributional variations in marine
825        crenarchaeotal membrane lipids: a new tool for reconstructing ancient sea water temperatures?
826        Earth and Planetary Science Letters, 204, 15 265-274, 2002.

827    Seki, O., Bendle, J. A., Harada, N., Kobayashi, M., Sawada, K., Moossen, H., Sakamoto, T.: Assessment
828        and calibration of $TEX_{86}$ paleothermometry in the Sea of Okhotsk and sub-polar North Pacific
829        region: Implications for paleoceanography. Progress in Oceanography, 126, 254–266, 2014.

Climate
of the Past
Discussions

830 Siliakus, M., van der Oost, J., Kengen, S.W.M.: Adaptations of archaeal and bacterial membranes to
831         variations in temperature, pH and pressure. Extremophiles, 21, 651 – 670, 2017.

832 Sluijs A, Schouten S, Pagani M, Woltering, M., Brinkhuis, H., Sinninghe Damsté, J.S., Dickens, G.R.,
833         Huber, M., Reichart, G., Stein, R., Matthiessen, J., Lourens, L.J., Pedentchouk, N., Backman, J.,
834         Moran, K. and the Expedition 320 Scientists: Subtropical arctic ocean temperatures during the
835         Palaeocene/Eocene thermal maximum. Nature 441, 610–613, 2006.

836 Sluijs, A., Schouten, S., Donders, T. H., Schoon, P. L., Rohl, U., Reichart, G.-J., Sangiorgi, F., Kim, J.-H.,
837         Sinninghe Damsté, J. S., Brinkhuis, H.: Warm and wet conditions in the Arctic region during
838         Eocene Thermal Maximum 2, Nature Geosci, 2, 777-780, 2009.

839 Taylor, K. W. R., Willumsen, P. S., Hollis, C. J., Pancost, R. D.: South Pacific evidence for the long-term
840         climate impact of the Cretaceous/Paleogene boundary event, Earth-Science Reviews, 179, 287-302,
841         2018.

842 Taylor, K. W., Huber, M., Hollis, C. J., Hernandez-Sanchez, M. T., Pancost, R. D.: Re-evaluating modern
843         and Palaeogene GDGT distributions: Implications for SST reconstructions, Global and Planetary
844         Change, 108, 158-174, 2013.

845 Tierney, J. E.: GDGT Thermometry: Lipid Tools for Reconstructing Paleotemperatures. Retrieved from
846         https://www.geo.arizona.edu/~jesst/resources/TierneyPSP_GDGTs.pdf, 2012

847 Tierney, J. E., and Tingley, M. P.: A Bayesian, spatially-varying calibration model for the $TEX_{86}$ proxy.
848         Geochimica et Cosmochimica Acta, 127, 83-106, 2014.

849 Tierney, J. E., and Tingley, M. P.: A $TEX_{86}$ surface sediment database and extended Bayesian calibration,
850         Scientific data, 2, 150029, 2015.

851 Williams, C.K.I., and Rasmussen, C.E.: Gaussian processes for machine learning. MIT Press Cambridge,
852         MA, 2006.

853 Wuchter, C., Schouten, S., Coolen, M. J. L., and Sinninghe Damsté, J. S.: Temperature-dependent variation
854         in the distribution 30 of tetraether membrane lipids of marine Crenarchaeota: Implications for
855         $TEX_{86}$ paleothermometry, Paleoceanography and Paleoclimatology. doi:10.1029/2004PA001041,
856         2004.

857 Zhang, Y. G., and Liu, X.: Export Depth of the $TEX_{86}$ Signal. Paleoceanography and Paleoclimatology.
858         doi.org/10.1029/2018PA003337, 2018.

859 Zhang, Y. G., Pagani, M., Wang, Z.: Ring Index: A new strategy to evaluate the integrity of $TEX_{86}$
860         paleothermometry, Paleoceanography, 31, 220-232, 2016.

861 Zhang, Y. G., Zhang, C. L., Liu, X.-L., Li, L., Hinrichs, K.-U., Noakes, J. E.: Methane Index: a tetraether
862         archaeal lipid 15 biomarker indicator for detecting the instability of marine gas hydrates, Earth and
863         Planetary Science Letters, 307, 525- 534, 2011.

864 Zhang, Y.G., Pagani, M., Liu, Z.: A 12-million-year temperature history of the tropical Pacific
865         Ocean: Science, 343, 84-86, 2014.
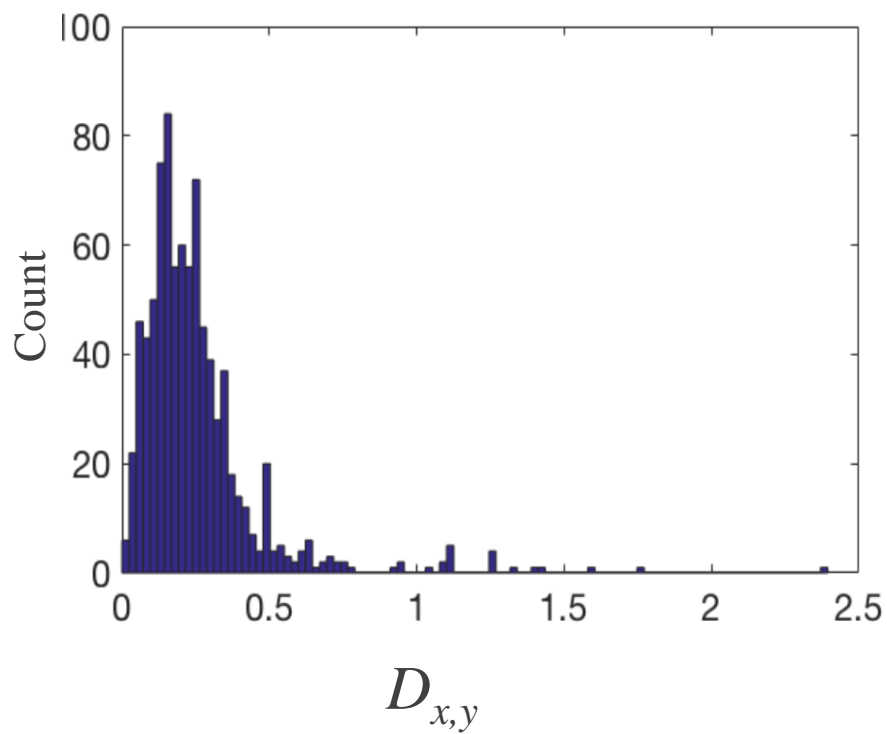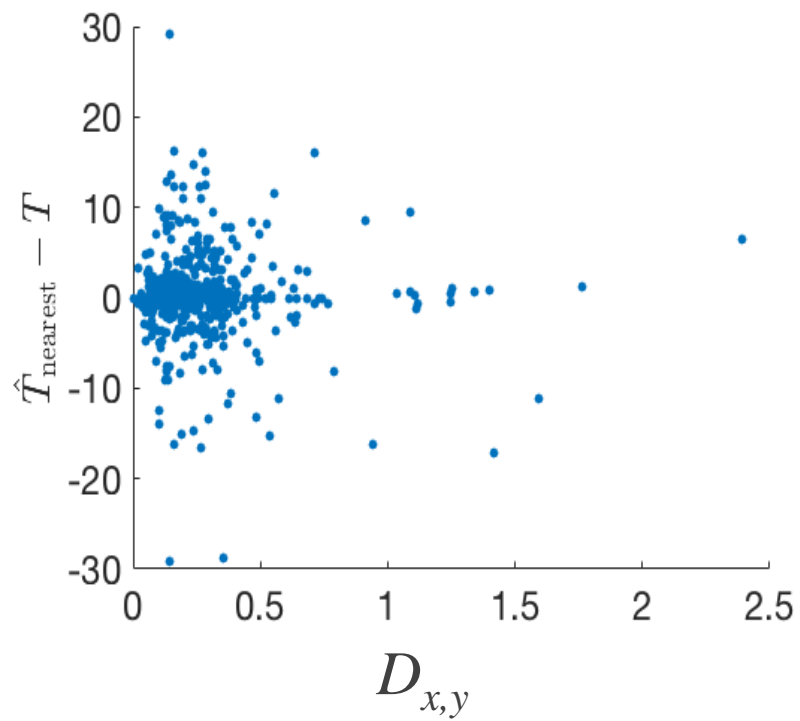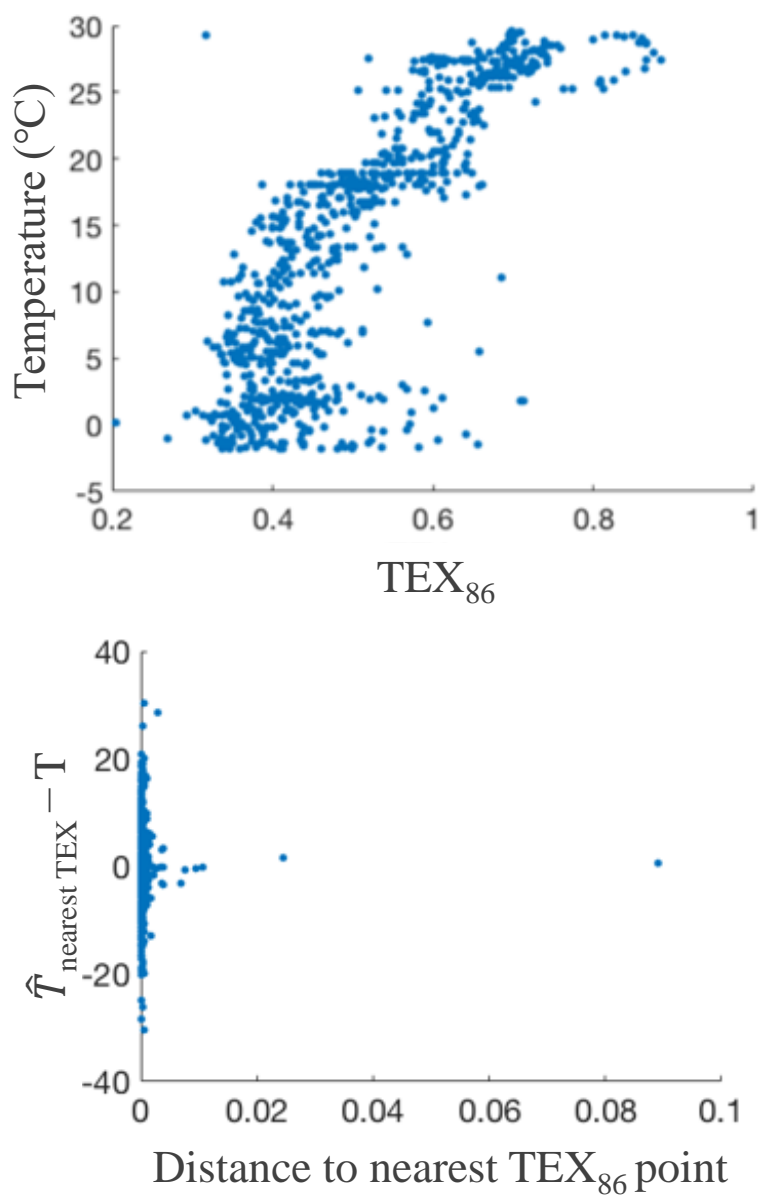
866

867

Figure 1

Climate
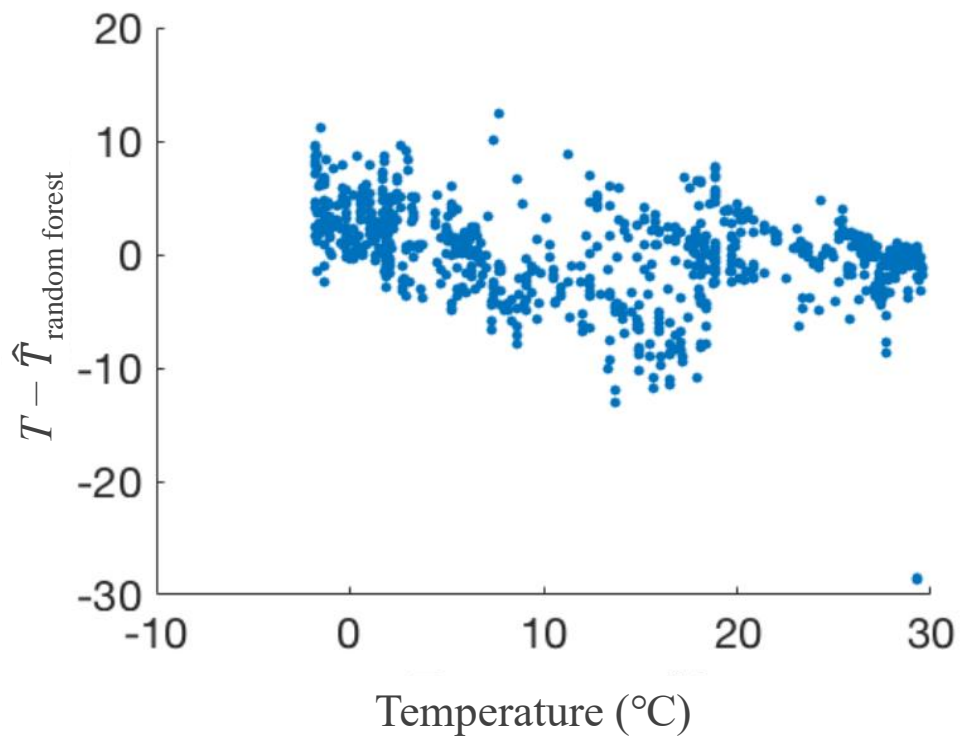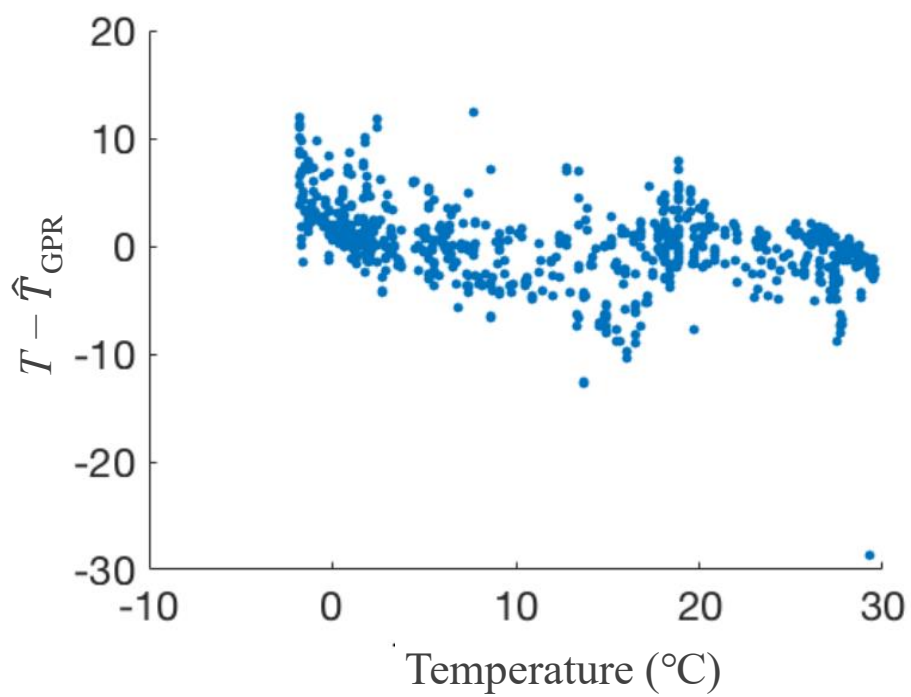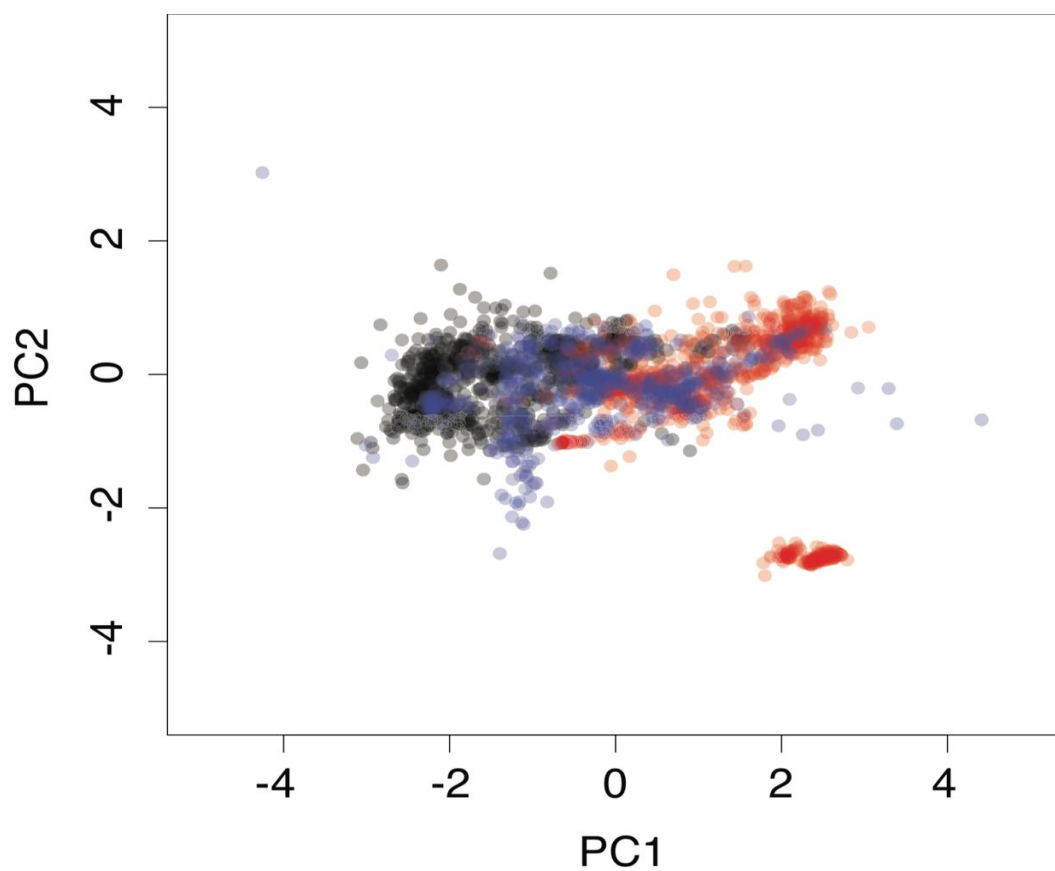of the Past
Discussions

Open Access

EGU

Figure 2

Figure 3
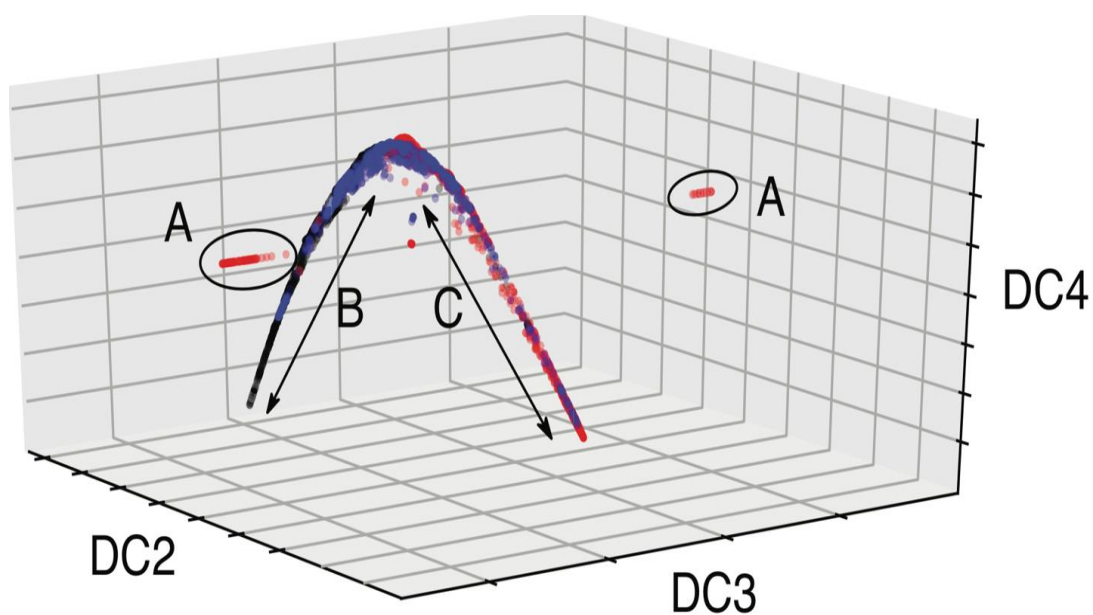
Figure 4

Figure 5

Figure 6
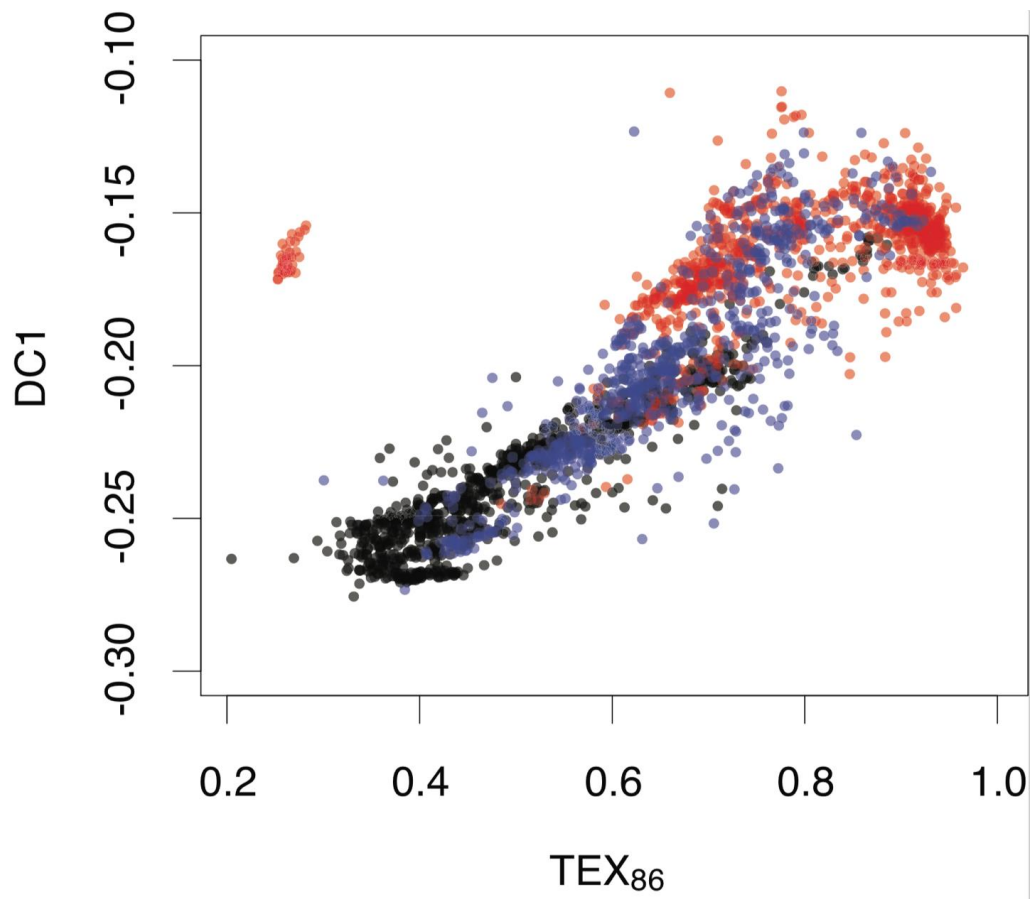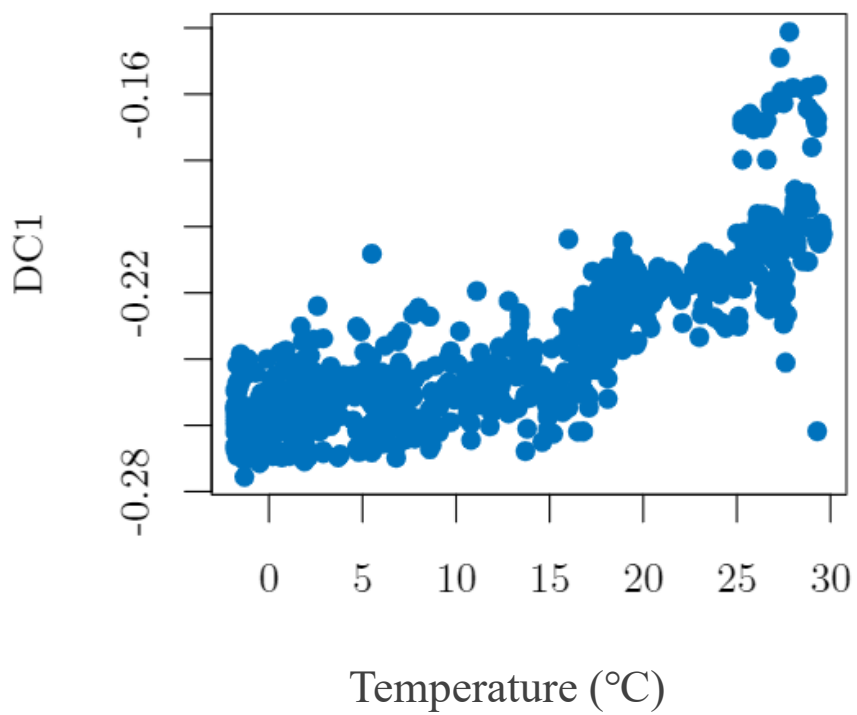
Figure 7

Figure 8

Climate
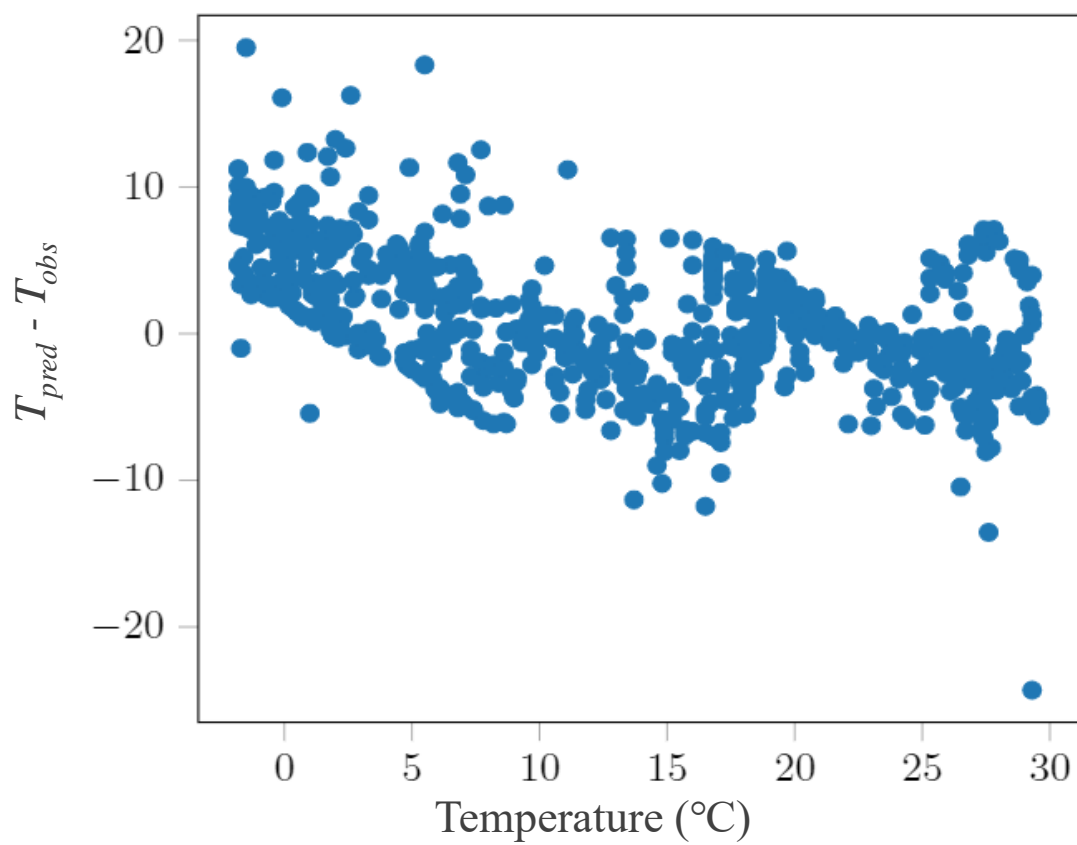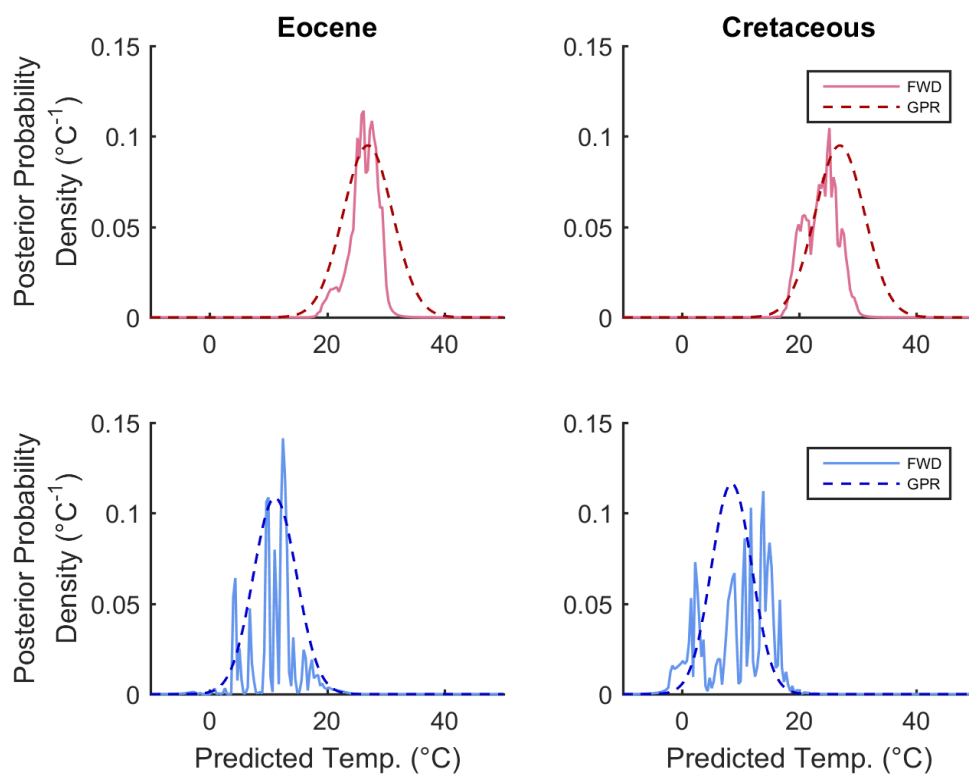of the Past

Discussions

Open Access

EGU

Figure 9

Figure 10

Figure 11

Figure 12

# Figure 13

Figure 14