

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

2

3 Tom Dunkley Jones¹, Yvette L. Eley¹, William Thomson², Sarah E. Greene¹, Ilya Mandel^{3,4,5}, Kirsty Edgar¹
4 and James A. Bendle¹

5

6 **Affiliations:**

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

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

10 ³School of Physics and Astronomy, Monash University, Clayton, Vic. 3800, Australia

11 ⁴The ARC Centre of Excellence for Gravitational Wave Discovery — OzGrav, Australia

12 ⁵Birmingham Institute for Gravitational Wave Astronomy and School of Physics and Astronomy,
13 University of Birmingham, B15 2TT, Birmingham, UK

14

15 **Abstract**

16

17 In the modern oceans, the relative abundances of Glycerol dialkyl glycerol tetraether (GDGTs) compounds
18 produced by marine archaeal communities show a significant dependence on the local sea surface
19 temperature at the site of deposition. When preserved in ancient marine sediments, the measured
20 abundances of these fossil lipid biomarkers thus have the potential to provide a geological record of long-
21 term variability in planetary surface temperatures. Several empirical calibrations have been made between
22 observed GDGT relative abundances in late Holocene core top sediments and modern upper ocean
23 temperatures. These calibrations form the basis of the widely used TEX₈₆ palaeothermometer. There are,
24 however, two outstanding problems with this approach, first the appropriate assignment of uncertainty to
25 estimates of ancient sea surface temperatures based on the relationship of the ancient GDGT assemblage to
26 the modern calibration data set; and second, the problem of making temperature estimates beyond the range
27 of the modern empirical calibrations (>30 °C). Here we apply modern machine-learning tools, including
28 Gaussian Process Emulators and forward modelling, to develop a new mathematical approach we call
29 OPTiMAL (Optimised Palaeothermometry from Tetraethers via MAchine Learning) to improve
30 temperature estimation and the representation of uncertainty based on the relationship between ancient
31 GDGT assemblage data and the structure of the modern calibration data set. We reduce the root mean
32 square uncertainty on temperature predictions (validated using the modern data set) from $\sim\pm 6$ °C using
33 TEX₈₆ based estimators to ± 3.6 °C using Gaussian Process estimators for temperatures below 30 °C. We
34 also provide a new quantitative measure of the distance between an ancient GDGT assemblage and the

35 nearest neighbour within the modern calibration dataset, as a test for significant non-analogue behaviour.
36 Finally, we advocate caution in the use of temperature estimates beyond the range of the modern empirical
37 calibration dataset, given the lack of a robust predictive biological model or extensive and reproducible
38 mesocosm experimental data in this elevated temperature range.

39

40 **1. Introduction**

41

42 Glycerol dibiphytanyl glycerol tetraethers (GDGTs) are membrane lipids consisting of isoprenoid carbon
43 skeletons ether-bound to glycerol (Schouten et al., 2013). In marine systems they are primarily produced
44 by ammonia oxidising marine Thaumarchaeota (Schouten et al., 2013). In modern marine core top
45 sediments, the relative abundance of GDGT compounds with more ring structures increases with the mean
46 annual sea surface temperature (SST) of the overlying waters (Schouten et al., 2002). This trend is most
47 likely driven by the need for increased cell membrane stability and rigidity at higher temperatures
48 (Sinninghe Damsté et al., 2002). On this basis, the TEX₈₆ (tetraether index of tetraethers containing 86
49 carbon atoms) ratio was derived to provide an index to represent the extent of cyclisation (Eq. 1; where
50 GDGT-x represents the fractional abundance of GDGT-x determined by liquid chromatography mass
51 spectrometry (LC-MS) peak area, and cren' is the peak area of the isomer of crenarchaeol) (Schouten et
52 al., 2002; Liu et al. 2018) and was shown to be positively correlated with mean annual SSTs:

53

$$54 \text{TEX}_{86} = (\text{GDGT-2} + \text{GDGT-3} + \text{cren}') / (\text{GDGT-1} + \text{GDGT-2} + \text{GDGT-3} + \text{cren}') \text{ (Eq. 1)}$$

55

56 Early applications of TEX₈₆ to reconstruct ancient SSTs were promising, especially in providing
57 temperature estimates in environments where standard carbonate-based proxies are hampered by poor
58 preservation (Schouten et al., 2003; Herfort et al., 2006; Schouten et al., 2007; Huguet et al., 2006; Sluijs
59 et al., 2006; Brinkhuis et al., 2006; Pearson et al., 2007; Sluijs et al., 2009). The TEX₈₆ approach also
60 extended beyond the range of the widely used alkenone-based U^{k'}₃₇ thermometer, in both temperature space,
61 where U^{k'}₃₇ saturates at ~28°C (Brassell, 2014; Zhang et al., 2017), and back into the early Cenozoic (Bijl
62 et al., 2009; Hollis et al., 2009; Bijl et al., 2013; Inglis et al., 2015) and Mesozoic (Schouten et al., 2002;
63 Jenkyns et al., 2012; O'Brien et al., 2017) where haptophyte-derived alkenones are typically absent from
64 marine sediments (Brassell, 2014). Initially, TEX₈₆ was converted to SSTs using the core-top calibration
65 (Schouten et al. 2002) (Eq. 2):

66

$$67 \text{TEX}_{86} = 0.015 * \text{SST} + 0.287 \text{ (Eq. 2)}$$

68

69 However as the number and range of applications of TEX_{86} palaeothermometry grew, concerns arose about
 70 proxy behaviour at both the high (Liu et al., 2009) and low (Kim et al., 2008) temperature ends of the
 71 modern calibration. In response to these observations, a new expanded modern core top dataset (Kim et al.,
 72 2010) was used to generate two new indices – TEX_{86}^L (Eq. 3), an exponential function that does not include
 73 the crenarchaeol regio-isomer and was recommended for use across the entire temperature range of the new
 74 core top data (-3 to 30 °C, particularly when SSTs are lower than 15 °C), and TEX_{86}^H (Eq. 4), also
 75 exponential, and recommended for use when SSTs exceeded 15 °C (Kim et al., 2010). TEX_{86}^L also excludes
 76 GDGT abundance data from the high-temperature regimes of the Red Sea, which are somewhat anomalous
 77 and likely related to salinity effects on community composition in this region (Trommer et al., 2009, Kim
 78 et al. 2010).

79

$$80 \quad TEX_{86}^L = \log \left(\frac{[GDGT2]}{[GDGT1]+[GDGT2]+[GDGT3]} \right) \quad \text{Eq. 3}$$

81

82

$$83 \quad TEX_{86}^H = \log \left(\frac{[GDGT2]+[GDGT3]+[Cren']}{[GDGT1]+[GDGT2]+[GDGT3]+[Cren']} \right) \quad \text{Eq. 4}$$

84

85 Despite the recommendations of Kim et al. (2010), both TEX_{86}^H and TEX_{86}^L were widely used and tested
 86 across a range of temperatures and palaeoenvironments, including comparisons against other
 87 palaeotemperature proxy systems (Hollis et al. 2012; Lunt 2012 Dunkley Jones et al. 2013; Zhang et al.,
 88 2014; Seki et al., 2014; Douglas et al., 2014; Linnert et al., 2014; Hertzberg et al., 2016). The rationale was
 89 that both TEX_{86}^L and TEX_{86}^H were calibrated across a full temperature range, with the exception of the
 90 inclusion or exclusion of Red Sea core-top data. The difference in model fit between the two proxy
 91 formulations to the calibration dataset was also minor (Kim et al. 2010). In certain environments, however,
 92 TEX_{86}^L was subject to significant variability in derived temperatures that were not apparent in TEX_{86}^H
 93 (Taylor et al., 2013). This was mostly due to changing GDGT2 to GDGT3 ratios, which strongly influence
 94 TEX_{86}^L , and may be related to local non-thermal environmental conditions at the site of GDGT production,
 95 and deep-water lipid production, (Taylor et al., 2013). As a result, TEX_{86}^L is no longer regarded as an
 96 appropriate tool for palaeotemperature reconstructions, except in limited Polar conditions (Kim et al., 2010;
 97 Tierney, 2012).

98

99 Three fundamental issues have troubled the TEX_{86} proxy. The first is a concern about undetected non-
 100 analogue palaeo-GDGT assemblages, for which the modern calibration data set is inadequate to provide a
 101 robust temperature estimation. Although various screening protocols, with independent indices and

102 thresholds, have been proposed to test for an excessive influence of terrestrial lipids (Branched and
103 Isoprenoid Tetraether, BIT index; Hopmans et al., 2004), within sediment methanogenesis (Methane Index,
104 ‘MI’; Zhang et al., 2011) and non-thermal effects such as nutrient levels and archaeal community structure
105 to impact the weighted average of cyclopentane moieties (Ring Index, ‘RI,’ Zhang et al., 2016), these do
106 not provide a fundamental measure of the proximity between GDGT abundance distributions in the modern,
107 and ancient GDGT abundance distributions recorded in sediment samples. The fundamental question
108 remains – are measured ancient assemblages of GDGT compounds anything like the modern assemblages,
109 from which palaeotemperatures are being estimated? Understanding this question cannot easily be
110 addressed with the use of indices – TEX₈₆ itself, or BIT and MI – that collapse the dimensionality of GDGT
111 abundance relationships onto a single axis of variation.

112
113 Second, from the earliest applications of the TEX₈₆ 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 °C), 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 °C. However, this problem of model choice, and its impact on temperature
118 estimation beyond the modern calibration range, persists (Hollis et al. 2019), 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 TEX₈₆ on SSTs, and the effect of these models on temperature estimates over 30 °C.

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., 2019). While there is a
124 basic underlying trend for more rings within GDGT structures at higher temperatures (Zhang et al. 2015;
125 Qin et al., 2015), the lack of a uniform response to archaeal GDGT production in response to increasing
126 growth temperatures (e.g., Elling et al., 2015; Qin et al., 2015) suggests that this does not easily translate
127 into a simple linear model at the community scale (i.e. the core top calibration dataset). Wuchter et al.
128 (2004) and Schouten et al. (2007) show a compiled linear calibration of TEX₈₆ against incubation
129 temperature (up to 40°C in the case of Schouten et al., 2007) based on strains that were enriched from
130 surface seawater collected from the North Sea and Indian Ocean respectively. Like Qin et al. (2015), we
131 note the *non*-linear nature of the individual experiments in Wuchter et al. (see Fig. 5 in Wuchter et al. 2004).
132 Moreover, the relatively lower Cren’ in these studies yield a very different intercept and slope compared to
133 core-top calibrations (e.g. Kim et al. 2010) making direct comparisons problematic.

134

135 More recently, Elling et al. (2015) studied three different strains (*N. maritimus*, NAOA6, NAOA2) isolated
136 from open ocean surface waters (South Atlantic) whilst Qin et al., (2015) studied a culture of *N. maritimus*
137 and three *N. maritimus*-like strains isolated from Puget Sound. All strains are of marine, mesophilic,
138 Thaumarchaeota within Marine Group 1 (equivalent to Crenarchaeota Group 1). Both of these papers
139 clearly demonstrate distinctly different responses of membrane lipid composition to temperature in these
140 strains, whilst Qin et al. (2015) additionally show that oxygen concentration is at least as important as
141 temperature in controlling TEX₈₆ values in culture. The impact of Thaumarchaeota community change on
142 TEX₈₆ in palaeoclimate studies is further suggested by the downcore study of Polik et al (2019). All of these
143 culture studies, made on marine, mesophilic archaea demonstrate how community composition may have
144 a significant impact on measured environmental TEX₈₆ signatures.

145

146 It is clear from the above discussion that there is evidence for more complex responses in GDGT-production
147 to growth temperature in some instances, and across distinct strains of archaea (Elling et al., 2015). More
148 fundamentally, in natural systems, it is likely that aggregated GDGT abundance variations in response to
149 growth temperatures result from changing compositions of archaeal populations as well as the physiological
150 response of individual strains to growth temperature (Elling et al. 2015). For instance, a multiproxy study
151 of Mediterranean Pliocene-Pleistocene sapropels indicates that specific distributions of archaeal lipids
152 might be reflective of temporal changes in thaumarchaeal communities rather than temperature alone
153 (Polik et al., 2018). Indeed, the potential influence of community switching on GDGT composition can be
154 seen in mesocosm studies, with different species preferentially thriving at different growth temperatures
155 (e.g., Schouten et al., 2007). To use the responses of single, selected archaeal strains in culture to validate
156 a particular model of community-level responses to growth temperature is problematic even in the modern
157 system (Elling et al., 2015). For deep time applications it is even more difficult, where there is no
158 independent constraint on the archaeal strains dominating production or their evolution through time (Elling
159 et al. 2015). What is notable, however, is that the Ring Index (RI) - calculated using all commonly measured
160 GDGTs (Zhang et al., 2016) – has a more robust relationship with culture temperature between archaeal
161 strains than TEX₈₆, indicating a potential loss of information within the TEX₈₆ index (Elling et al. 2015).

162

163 Finally, the original uses of the TEX₈₆ proxy had a relatively poor representation of the true uncertainty
164 associated with palaeotemperature estimates, as they included no assessment of non-analogue behavior
165 relative to the modern core-top data. Instead, uncertainty was typically based on the residuals on the modern
166 calibration, with no reference to the relationship between GDGT distributions of an ancient sample and the
167 modern calibration data. An improved Bayesian uncertainty model “BAYSPAR” is now in widespread use
168 for SST estimation, which models TEX₈₆ to SSTs regression parameters, and associated uncertainty, as

169 spatially varying functions (Tierney and Tingley, 2015). The Bayesian approach, as with all approaches
170 based on the TEX₈₆ index, however, still does not include an uncertainty that reflects how well modelled
171 ancient GDGT assemblages are by the modern calibration – i.e. the degree to which they are non-analogue
172 - as it still functions on one-dimensional TEX₈₆ index values.

173
174 All empirical calibrations of GDGT-based proxies assume that mean annual SST is the master variable on
175 GDGT assemblages both today and in the past. Mean annual SST, however, is strongly correlated with
176 many other environmental variables (e.g., seasonality, pH, mixed layer depth, and productivity). In the
177 modern calibration dataset, mean annual SST shows the strongest correlation with TEX₈₆ index (Schouten
178 et al., 2002), but this does not preclude an important (but undetectable) influence of these other
179 environmental variables. The use of empirical GDGT calibrations to infer ancient sea surface temperatures
180 thus implicitly assumes that the relationships between mean annual SST and all other GDGT-influencing
181 variables are invariant through time. This assumption is inescapable until, and unless, a more complete
182 biological mechanistic model of GDGT production emerges.

183
184 Here, we return to the primary modern core-top GDGT assemblage data (Tierney and Tingley, 2015), and
185 systematically explore the relationships between the modern GDGT distributions and surface ocean
186 temperatures using powerful mathematical tools. These tools can investigate correlations without prior
187 assumptions on the best form of relationship or *a priori* selection of GDGT compounds to be used. This
188 analysis is then extended through the exploration of the relationships between the modern core top GDGT
189 distributions and two compilations of ancient GDGT datasets, one from the Eocene (Inglis et al. 2015) and
190 one from the Cretaceous (O’Brien et al. 2017). We explore simple metrics to answer the fundamental
191 question – are modern core-top GDGT distributions good analogues for ancient distributions? We propose
192 the first robust methodology to answer this question, and so screen for significantly non-analogue palaeo-
193 assemblages. From this, we go on to derive a new machine learning approach ‘OPTiMAL’ (Optimised
194 Palaeothermometry from Tetraethers via MAchine Learning) for reconstructing SSTs from GDGT
195 datasets, which outperforms previous GDGT palaeothermometers and includes robust error estimates that,
196 for the first time, accounts for model uncertainty.

197

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

199
200 Our new analyses use the modern core-top data compilation, and satellite-derived estimates of SSTs, of
201 Tierney and Tingley (2015) as well as compilations of Eocene (Inglis et al. 2015) and Cretaceous (O’Brien
202 et al. 2017) GDGT assemblages. Within these fossil assemblages, only data points with full characterisation

203 of individual GDGT relative abundances were used. We also note that, in the first instance, all available
204 fossil assemblage data were included, although later comparisons between BAYSPAR and our new
205 temperature predictor excludes fossil data that was regarded as unreliable based on standard pre-screening
206 indices, as noted within the original compilations (Inglis et al. 2015; O'Brien et al. 2017). All data used in
207 this study are tabulated in the supplementary information.

208

209 In order to enable meaningful comparison between new and existing temperature predictors, we use the
210 following consistent procedure for evaluating all predictors throughout this paper. We divide the modern
211 core-top data set of 854 data points into 85 validation data points (chosen randomly) and 769 calibration
212 points (as we require fractional abundances for all 6 commonly measured GDGTs, we excluded those data
213 points for which these values were not reported). We calibrate the predictor on the calibration points, and
214 then judge its performance on the validation points using the root mean square error:

215

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

(Eq. 5)

217

218

219 where the sum is taken over each of $N_v = 85$ validation points, T is the known measured temperature (which
220 we refer to as the true temperature) and \hat{T} is the predicted temperature. For conciseness, we refer to δT as
221 the predictor standard error. It is useful to compare the accuracy of the predictor to the standard deviation
222 of all temperatures in the data set σT , which corresponds to using the mean temperature as the predictor in
223 Equation 1; for the modern data set, $\sigma T = 10.0$ °C. The coefficient of determination, R^2 , provides a measure
224 of the fraction of the fluctuation in the temperature explained by the predictor. To facilitate performance
225 comparisons between different methods of predicting temperature, we use the same subset of validation
226 points for all analyses. To avoid sensitivity to the choice of validation points, we repeat the calibration-
227 validation procedure for 10 random choices from the validation dataset.

228

229 *2.1 Nearest neighbours*

230

231 We begin with an agnostic approach to using some combination of the proportions of each of the six
232 observables - GDGT-0, GDGT-1, GDGT-2, GDGT-3, crenarchaeol and cren', which we will jointly refer
233 to as GDGTs - to predict sea surface temperatures. Whatever functional form the predictor might take, it
234 can only provide accurate temperature predictions if nearby points in the six-dimensional observable space

235 - i.e. the distribution of all of the six commonly reported GDGTs - can be translated to nearby points in
 236 temperature space. Conversely, if nearby points in the observable space correspond to vastly different
 237 temperatures, then no predictor, regardless of which combination of GDGTs are used, will be able to
 238 provide a useful temperature estimate. In other words, the structuring of GDGT distributions within multi-
 239 dimensional space, must have some correspondence to the temperatures of formation (or rather the mean
 240 annual SSTs used for standard calibrations).

241
 242 We therefore consider the prediction offered by the temperature at the nearest point in the GDGT parameter
 243 space. Of course, nearness depends on the choice of the distance metric. For example, it may be that sea
 244 surface temperatures are very sensitive to a particular GDGT, so even a small change in that GDGT
 245 corresponds to a significant distance, and rather insensitive to another, meaning that even with a large
 246 difference in the nominal value of that GDGT the distance is insignificant. In the first instance, we use a
 247 very simple Euclidian distance estimate $D_{x,y}$ where the distance along each GDGT is normalised by the total
 248 spread in that GDGT across the entire data set. This normalisation ensures that a dimensionless distance
 249 estimate can be produced even when observables have very different dynamical ranges, or even different
 250 units. Thus, the normalised distance D between parameter data points x and y is

$$D_{x,y}^2 \equiv \sum_{i=0}^5 \frac{(GDGT_i(x) - GDGT_i(y))^2}{var(GDGT_i)} \quad (\text{Eq. 7})$$

251
 252
 253
 254 We show the distribution of nearest distances of points in the modern data set, excluding the sample itself,
 255 in (Fig. 1).

256
 257 The nearest-sample temperature predictor is $\hat{T}_{\text{nearest}}(x) = T(y)$ where y is the nearest point to x over the
 258 calibration data set, i.e., one that minimises $D_{x,y}$. Fig. 2 shows the scatter in the predicted temperature when
 259 using the temperature of the nearest data point to make the prediction. Overall, the failure of the nearest-
 260 neighbour predictor to provide accurate temperature estimates even when the normalised distance to the
 261 nearest point is small, $D_{x,y} \leq 0.5$, casts doubt on the possibility of designing an accurate predictor for
 262 temperature based on GDGT observations. This is most likely due to additional environmental controls on
 263 GDGT abundance distributions in natural systems, in particular the water depth (Zhang and Liu, 2018),
 264 nutrient availability (Hurley et al., 2016; Polik et al., 2018; Park et al., 2018), seasonality, growth rate
 265 (Elling et al., 2014; Hurley et al., 2016) and ecosystem composition (Polik et al., 2018), that obscure a
 266 predominant relationship to mean annual SSTs.

267
 268

269 On the other hand, the standard error for the nearest-neighbour temperature predictor is $\delta T_{\text{nearest}} = 4.5 \text{ }^\circ\text{C}$.
270 This is less than half of the standard deviation σT in the temperature values across the modern data set.
271 Thus, the temperatures corresponding to nearby points in GDGT observable space also cluster in
272 temperature space. Consequently, there is hope that we can make some useful, if imperfect, temperature
273 predictions. The value of $\delta T_{\text{nearest}}$ will also serve as a useful benchmark in this design: while we may hope
274 to do better by, say, suitably averaging over multiple nearby calibration points rather than adopting the
275 temperature at one nearest point as a predictor, any method that performs worse than the nearest-neighbour
276 predictor is clearly suboptimal.

277

278 *2.2 TEX₈₆ and Bayesian applications*

279

280 The TEX₈₆ index reduces the six-dimensional observable GDGT space to a single number. While this has
281 the advantage of convenience for manipulation and the derivation of simple analytic formulae for
282 predictors, as illustrated below, this approach has one critical disadvantage: it wastes significant information
283 embedded in the hard-earned GDGT distribution data. Fig. 3 illustrates both the advantage and
284 disadvantage of TEX₈₆. On the one hand, there is a clear correlation between TEX₈₆ and temperature (top
285 panel of Fig. 3), with a correlation coefficient of 0.81 corresponding to an overwhelming statistical
286 significance of 10^{-198} . On the other hand, very similar TEX₈₆ values can correspond to very different
287 temperatures. We can apply the nearest-neighbour temperature prediction approach to the TEX₈₆ value
288 alone rather than the full GDGT parameter space; this predictor yields a large standard error of $\delta T_{\text{nearestTEX86}}$
289 $= 8.0 \text{ }^\circ\text{C}$ (bottom panel of Fig. 3). While smaller than σT , this is significantly larger than $\delta T_{\text{nearest}}$ (Fig. 2),
290 consistent with the loss of information in TEX₈₆. We therefore do not expect other predictors based on
291 TEX₈₆ to perform as well as those based on the full available data set.

292

293 Indeed, this is what we find when we consider predictors of the form $\hat{T}_{1/\text{TEX}} = a + b/\text{TEX}_{86}$ and $\hat{T}_{\text{TEXH}} = c$
294 $+ d \log_{\text{TEX}_{86}}$ (Liu et al., 2009; Kim et al., 2010), i.e., the established relationships between GDGT
295 distributions and SST. We fit the free parameters a , b , c , and d by minimising the sum of squares of the
296 residuals over the calibration data sets (least squares regression). We find that $\delta T_{1/\text{TEX}} = 6.1 \text{ }^\circ\text{C}$ (note that
297 this is slightly better than using the fixed values of a and b from (Kim et al., 2010), which yield $\delta T_{1/\text{TEX}} =$
298 $6.2 \text{ }^\circ\text{C}$). We note that the corresponding R^2 value associated with these TEX₈₆ based predictors is 0.64,
299 which is lower than the R^2 values in Kim et al. (2010). We attribute this to the fact that we are using a larger
300 dataset based on Tierney and Tingley (2015), including data from the Red Sea (Kim et al. 2010).

301

302 Tierney and Tingley (2014) proposed a more sophisticated approach to obtaining the transfer function from
303 TEX₈₆ to temperature, continuing to use simple linear regression, but with the addition of Gaussian
304 processes to model spatial variability in the temperature-TEX₈₆ relationship and working with a forward
305 model which is subsequently inverted to produce temperature predictions. This forward model
306 ‘BAYSPAR’ is capable of generating an infinite number of calibration curves relating TEX₈₆ to sea surface
307 temperatures (Tierney and Tingley, 2014). In order to derive a calibration for a specific dataset, the user
308 edits a range of parameters which vary depending on whether the dataset in question is from the relatively
309 recent past or deep time (Tierney and Tingley, 2014). For deep time applications, the authors propose a
310 modern analogue-type approach, in which they search the modern data for 20° x 20° grid boxes containing
311 ‘nearby’ TEX₈₆ measurements and subsequently apply linear regression models calibrated on the analogous
312 samples for making predictions.

313
314 However, along with the simpler TEX₈₆-based models described above, this approach still suffers from the
315 reduction of a six-dimensional data set to a single number. Therefore, it is not surprising that even the
316 simplest nearest-neighbour predictor (such as the one described above) that makes use of the full six-
317 dimensional dataset outperforms single-dimensional forward modelling approaches. Additionally,
318 uncertainty estimates do not account for the fact that TEX₈₆ is, fundamentally, an empirical proxy, and so
319 its validity outside the range of the modern calibration is not guaranteed. This is a fundamental issue for
320 attempts to reconstruct surface temperatures during Greenhouse climate states, when tropical and sub-
321 tropical SSTs were likely hotter than those observed in the modern oceans.

322 323 *2.3 Machine learning Approaches – Random Forests*

324
325 There are a number of options to improve on nearest-neighbour predictions using machine learning
326 techniques such as artificial neural networks and random forests. These flexible, non-parametric models
327 would ideally be based on the underlying processes driving the GDGT response to temperature, but since
328 these processes remain unconstrained at present, we choose to deploy models which can reasonably reflect
329 predictive uncertainty and will be sufficiently adaptable in future (as new information regarding controls
330 on GDGTs emerge). These machine learning approaches are all based on the idea of training a predictor by
331 fitting a set of coefficients in a sufficiently complex multi-layer model in order to minimise residuals on
332 the calibration data set. As an example of the power of this approach, we train a random forest of decision
333 trees with 100 learning cycles using a least-squares boosting to fit the regression ensemble. Figure 4 shows
334 the prediction accuracy for this random forest implementation. This machine learning predictor yields δT
335 = 4.1 °C degrees, outperforming the naive nearest-neighbour predictor by effectively applying a suitable

336 weighted average over multiple near neighbours. This corresponds to a very respectable $R^2 = 0.83$, meaning
337 that 83% of the variation in the observed temperature is successfully explained by our GDGT-based model.

338

339 *2.4 Gaussian Process Regression*

340

341 One downside of the random forest predictor is the difficulty of accurately estimating the uncertainty on
342 the prediction (Mentch and Hooker, 2016), although this is possible with, e.g., a bootstrapping approach
343 (Coulston et al., 2016). Fortunately, Gaussian process (GP) regression provides a robust alternative. For
344 full details on GP regression refer to Williams and Rasmussen (2006) and Rasmussen and Nickisch (2010).
345 Loosely, the objective here is to search among a large space of smoothly varying functions of GDGT
346 compositions for those functions which adequately describe temperature variability. This, essentially, is a
347 way of combining information from all calibration data points, not just the nearest neighbours, assigning
348 different weights to different calibration points depending on their utility in predicting the temperature at
349 the input of interest. The trained Gaussian process learns the best choice of weights to fit the data. Typically,
350 the GP will give greater weight to closer points, but, as we discuss below, it will learn the appropriate
351 distance metric on the multi-dimensional GDGT input space.

352

353 The weighting coefficients learned by the GP emulator represent a covariance matrix on the GDGT
354 parameter space. We can use this as a distance metric to provide meaningfully normalised distances
355 between points, removing the arbitrariness from the nearest neighbour distance ($D_{x,y}$) definition used earlier,
356 and this is the basis of the D_{nearest} metric described below. If the temperature is insensitive to a particular
357 GDGT input coordinate (i.e., the value of that input has a minimal effect on the temperature) then points
358 within GDGT space that have large differences in absolute input values in that coordinate are still near. We
359 find that Cren has very limited predictive power, and so points with large Cren differences are close in term
360 of the normalised distance. Conversely, if the temperature is sensitive to small changes in a particular
361 GDGT variant, then points with relatively nearby absolute input values in that coordinate are still distant.
362 We find that most GDGT parameters other than Cren are comparably useful in predicting temperature, with
363 GDGT-0 and GDGT-3 marginally the most informative. We considered whether interdependency of
364 percentage GDGT data could influence our calculations. Our analysis suggests that there are only five free
365 parameters. Machine learning tools should be able to pick up this correlation and effectively ignore one of
366 the parameters (or one parameter combination). For example, we do find that the GP emulator has a very
367 broad kernel in at least one dimension, signaling this. In principle, we could have considered only five of
368 six parameters. The smaller scale of some of the parameters is automatically accounted for by the trained
369 kernel size in GP regression, or by normalising to the appropriate dynamical range in our initial

370 investigation. In short, the accuracy of Gaussian process regression is not adversely affected by correlations
371 between inputs (Rasmussen & Williams, 2006). Significantly correlated inputs that do not bring in new
372 predictive power are appropriately down-weighted.

373

374 We use a Gaussian process model with a squared exponential kernel with automatic relevance
375 determination (ARD) to allow for a separate length scale for each GDGT predictor. We fit the GP
376 parameters with an optimiser based on quasi-Newton approximation to the Hessian. Prediction accuracy is
377 shown in Figure 5, and we find that $\delta T = 3.72$ °C, which is a substantial improvement over the existing
378 indices, at least on the modern data. As mentioned, the GP framework provides a natural quantification of
379 predictive uncertainty, which includes uncertainty about the learned function. This is in contrast to, for
380 example, the TEX₈₆ proxy, whereby the uncertainty associated with the selection of the particular functional
381 form used for predictions is ignored. While Tierney & Tingley (2014) also use Gaussian processes to model
382 uncertainty, they model spatial variability in the TEX₈₆-temperature relationship with a Gaussian process
383 prior. While this is a valuable approach to understand regional effects in the TEX₈₆-temperature
384 relationship, it does not deal with the 'non-analogue' situations we are concerned with in this paper.

385

386 *2.5 Data Structure*

387

388 The random forest (Section 2.3) and GPR approaches (Section 2.4) are agnostic about any underlying bio-
389 physical model that might impart the observed temperature-dependence on GDGT relative abundances
390 produced by archaea. They are essentially optimized interpolation tools for mapping correlations between
391 temperature and GDGT abundances within the range of the modern calibration data set; they can make no
392 sensible inference about the behavior of this relationship outside of the range of this training data. To move
393 from interpolation within, to extrapolation beyond, the modern calibration requires an understanding of,
394 and model for, the temperature-dependence of GDGT production. To explore these relationships and the
395 extent to which the ancient and modern data reside in a coherent relationship within GDGT space, we
396 employed two forms of dimensionality reduction to enable visualisation of the data in two or three
397 dimensions. The fundamental point is that if temperature is the dominant control, all of the data should lie
398 approximately on a one-dimensional curve in GDGT space, and the arclength along this curve should
399 correspond to temperature; we will revisit this point below.

400

401 We first employed a version of principal component analysis (PCA) tailored to compositional data
402 (Aitchison, 1982, 1983; Aitchison and Greenacre, 2002; Filzmoser et al., 2009a; Filzmoser et al., 2009b;
403 Filzmoser et al., 2012). Taking into account the compositional nature of the data is important because the

404 sum-to-one constraint induces correlations between variables which are not accounted for by classical PCA.
405 Furthermore, apparently nonlinear structure in Euclidean space often corresponds to linearity in the simplex
406 (i.e. the restricted space in which all elements sum to one) (Egozcue et al., 2003). Figure 6 shows the
407 modern, Eocene and Cretaceous data projected onto the first two principal components. Aside from the
408 obvious outlying cluster of Cretaceous data, characterised by GDGT-3 fractions above 0.6, the bulk of the
409 data occupy a two-dimensional point cloud with a small amount of curvature. The large majority of the
410 Cretaceous data has more positive PC1 values relative to the modern data.

411
412 We also explored the data using diffusion maps (Coifman et al., 2005; Haghverdi et al., 2015), a nonlinear
413 dimensionality reduction tool designed to extract the dominant modes of variability in the data. Such
414 diffusion maps have been successfully used to infer latent variables that can explain patterns of gene
415 expression. In the case of biological organisms, this latent variable is commonly developmental age (called
416 pseudo-time) (Haghverdi et al., 2016). In our case, the assumption would be that this latent variable
417 corresponds to temperature. Inspection of the eigenvalues of the diffusion map transition matrix suggests
418 that four diffusion components are adequate to represent the data; we plot the second, third and fourth of
419 these components in Figure 7 for the modern and ancient data. The separate clusters marked 'A' are the
420 outlying Cretaceous points with high GDGT-3 values. The bulk of the modern data lies on the branch
421 marked 'B', while the bulk of the Cretaceous data lies on the branch marked 'C'. Notably, the majority of
422 the modern points lying on branch C are from the Red Sea, which suggests that the Red Sea data is essential
423 for understanding ancient climates (particularly Cretaceous climates).

424
425 The relationship between the first diffusion component and TEX_{86} for all data is shown in Figure 8. There
426 is a clear correlation, despite the presence of some outlying Cretaceous points, some of which are not shown
427 because they lie so far outside the majority data range within this projection. This suggests that TEX_{86} is,
428 in one sense, a natural one-dimensional representation of the data. We also plot the first diffusion
429 component for the modern data as a function of temperature (Figure 9). We see a similar pattern emerging
430 to that displayed by TEX_{86} - there is little sensitivity to temperature below 15 °C, and between ~20 and 25
431 °C. An interesting avenue for future research might be to explore the temperature-GDGT system from a
432 dynamical systems perspective, i.e. use simple mechanistic mathematical models to explore the
433 temperature-dependence of steady-state GDGT distributions. It may be that such models suggest that only
434 a few steady-states exist, and that temperature is a bifurcation parameter, i.e. it controls the switch between
435 the steady states. Note also the downward slope in the residual pattern in Figure 4 between 0 and 15-17
436 degrees celsius, and again at higher temperatures. This pattern is consistent with predictions that are biased
437 towards the centre of each 'cluster', i.e. a system which is not very sensitive to temperature, but can

438 distinguish between high and low temperatures reasonably well. This observation also links to recent culture
439 studies (Elling et al., 2015) and Pliocene-Pleistocene sapropel data (Polik et al., 2018), which support the
440 existence of discrete populations with unique GDGT-temperature relationships and that temporal changes
441 in population over time can drive changes in TEX_{86} .

442

443 *2.6 Forward Modelling*

444

445 Based on the analysis of the combined modern and ancient data structure outlined above, there appears to
446 be some consistency to underlying trends in the overall variance of GDGT relative abundances. These
447 trends provide some hope that models of this variance, and its relationship to sea surface temperature, within
448 the modern dataset could be developed to predict ancient SSTs. TEX_{86} and BAYSPAR are such models,
449 but they are limited by, first, the reduction of six-dimensional GDGT space to a one-dimensional index;
450 and second, by an *ad hoc* model choice – linear, exponential – that does not account for uncertainty in
451 model fit to the modern calibration data, and the resultant uncertainty in the estimation of ancient SSTs
452 relating to model choice. To overcome these issues, we develop a forward model based on a multi-output
453 Gaussian Process (Alvarez et al., 2012), which models GDGT compositions as functions of temperature,
454 accounting for correlations between GDGT measurements. This model is then inverted to obtain
455 temperatures which are compatible with a measured GDGT composition. In simple terms, we posit that a
456 measured GDGT composition is generated by some unknown function of temperature and corrupted by
457 noise, which may be due to measurement error or some unmodelled particularity of the environment in
458 which the sample was generated. We proceed by defining a large (in this case infinite) set of functions of
459 temperature to explore and compare them to the available data, throwing away those functions which do
460 not adequately fit the data. This means, of course, that the behaviour of the functions we accept is allowed
461 to vary more widely outside the range of the modern data than within it. With no mechanistic underpinning,
462 choosing only one function (such as the inverse of TEX_{86}) based on how well it fits the modern data grossly
463 underestimates our uncertainty about temperature where no modern analogue is available.

464

465 The forward modelling approach is similar to that of Haslett et al. (2006), who argue that it is preferable to
466 model measured compositions as functions of climate, before probabilistically inverting the model to infer
467 plausible climates given a composition. The cost of modelling the data in this more natural way is the loss
468 of degrees of freedom -- we are now attempting to fit a one-dimensional line through a multidimensional
469 point cloud rather than fit a multidimensional surface to the GDGT data, which means that the predictive
470 power of the model suffers, at least on the modern data. The existing BAYSPAR calibration also specifies
471 the model in the forward direction, however while BAYSPAR does model spatial variability it assumes a

472 monotonic relationship between TEX and SST, only accounting for uncertainties on the parameters within
473 the model, rather than any systematic uncertainty in the model itself. As with all GP models, the choice of
474 kernel has a substantial impact on predictions (and their associated uncertainty) outside the range of the
475 modern data, where predictions revert to the prior implied by the kernel. Given that we have no mechanistic
476 model for the data generating process, we recommend the use of kernels which do not impose strong prior
477 assumptions on the form of the GDGT-temperature relationship (e.g. kernels with a linear component) and
478 thus reasonably represent model uncertainty outside the range of the modern data. We choose a zero-mean
479 Matern 3/2 kernel for the applications below. Note, however, that since we are working in ilr-transformed
480 coordinates, this corresponds to a prior assumption of uniform compositions at all temperatures, i.e. all
481 components are equally abundant.

482
483 The residuals for the forward model are shown in Figure 10. The clear pattern in the residuals does not
484 necessarily indicate model misspecification, since no explicit noise model is specified for temperatures.
485 Predictive distributions are to be interpreted in the Bayesian sense, in that they represent a 'degree of belief'
486 in temperatures given the model and the modern data. The residual pattern is similar to that of the random
487 forest (Figure 4) with two clear downward slopes, suggesting again that the data are clustered into
488 temperatures above and below 16-17 °C, and that predictions tend towards temperatures at the centres of
489 these clusters.

490
491 An advantage of the forward modelling approach is that the inversion can incorporate substantive prior
492 information about temperatures for individual data points. In particular, other proxy systems can be used to
493 elicit prior distributions over temperatures to constrain GDGT-based predictions, particularly when
494 attempting to reconstruct ancient climates with no modern analogue in GDGT-space. We emphasise that
495 outside the range of the modern data, the utility of the models is almost solely due to the prior information
496 included in the reconstruction. At present, the only priors being used in the forward model prescribe a
497 reasonable upper limit and lower limit on temperatures (see Supplementary Information). The only way to
498 improve these reconstructions will be for future iterations to incorporate prior information from other
499 proxies. It is worth noting that the predictive uncertainty, while reasonably well-described by the standard
500 deviation in cases where ancient data lie quite close to the modern data in GDGT space, can be highly
501 multimodal (Fig. 11). This is the case when estimates are significantly outside of the modern calibration
502 dataset, such as low latitude data in the Cretaceous, or where there is considerable scatter in the modern
503 calibration data, for example in the low temperature range (<5 °C).

504

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

506

507 In principle, the predictors described above can be applied directly to ancient data, such as data from the
508 Eocene or Cretaceous (Inglis et al., 2015; O'Brien et al., 2017). In practice, one should be careful with
509 using models outside their domain of applicability. The machine learning tools described above, which are
510 ultimately based on the analysis of nearby calibration data in GDGT space, are fundamentally designed for
511 *interpolation*. To the extent that ancient data occupy a very different region in GDGT space, *extrapolation*
512 is required, which the models do not adequately account for. The divergence between modern calibration
513 data and ancient data is evident from Fig. 12, which shows histograms of minimum normalised distances
514 between 'high quality' Eocene/Cretaceous data points (those that passed the screening tests applied by
515 O'Brien et al., 2017 and Inglis et al., 2015) and the nearest point in the full modern data set. We strongly
516 recommend the use of the weighted distance metric (D_{nearest}) as a screening method to determine whether
517 the modern core top GDGT assemblage data is an appropriate basis for ancient SST estimation on a case-
518 by-case basis. Note that this distance measure is weighted by the scale length of the relevant parameter as
519 estimated by the Gaussian process emulator in order to quantify the relative position of ancient GDGT
520 assemblages to the modern core-top data. By using the GP-estimated covariance as the distance metric, we
521 account for the sensitivity of different GDGT components to temperature. Our inference is that samples
522 with $D_{\text{nearest}} > 0.5$, *regardless of the calibration model or approach applied*, are unlikely to generate
523 temperature estimates that are much better than informed guesswork. In these instances, in both our GPR
524 and Fwd models, the constraints provided by the modern calibration data set are so weak that estimates of
525 temperature have large uncertainty bands that are dictated by model priors; i.e. are unconstrained by the
526 calibration data (e.g., Figure 13 and Figure 14). This uncertainty is not apparent from estimates generated
527 by BAYSPAR or TEX_{86}^H models, although the underlying and fundamental lack of constraints are the same.
528 While 93% of validation data points in the modern data have $D_{\text{nearest}} < 0.5$, this is the case for only 33% of
529 Eocene samples and 3% for Cretaceous samples.

530

531 Where ancient GDGT distributions lie far from the modern calibration data set ($D_{\text{nearest}} > 0.5$), we argue that
532 there is no suitable set of modern analogue GDGT distributions from which to infer growth temperatures
533 for this ancient GDGT distribution. Both the GPR and Fwd models revert to imposed priors once the
534 distance from the modern calibration dataset increases. We propose that this is more rigorous and justified
535 model behavior than extrapolation of TEX_{86} or BAYSPAR predictors to non-analogue samples far from
536 the modern calibration data. As a result, the predictive models can only be applied to a subset of the Eocene
537 and Cretaceous data. We also note that there are two broad, non-mutually-exclusive categories of samples
538 that lie far from the modern calibration dataset ($D_{\text{nearest}} > 0.5$), the first are samples that seem to lie 'beyond'
539 the temperature-GDGT calibration relationship, likely with (unconstrained) GDGT formation temperatures

540 higher than the modern core-top calibrations; the second are samples with anomalous GDGT distributions
541 lying on the margins of, or far away from the main GDGT clustering in 6-dimensional space (see outliers
542 in Fig. 8).

543
544 Given the (current) limit on natural mean annual surface ocean temperatures of ~30 °C, extending the
545 GDGT-temperature calibration might be possible through, 1) integration of full GDGT abundance
546 distributions produced in high temperature culture, mesocosm or artificially warmed sea surface
547 conditions into the models; followed by, 2) validation through robust inter-comparisons of any new
548 GDGT palaeothermometer for high temperatures conditions with other temperature proxies from past
549 warm climate states. As discussed in the introduction, the first approach is limited by the ability of culture
550 or mesocosm experiments to accurately represent the true diversity and growth environments and
551 dynamics of natural microbial populations. Such studies clearly indicate a more complex, community-
552 scale control on changing GDGT relative abundances to growth temperatures (e.g., Elling et al., 2015).
553 Community-scale temperature dependency can be modelled relatively well with analyses of natural
554 production preserved in core-top sediments, especially with more sophisticated model fitting, including
555 the GPR and Fwd model presented here. Above ~30°C, however, the behavior of even single strains of
556 mesophilic archaea are not well-constrained by culture experiments, and the natural community-level
557 responses above this temperature are, so far, completely unknown. While there is evidence for the
558 temperature-sensitivity of GDGT production by thermophilic and acidophilic archaea in older papers (de
559 Rosa et al., 1980; Gliozzi et al., 1983), recent work, characterised by more precise phylogenetic and
560 culturing techniques show a more complex relationship between GDGT production and temperature.
561 Elling et al., (2017) highlight that there is no correlation between TEX₈₆ and growth temperature in a
562 range of phylogenetically different thaumarchaeal cultures - including thermophilic species. Bale et al.
563 (2019) recently cultured *Candidatus nitrosotenuis uzonensis* from the moderately thermophilic order
564 Nitrosopumilales (that contains many mesophilic marine strains). They found no correlation between
565 TEX₈₆ calibrations (either the Kim et al., core-top or Wuchter et al. 2004 and Schouten et al., 2008
566 mesocosm calibrations) with membrane lipid composition at different growth temperatures (37°C, 46°C,
567 and 50°C) and found that phylogeny generally seems to have a stronger influence on GDGT distribution
568 than temperature. In view of these existing data, we see no robust justification at present for the
569 extrapolation of modern core-top calibration data sets into the unknown above 30 °C, although the
570 coherent patterns apparent across GDGT space, between modern, Eocene and Cretaceous data (Figure 7),
571 do provide some grounds for hope that the extension of GDGT palaeothermometry beyond 30°C might be
572 possible in future.

573

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

575

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

584

585 Following Tierney and Tingley (2014) we use a reduced calibration data set, with the exclusion of Arctic
586 data with observed SSTs less than 3°C (“NoNorth / TT13” of Tierney and Tingley (2014)) but with the
587 inclusion of additional core top data from Seki et al. (2014). Full details of this calibration dataset are
588 provided in the Supplementary Information; to distinguish from the original OPTiMAL calibration data,
589 which included the Arctic data <3°C, we refer to the original data as “Op1” and the new calibration dataset
590 as “Op3”. An “Op2” is also available, which is the same as Op1 except that it excludes the Seki et al. (2014)
591 data. In sensitivity tests to a range of applications across Quaternary and deep-time datasets, calibration
592 Op1 and Op2 performed in almost identical fashion. The performance of Op1 and Op3 were very similar
593 in most applications, except in applications to the paleo-Arctic (see below), where the inclusion of modern
594 Arctic calibration data (Op1) provided closer calibration constraints to the paleo-data. Although
595 superficially this may be regarded as beneficial, in these instances the paleo-data have previously been
596 rejected because of a potential bias by non-marine inputs indicated by high BIT indices (Sluijs et al. 2020).
597 In this case, either the modern Arctic calibration data is impacted by similar non-thermal processes,
598 generating unusual GDGT abundance patterns, which are not appropriate to use for SST calibration, or,
599 there could be some consistency between the modern and ancient GDGT production by marine archaea in
600 the Arctic which may help in the understanding of GDGT-based paleothermometry in this unusual
601 environment (Sluijs et al. 2020). The D_{nearest} methodology may prove useful in quantifying analogue and
602 non-analogue behavior through time in such conditions. For the purposes of this study, however, we take
603 the conservative approach, and one that maintains a more consistent calibration basis with BAYSPAR, by
604 using OPTiMAL calibration Op3 in the remainder of this discussion, and recommend its use in future
605 applications of OPTiMAL.

606

607 To investigate the behaviour of the new OPTiMAL model, we compare temperature predictions including
608 uncertainties for the Eocene and Cretaceous datasets, made by OPTiMAL and the BAYSPAR methodology
609 of Tierney and Tingley (2014) (Figures 13 and 14), using the default priors specified in the model code for
610 the BAYSPAR estimation. The OPTiMAL model systematically estimates slightly cooler temperatures
611 than BAYSPAR, with the biggest offsets below ~ 15 °C (Figure 13). Fossil GDGT assemblages that fail the
612 D_{nearest} test are shown in grey, which clearly illustrate the regression to the mean in the OPTiMAL model,
613 whereas BAYSPAR continues to make SST predictions up to and exceeding 40 °C for these “non-analogue”
614 samples due to the fact that BAYSPAR assumes that higher TEX_{86} values equate to higher temperatures as
615 part of the functional form of the model, whereas the GPR model is agnostic on this. A comparison of error
616 estimation between OPTiMAL and BAYSPAR is shown in Figure 14. For most of the predictive range
617 below the D_{nearest} cut-off of 0.5, OPTiMAL has smaller predicted uncertainties than BAYSPAR, especially
618 in the lower temperature range. As D_{nearest} increases, i.e. as the fossil GDGT assemblage moves further from
619 the constraints of the modern calibration dataset, the error on OPTiMAL increases, until it reaches the
620 standard deviation of the modern calibration dataset (i.e., is completely unconstrained). In other words,
621 OPTiMAL generates maximum likelihood SSTs with robust confidence intervals, which appropriately
622 reflect the relative position of an ancient sample used for SST estimation and the structure of the modern
623 calibration data set. Where there are strong constraints from near analogues in the modern data,
624 uncertainties will be small, where there are weak constraints, uncertainty increases. In contrast, while
625 uncertainty bounds do increase when BAYSPAR is used to extrapolate beyond the modern calibration, they
626 are not as large as Optimal because BAYSPAR assumes a linear increase in SST at higher TEX values.

627
628 We also provide an initial assessment of the inter-relationship between standard screening indices and
629 D_{nearest} , for the Eocene and Cretaceous compilations where the data are available to calculate these measures
630 (Figure 15). For ease of comparison between Eocene and Cretaceous datasets and visualization of the
631 majority of the data, extreme outliers ($D_{\text{nearest}} > 4.0$) are not shown. The metrics include the BIT index
632 (Hopmans et al., 2004; Weijers et al., 2006), the Methane Index (MI; Zhang et al., 2011), the deviation
633 between TEX_{86} and the Ring Index (ΔRI ; Zhang et al., 2016) and the %GDGT-0 (Blaga et al., 2009;
634 Sinninghe Damsté et al., 2012). The standard screening levels for each of these metrics, as used in previous
635 paleo-compilations (O’Brien et al. 2017), are shown in the blue shaded areas on Figure 15 (BIT > 0.5; MI
636 > 0.5; $\Delta\text{RI} > 0.3$; %GDGT-0 > 67%) – data points within these areas fail the standard screening. Also shown
637 on Figure 15 is the region where data pass our D_{nearest} screening requirement (grey shaded vertical region).
638 In nearly all cases GDGT assemblages that fail these traditional screening tests also have D_{nearest} values that
639 exceed 0.5 – i.e. “abnormal” GDGT assemblages are well screened D_{nearest} . The main exception to this is
640 the BIT index in the Eocene data set, where 15 samples have high BIT values (>0.5) but have GDGT

641 assemblages that are close to modern analogues in the calibration dataset ($D_{\text{nearest}} < 0.5$). Of these samples,
642 9 are from the Arctic Ocean between the PETM and ETM2, an interval noted for its relatively high BIT
643 index values (Sluijs et al. 2020), 3 are from the Eocene-Oligocene transition of ODP Site 1218 (eastern
644 Equatorial Pacific) (Liu et al. 2009), 2 are from the middle Eocene of Seymour Island (Douglas et al. 2014),
645 and 1 is from the late Eocene of DSDP Site 511, which has been already noted as an individual sample with
646 anomalous high BIT in this dataset (Liu et al. 2009; Inglis et al. 2015). Although high BIT at ODP Site
647 1218 has been inferred to represent “relatively high terrestrial input” (Inglis et al. 2015) this seems unusual
648 for a fully pelagic site situated on oceanic crust >3000 km away from the nearest continental landmass.
649 Interpreting high BIT values as exclusively caused by terrestrial organic components appears problematic
650 in this instance, especially as $D_{\text{nearest}} < 0.5$ give some assurance that these GDGT assemblages from ODP
651 Site 1218 are well-modelled by the modern calibration dataset. GDGT assemblages from Seymour Island
652 associated with high BIT values (> 0.4) appear to have an impact on the $\text{TEX}_{86}^{\text{H}}$ SST proxy (Inglis et al.
653 2015), but the 2 samples that fail BIT (> 0.5) but pass D_{nearest} (< 0.5) give OPTiMAL SSTs consistent (5-
654 6°C) with the SSTs from samples that pass all other screening and D_{nearest} ($\sim 4\text{-}7^{\circ}\text{C}$). In summary, the
655 relationship between D_{nearest} and BIT suggests that BIT is not always closely coupled to GDGT assemblages
656 that are strongly divergent from the modern calibration dataset.

657
658 With respect to the other screening indices there are clear indications that increased distance from the
659 modern calibration (increased D_{nearest}) is associated with a trend towards the “thresholds of failure” in the
660 screening indices. This pattern is most clear with the ΔRI in both the Cretaceous and the Eocene data, as
661 increasing numbers of samples fail ΔRI as D_{nearest} increases. This supports ΔRI as a robust methodology for
662 identifying samples that strongly diverge from the expected temperature-dependence of GDGT
663 assemblages as modelled by TEX_{86} in the modern calibration dataset. There are, however, samples that pass
664 $D_{\text{nearest}} < 0.5$ but fail ΔRI in both the Eocene and Cretaceous datasets – these must have “near neighbours”
665 in the modern calibration data, but yet have a temperature-sensitivity that is less well-modelled by TEX_{86}
666 (divergence between RI and TEX_{86}). Conversely there are many Eocene and Cretaceous data points with
667 $\Delta\text{RI} < 0.3$, but which fail D_{nearest} (> 0.5). These data most likely represent GDGT assemblages formed at
668 high temperatures, beyond the range of the modern calibration data.

669
670 To investigate these behaviours requires the publication of the full range GDGT abundance data. Whilst
671 key compilations of Eocene and Cretaceous GDGT data have strongly encouraged the release of such
672 datasets (Lunt et al. 2012; Dunkley Jones et al. 2013; Inglis et al. 2015; O’Brien et al. 2017), most Neogene
673 studies only publish TEX_{86} values. Without full GDGT assemblage data neither OPTiMAL nor other

674 detailed assessments of GDGT behaviour and type can be made, and we would strongly encourage authors,
675 reviewers and editors to ensure the publication of full GDGT assemblages in future.

676

677 Finally, to test the behavior of OPTiMAL within established SST time series, we provide three examples
678 two from the late Pleistocene to Holocene (Figure 16) and one from the Eocene (Figures 17 and 18). For
679 the Pleistocene to Holocene examples OPTiMAL SSTs are shown against estimates from BAYSPAR and
680 the alkenone-based $U^{k'}_{37}$ temperature proxy. The first of these timeseries is from GeoB 7702-3 in the
681 Eastern Mediterranean and spans the last 26 kyr, including data spanning Termination I (Castañeda et al.,
682 2010). The second is from ODP Site 1146 in the South China Sea and spans the last 350 kyr (Thomas et al.
683 2014). In both records the long-term dynamics are consistent between the independent $U^{k'}_{37}$ SST proxy
684 and both BAYSPAR and OPTiMAL. In the Eastern Mediterranean OPTiMAL SSTs are slightly cooler in
685 the glacial and warmer in the Holocene than the other proxies. In the South China Sea, OPTiMAL is again
686 cooler than BAYSPAR during glacial intervals, but at this location is in closer agreement than BAYSPAR
687 with the $U^{k'}_{37}$ SST proxy through most of the record. In both these examples, we show the 5th and 95th
688 percentiles for OPTiMAL and those reported by the BAYSPAR methodology.

689

690 The final example is from the latest Paleocene to early Eocene of IODP Expedition 302 Hole 4A on
691 Lomonosov Ridge (Sluijs et al. 2006; Sluijs et al. 2009; Sluijs et al. 2020). This site is useful as it has been
692 the focus of detailed reassessment and reanalysis, using most of the available screening methodologies to
693 detect aberrant GDGT assemblages (Sluijs et al. 2020). Here we use this recently published data to compare
694 the new $D_{nearest}$ screening metric against multiple other screening protocols (Figure 17). We also show both
695 $D_{nearest}$ values and OPTiMAL SST estimates for two models – one with modern Arctic data with $SST < 3^{\circ}C$
696 included in the calibration (OPTiMAL_{Arctic}; equivalent to calibration dataset Op1 first present by Eley et al.
697 2019) and one with this data excluded (OPTiMAL_{noArctic}; equivalent to the new calibration dataset Op3). It
698 is clear from the pattern of $D_{nearest}$ for these two options, that the inclusion of modern Arctic data provides
699 more calibration data that are closer to the Eocene paleo-Arctic, to the extent that substantially more
700 samples pass the $D_{nearest} < 0.5$ constraint, especially in pre-ETM2 interval from ~372 to 376 mcd. This
701 interval contains, however, samples with the highest BIT values of the succession (> 0.4), and elevated ΔRI
702 (> 0.3). With these other “warning signs” concerning the reliability of GDGT assemblages for SST
703 estimation in this interval, the relatively low $D_{nearest}$ values are most likely to represent some similarity in
704 the non-thermal controls on GDGT assemblages between the modern and paleo-Arctic. More work needs
705 to be done to constrain the reliability of temperature-dependence and archaeal GDGT production in these
706 modern high latitude systems before we can have confidence in their inclusion in calibration datasets for
707 paleo-SST estimation. It is on the basis that we recommend users of OPTiMAL use the the “noArctic”

708 (Op3) calibration for the time being. The OPTiMAL methodology does, however, offer a simple means to
709 integrate new robust calibration data, and a method to explore the distance relationships between modern
710 and ancient GDGT production.

711
712 Considering the “noArctic” D_{nearest} and OPTiMAL SSTs for Exp. 302 Hole 4A, it is clear that of all the
713 screening methods, D_{nearest} shows the strongest similarity to ΔRI – with high (“failure”) values in the pre-
714 PETM and then again between ~371 and 376 mcd, and even picking up the same short-lived “failure”
715 intervals, or spikes, between 368 and 371 mcd. SST estimates based on OPTiMAL show broadly similar
716 trends to $\text{TEX}_{86}^{\text{H}}$ and BAYSPAR, with a warm PETM, cooling post-PETM and then warming again into
717 ETM2. It should be noted, however, that peak temperatures for OPTiMAL are ~5°C cooler than $\text{TEX}_{86}^{\text{H}}$
718 and BAYSPAR (e.g. PETM SSTs <20°C for OPTiMAL and > 25°C for $\text{TEX}_{86}^{\text{H}}$ and BAYSPAR), and show
719 more cooling post-PETM, with SST estimates of ~10°C (OPTiMAL_{noArctic}) as opposed to ~20°C for $\text{TEX}_{86}^{\text{H}}$
720 and BAYSPAR.

721

722 5. Conclusions

723

724 Although the fundamental issue of non-analogue behaviour is a key problem for GDGT-temperature
725 estimation, it has an undue impact on the community’s general confidence in this method. In part, this is
726 because these issues have not been clearly stated and circumscribed - rather they have been allowed to erode
727 confidence in the GDGT-based methodology through the use of GDGT-based palaeothermometry far
728 outside the modern constraints on the behavior of this system. The use of GDGT abundances to estimate
729 temperatures in clearly non-analogue conditions is, at present, problematic on the basis of the available
730 calibration constraints or a good understanding of underlying biophysical models. We hope that this study
731 prompts further investigations that will improve these constraints for the use of GDGTs in deep-time
732 paleoclimate studies, where they clearly have substantial potential as temperature proxies. Temperature
733 estimates based on fossil GDGT assemblages that are within range of, or similar to, modern GDGT
734 calibration data, do, however, rest on a strong, underlying temperature-dependence observed in the
735 empirical data. With no effective means of separating the “good from the bad” can lead to either false
736 confidence and inappropriate inferences in non-analogue conditions, or a false pessimism when ancient
737 samples are actually well constrained by modern core-top assemblages.

738

739 In this study, we apply modern machine-learning tools, including Gaussian Process Emulators and forward
740 modelling, to improve temperature estimation and the representation of uncertainty in GDGT-based SST
741 reconstructions. Using our new nearest neighbour test, we demonstrate that >60% of Eocene, and >90% of

742 Cretaceous, fossil GDGT distribution patterns differ so significantly from modern as to call into question
743 SSTs derived from these assemblages. For data that does show sufficient similarity to modern, we present
744 OPTiMAL, a new multi-dimensional Gaussian Process Regression tool which uses all six GDGTs (GDGT-
745 0, -1, -2, -3, Cren and Cren') to generate an SST estimate with associated uncertainty. The key advantages
746 of the OPTiMAL approach are: 1) that these uncertainty estimates are intrinsically linked to the strength of
747 the relationship between the fossil GDGT distributions and the modern calibration data set, and 2) by
748 considering all GDGT compounds in a multi-dimensional regression model it avoids the dimensionality
749 reduction and loss of information that takes place when calibrating single parameters (TEX_{86}) to
750 temperature. The methods presented above make very few assumptions about the data. We argue that such
751 methods are appropriate with the current absence of any reasonable mechanistic model for the data
752 generating process, in that they reflect model uncertainty in a natural way. Finally, we note the potential
753 for multi-proxy machine learning approaches, synthesising data from other palaeothermometers with
754 independent uncertainties and biases, to improve calibration of ancient GDGT-derived SST reconstructions.

755
756

757 **Acknowledgements:**

758 TDJ, JAB, IM, KME and YE acknowledge NERC grant NE/P013112/1. SEG was supported by NERC
759 Independent Research Fellowship NE/L011050/1 and NERC large grant NE/P01903X/1. WT
760 acknowledges the Wellcome Trust (grant code: 1516ISSFFEL9, www.wellcome.ac.uk/) for funding a
761 parameterisation workshop at the University of Birmingham (UK). WT, TDJ and IM would like to thank
762 the BBSRC UK Multi-Scale Biology Network Grant No. BB/M025888/1. IM is a recipient of the
763 Australian Research Council Future Fellowship, FT190100574. We would like to thank the extensive and
764 constructive reviews of Jessica Tierney, Yige Zhang, Huan Yang and an anonymous reviewer, as well as
765 the great patience of Editor Alberto Reyes. We would also like to give especial thanks to Elizabeth
766 Thomas and Isla Castaneda for digging up primary GDGT data when access to labs was far from straight
767 forward.

768
769

770 **Figure Captions:**

771

772 **Figure 1.** A histogram of the normalised distance to the nearest neighbour in GDGT space ($D_{x,yt}$) for all
773 samples in the modern calibration dataset of Tierney and Tingley (2015).

774

775 **Figure 2.** The error of the nearest-neighbour temperature ($D_{x,y}$) predictor, for modern core-top data, as a
776 function of the distance to the nearest calibration sample.

777

778 **Figure 3.** Top: The temperature of the modern data set as a function of the TEX_{86} value, showing a clear
779 linear correlation between the two, but also significant scatter. Bottom: the error of the predictor based on
780 the nearest TEX_{86} calibration point.

781

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

783

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

786

787 **Figure 6.** Modern and ancient data projected onto the first two compositional principal components. Black:
788 Modern; Blue: Eocene (Inglis et al., 2015); Red: Cretaceous (O'Brien et al., 2017).

789

790 **Figure 7.** Diffusion map projection of the modern and ancient data. Black: Modern; Blue: Eocene (Inglis
791 et al., 2015); Red: Cretaceous (O'Brien et al., 2017). Separate clusters marked 'A' are the outlying
792 Cretaceous points with high GDGT-3 values. Branch 'B' is dominated by modern data points; branch 'C'
793 by Cretaceous data.

794

795 **Figure 8.** The first diffusion component as a function of TEX_{86} . Some outlying points have been excluded
796 from the plot for the purposes of visualisation. Black: Modern; Blue: Eocene (Inglis et al., 2015); Red:
797 Cretaceous (O'Brien et al., 2017).

798

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

800

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

802

803 **Figure 11.** The posterior distributions over temperature from the forward model for selected examples of
804 high and low temperature, Eocene and Cretaceous, data points. The Gaussian error envelope from the GPR
805 model is shown for comparison.

806

807 **Figure 12.** A histogram of normalised distances to the nearest sample in the modern data set for Eocene
808 and Cretaceous data, excluding samples that had been screened out in previous compilations using BIT, MI
809 and RI following the approach of (Inglis et al., 2015; O'Brien et al., 2017).

810

811 **Figure 13.** Comparison of temperature estimates for the BAYSPAR and the OPTiMAL GPR model, greyed
812 out data fails the $D_{nearest}$ test (>0.5), and the colour scaling reflects $D_{nearest}$ values for those datapoints that
813 pass. Note that outside of the constraints of the modern calibration (training) dataset, ($D_{nearest}$ test >0.5) the
814 GPR model temperature estimates revert to the mean value of the calibration dataset, with an uncertainty
815 that reverts to the standard deviation of the training data.

816
817 **Figure 14.** Inter-comparison of temperature estimates and standard errors (y-axis) for compiled Eocene
818 and Cretaceous data calculated using OPTiMAL (top) and BAYSPAR (bottom). Greyed out data fails the
819 $D_{nearest}$ test (>0.5), and the colour scaling reflects $D_{nearest}$ values for those datapoints that pass. The black
820 dashed line shows the $D_{nearest}$ threshold (>0.5).

821
822 **Figure 15.** Comparison of $D_{nearest}$ against standard screening indices, BIT and MI index, ΔRI and
823 %GDGT-O for the Eocene (Inglis et al., 2015) and Cretaceous (O'Brien et al., 2017) datasets. Blue
824 shaded regions show the standard cut-off points for these indices (see text); grey shaded region highlights
825 data that are below the $D_{nearest}$ threshold of 0.5. The outlined black box is the region of data that fails
826 traditional screening indices but passes $D_{nearest}$ (<0.5).

827
828 **Figure 16.** Late Pleistocene to Holocene GDGT-derived OPTiMAL palaeotemperatures compared to
829 BAYSPAR and U^{k}_{37} SSTs. Shaded regions represent reported 5th and 95th percentile confidence intervals.
830 Top panel - Eastern Mediterranean data from core GeoB 7702-3 (Castaneda et al. 2010); bottom panel –
831 South China Sea data from ODP Site 1146 (Thomas et al. 2014).

832
833 **Figure 17.** Comparison of GDGT screening indices, TEX_{86}^H , BAYSPAR and OPTiMAL SSTs from the
834 Eocene Arctic Site IODP Expedition 302 Hole 4A. Data and figures modified from the most recent
835 reassessment by Sluijs et al. (2020).

836
837 **References:**

838 Aitchison, J.: The Statistical Analysis of Compositional Data. J. R. Stat. Soc. Series B Stat. Methodol. 44,
839 139–160, 1982.

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

841 Aitchison, J., Greenacre, M.: Biplots of compositional data. J. R. Stat. Soc. Ser. C Appl. Stat. 51, 375–
842 392, 2002.

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

845 Bale, N. J., Palatinszky, M., Rijpstra, I. C., Herbold, C. W., Wagner, M., Sinnighe Damste, J. S.:
846 Membrane lipid composition of the moderately thermophilic ammonia-oxidizing Archaeon

847 “*Candidatus Nitrosotenus uzonensis*” at different growth temperatures, Applied and
848 Environmental Microbiolpogy, DOI: 10.1128/AEM.01332-19, 2019.

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

851 Bijl, P. K., Bendle, J.A.P., Bohaty, S.M., Pross, J., Schouten, S., Tauxe, L., Stickley, C., McKay, R.M.,
852 Röhl, U., Olney, M., Slujis, A., Escutia, C., Brinkhuis, H. and Expedition 318 Scientists.: Eocene
853 cooling linked to early flow across the Tasmanian Gateway. *Proc. Natl. Acad. Sci. U.S.A.*, 110,
854 9645–9650, 2013.

855 Brassell, S. C.: Climatic influences on the Paleogene evolution of alkenones, *Paleoceanography*, 29, 255-
856 272, doi:10.1002/2013PA002576, 2014.

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

862 Cadillo-Quiroz, H., Didelot, X., Held, N. L., Herrera, A., Darling, A., Reno, M. L., Krause, D. J.,
863 Whitaker, R. J.: Patterns of gene flow define species of thermophilic Archaea, *PLOS*, Cadillo-
864 Quiroz, H. et al. (2012) *PLOS Biology*, <https://doi.org/10.1371/journal.pbio.1001265>, 2012.
865

866 Castañeda, I. S., E. Schefuß, J. Pätzold, J. S. Sinninghe Damsté, S. Weldeab, and S. Schouten.:
867 Millennial-scale sea surface temperature changes in the eastern Mediterranean (Nile River Delta
868 region) over the last 27,000 years, *Paleoceanography*, 25, PA1208, doi:10.1029/2009PA001740,
869 2010.
870

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

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

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

880 De Rosa, M., Esposito, E., Gambacorta, A., Nicolaus, B., Bu’Lock, J. D.: Effects of temperatures on ether
881 lipid composition of *Caldariella acidophila*, 19, 827 – 831, 1980.

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

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

- 888 Egozcue, J.J., Pawlowsky-Glahn, V., Mateu-Figueras, G., Barceló-Vidal, C.: Isometric Logratio
889 Transformations for Compositional Data Analysis. *Math. Geol.* 35, 279–300, 2003.
- 890 Eley, Y. L., Thompson, W., Greene, S. E., Mandel, I., Edgar, K., Bendle, J. A., and Dunkley Jones, T.:
891 OPTiMAL: A new machine learning approach for GDGT-based palaeothermometry, *Clim. Past*
892 *Discuss.*, <https://doi.org/10.5194/cp-2019-60>, in review, 2019.
- 893 Elling, F. J., Könneke, M., Lipp, J. S., Becker, K. W., Gagen, E. J., and Hinrichs, K.-U.: Effects of growth
894 phase on the membrane lipid composition of the thaumarchaeon *Nitrosopumilus maritimus* and
895 their implications for archaeal lipid 20 distributions in the marine environment, *Geochimica et*
896 *Cosmochimica Acta*, 141, 579-597, 2014.
- 897 Elling, F. J., Könneke, M., Mußmann, M., Greve, A., and Hinrichs, K.-U.: Influence of temperature, pH,
898 and salinity on membrane lipid composition and TEX₈₆ of marine planktonic thaumarchaeal
899 isolates, *Geochimica et Cosmochimica Acta*, 171, 238-255, 2015.
- 900 Elling, F.J., Konnecke, M., Nicol, G. W., Stieglmeier, M., Bayer, B., Spieck, E., de la Torre, J. R.,
901 Becker, K. W., Thomm, M., Prosser, J. I., Herndl, G., Schleper, C., Hinrichs, K-U.:
902 Chemotaxonomic characterisation of the thaumarchaeal lipidome, *Environmental Microbiology* 19,
903 2681–2700 , 2017.
- 904
- 905 Filzmoser, P., Hron, K., Reimann, C.: Principal component analysis for compositional data with outliers.
906 *Environmetrics* 20, 621–632, 2009a.
- 907 Filzmoser, P., Hron, K., Reimann, C., Garrett, R.: Robust factor analysis for compositional data. *Comput.*
908 *Geosci.* 35, 1854–1861, 2009b.
- 909 Filzmoser, P., Hron, K., Reimann, C.: Interpretation of multivariate outliers for compositional data.
910 *Comput. Geosci.* 39, 77–85, 2012.
- 911 Gliozzi, A., Paoli, G., De Rosa, M., Gambacorta, A.: Effect of isoprenoid cyclization on the transition
912 temperature of lipids in thermophilic archaeobacteria, *Biochimica et Biophysica Acta (BBA)*
913 *Biomembranes*, 735, 234 – 242, 1983.
- 914 Haghverdi, L., Buettner, F., Theis, F.J: Diffusion maps for high-dimensional single-cell analysis of
915 differentiation data. *Bioinformatics* 31, 2989–2998, 2015.
- 916 Haghverdi, L., Büttner, M., Wolf, F.A., Buettner, F., Theis, F.J.: Diffusion pseudotime robustly
917 reconstructs lineage branching. *Nat. Methods* 13, 845–848, 2016.
- 918 Haslett, J., Whiley, M., Bhattacharya, S., Salter-Townshend, M., Wilson, S.P., Allen, J.R.M., Huntley, B.,
919 Mitchell, F.J.G.: Bayesian palaeoclimate reconstruction. *J Royal Statistical Soc A* 169, 395–438,
920 2006.
- 921 Herfort, L., Schouten, S., Boon, J. P., and Sinninghe Damsté, J. S.: Application of the TEX₈₆ temperature
922 proxy to the southern North Sea, *Organic Geochemistry*, 37, 1715-1726, 2006.
- 923 Hertzberg, J. E., Schmidt, M. W., Bianchi, T. S., Smith, R. K., Shields, M. R., & Marcantonio, F.:
924 Comparison of eastern tropical Pacific TEX₈₆ and Globigerinoides ruber Mg/Ca derived sea surface
925 temperatures: Insights from the Holocene and Last Glacial Maximum. *Earth and Planetary Science*
926 *Letters*, 434, 320–332, 2016.

- 927 Hollis, C. J., Taylor, K. W. R., Handley, L., Pancost, R. D., Huber, M., Creech, J. B., Hines, B. R.,
928 Crouch, E. M., Morgans, H. E. G., Crampton, J. S., Gibbs, S., Pearson, P. N., and Zachos, J. C.:
929 Early Paleogene temperature history of the Southwest Pacific Ocean: Reconciling proxies and
930 models, *Earth and Planetary Science Letters*, 349–350, 53–66, 2012.
- 931 Hollis, C. J., Dunkley Jones, T., Anagnostou, E., Bijl, P. K., Cramwinckel, M. J., Cui, Y., Dickens, G. R.,
932 Edgar, K. M., Eley, Y., Evans, D., Foster, G. L., Frieling, J., Inglis, G. N., Kennedy, E. M.,
933 Kozdon, R., Laurentano, V., Lear, C. H., Littler, K., Meckler, N., Naafs, B. D. A., Pälike, H.,
934 Pancost, R. D., Pearson, P., Royer, D. L., Salzmann, U., Schubert, B., Seebeck, H., Sluijs, A.,
935 Speijer, R., Stassen, P., Tierney, J., Tripathi, A., Wade, B., Westerhold, T., Witkowski, C., Zachos,
936 J. C., Zhang, Y. G., Huber, M., and Lunt, D. J.: The DeepMIP contribution to PMIP4:
937 methodologies for selection, compilation and analysis of latest Paleocene and early Eocene climate
938 proxy data, incorporating version 0.1 of the DeepMIP database, *Geosci. Model Dev. Discuss.*,
939 <https://doi.org/10.5194/gmd-2018-309>, in review, 2019.
- 940 Hollis, C. J., Handley, L., Crouch, E. M., Morgans, H. E., Baker, J. A., Creech, J., Collins, K. S., Gibbs,
941 S. J., Huber, M., Schouten, S.: Tropical sea temperatures in the high-latitude South Pacific during
942 the Eocene. *Geology*, 37, 99–102, 2009.
- 943 Hopmans, E. C., Weijers, J. W. H., Schefuss, E., Herfort, L., Sinninghe Damsté, J. S., Schouten, S.: A
944 novel proxy for terrestrial organic matter in sediments based on branched and isoprenoid tetraether
945 lipids, *Earth and Planetary Science Letters*, 224, 107–116, 2004.
- 946 Huguet C, Kim J-H, Sinninghe Damsté J.S., Schouten S: Reconstruction of sea surface temperature
947 variations in the Arabian Sea over the last 23 kyr using organic proxies (TEX₈₆ and UK 0 37 .
948 *Paleoceanography* 21(3): PA3003, 2006.
- 949 Hurley, S. J., Elling, F. J., Könneke, M., Buchwald, C., Wankel, S. D., Santoro, A. E., Lipp, J.S.,
950 Hinrichs, K., Pearson, A.: Influence of ammonia oxidation rate on thaumarchaeal lipid composition
951 and the TEX₈₆ temperature proxy. *Proceedings of the National Academy of Sciences*, 113, 7762–
952 7767, 2016.
- 953 Inglis, G. N., Farnsworth, A., Lunt, D., Foster, G. L., Hollis, C. J., Pagani, M., Jardine, P. E., Pearson, P.
954 N., Markwick, P., Galsworthy, A. M. J., Raynham, L., Taylor, K. W. R., and Pancost, R. D.:
955 Descent toward the Icehouse: Eocene sea surface cooling inferred from GDGT distributions,
956 *Paleoceanography*, 30, 1000–1020, 2015.
- 957 Jenkyns H.C., Schouten-Huibers L., Schouten S., Damsté J.S.S.: Warm Middle Jurassic-Early Cretaceous
958 high-latitude sea-surface temperatures from the Southern Ocean. *Clim Past* 8 (1):215–226, 2012.
- 959 Kim, J.-H., Schouten, S., Hopmans, E. C., Donner, B., and Sinninghe Damsté, J. S.: Global sediment
960 core-top calibration of the TEX₈₆ paleothermometer in the ocean, *Geochimica et Cosmochimica*
961 *Acta*, 72, 1154–1173, 2008.
- 962 Kim, J.-H., van der Meer, J., Schouten, S., Helmke, P., Willmott, V., Sangiorgi, F., Koç, N., Hopmans, E.
963 C., and Sinninghe Damsté, J. S.: New indices and calibrations derived from the distribution of
964 crenarchaeal isoprenoid tetraether lipids: Implications for past sea surface temperature
965 reconstructions, *Geochimica et Cosmochimica Acta*, 74, 4639–4654, 2010.

- 966 Linnert, C., Robinson, S. A., Lees, J. A., Bown, P. R., Perez-Rodriguez, I., Petrizzo, M. R., Falzoni, F.,
967 Littler, K., Antonio Arz, J., Russell, E. E. : Evidence for global cooling in the Late Cretaceous.
968 Nature Communications, 5, 1–7, 2014.
- 969 Liu, X-L., Zhu, C., Wakeham, S.G., Hinrichs, K-U.: In situ production of branched glycerol dialkyl
970 glycerol tetraethers in anoxic marine water columns, Marine Chemistry, 166, 1 – 8, 2014.
- 971 Lunt, D. J., Dunkley Jones, T., Heinemann, M., Huber, M., LeGrande, A., Winguth, A., Loptson, C.,
972 Marotzke, J., Tindall, J., 15 Valdes, P., Winguth, C.: A model-data comparison for a multi-model
973 ensemble of early Eocene atmosphere-ocean simulations: EoMIP, Clim. Past Discuss., 8, 1229-
974 1273, 2012.
- 975 Mentch, L., Hooker, G.: Quantifying Uncertainty in Random Forests via Confidence Intervals and
976 Hypothesis Tests. Journal of Machine Learning Research, 17, 1-41, 2016.
- 977 O'Brien, C. L., Robinson, S. A., Pancost, R. D., Sinninghe Damsté, J. S., Schouten, S., Lunt, D. J.,
978 Alsenz, H., Bornemann, 20 A., Bottini, C., Brassell, S. C., Farnsworth, A., Forster, A., Huber, B.
979 T., Inglis, G. N., Jenkyns, H. C., Linnert, C., Littler, K., Markwick, P., McAnena, A., Mutterlose,
980 J., Naafs, B. D. A., Püttmann, W., Sluijs, A., van Helmond, N. A. G. M., Vellekoop, J., Wagner, T.,
981 Wrobel, N. E.: Cretaceous sea-surface temperature evolution: Constraints from TEX₈₆ and
982 planktonic foraminiferal oxygen isotopes, Earth-Science Reviews, 172, 224-247, 2017.
- 983 Park, E., Hefter, J., Fischer, G., Mollenhauer, G.: TEX₈₆ in sinking particles in three eastern Atlantic
984 upwelling regimes. Organic Geochemistry, 124, 151–163, 2018.
- 985 Pearson, P. N., van Dongen, B. E., Nicholas, C. J., Pancost, R. D., Schouten, S., Singano, J. M., Wade, B.
986 S.: Stable warm tropical climate through the Eocene Epoch. Geology, 35, 211-214, 2007.
- 987 Polik, C. A., Elling, F. J., Pearson, A.: Impacts of Paleoecology on the TEX₈₆ Sea Surface Temperature
988 Proxy in the Pliocene-Pleistocene Mediterranean Sea. Paleoceanography and Paleoclimatology, 33,
989 1472–1489, 2018.
- 990 Qin, W., Amin, S. A., Martens-Habbena, W., Walker, C. B., Urakawa, H., Devol, A. H., Ingalls, A.E.,
991 Moffett, J.W., Ambrust, E.V., Stahl, D. A.: Marine ammonia-oxidizing archaeal isolates display
992 obligate mixotrophy and wide ecotypic variation. Proceedings of the National Academy of
993 Sciences of the United States of America, 111, 12504–12509, 2014.
- 994 Qin, W., Carlson, L. T., Armbrust, E. V., Stahl, D. A., Devol, A. H., Moffett, J. W., Ingalls, A. E.:
995 Confounding effects of oxygen and temperature on the TEX 86 signature of marine
996 Thaumarchaeota. Proceedings of the National Academy of Sciences, 112, 10979–10984, 2015.
- 997 Rasmussen, C.E., Nickisch, H.: Gaussian Processes for Machine Learning (GPML) Toolbox. J. Mach.
998 Learn. Res. 11, 3011–3015, 2010.
- 999 Rasmussen, C. E. & Williams, C. K. I.: Gaussian Processes for Machine Learning, the MIT Press, 2006,
1000 ISBN 026218253X.
- 1001 Sangiorgi, F., van Soelen Els, E., Spofforth David, J. A., Pälke, H., Stickley Catherine, E., St. John, K.,
1002 Koç, N., Schouten, S., Sinninghe Damsté Jaap, S., Brinkhuis, H.: Cyclicity in the middle Eocene
1003 central Arctic Ocean sediment record: Orbital forcing and environmental response,
1004 Paleoceanography, 23, 10.1029/2007PA001487, 2008.

- 1005 Schouten, E., Hopmans, E.C., Forster, A., Van Breugel, Y., Kuypers, M.M.M., Sinninghe Damsté, J.S.:
 1006 Extremely high seasurface temperatures at low latitudes during the middle Cretaceous as revealed
 1007 by archaeal membrane lipids. *Geology*, 31, 1069–1072, 2003.
- 1008 Schouten, S., Forster, A., Panoto, F. E., and Sinninghe Damsté, J. S.: Towards calibration of the TEX₈₆
 1009 palaeothermometer for 20 tropical sea surface temperatures in ancient greenhouse worlds, *Organic*
 1010 *Geochemistry*, 38, 1537-1546, 2007.
- 1011 Schouten, S., Hopmans, E. C., Sinninghe Damsté, J. S.: The organic geochemistry of glycerol dialkyl
 1012 glycerol tetraether lipids: A review, *Organic Geochemistry*, 54, 19-61, 2013.
- 1013 Schouten, S., Hopmans, E. C., Schefuß, E., Sinninghe Damsté, J. S.: Distributional variations in marine
 1014 crenarchaeotal membrane lipids: a new tool for reconstructing ancient sea water temperatures?
 1015 *Earth and Planetary Science Letters*, 204, 15 265-274, 2002.
- 1016 Seki, O., Bendle, J. A., Harada, N., Kobayashi, M., Sawada, K., Moossen, H., Sakamoto, T.: Assessment
 1017 and calibration of TEX₈₆ paleothermometry in the Sea of Okhotsk and sub-polar North Pacific
 1018 region: Implications for paleoceanography. *Progress in Oceanography*, 126, 254–266, 2014.
- 1019 Sluijs A, Schouten S, Pagani M, Woltering, M., Brinkhuis, H., Sinninghe Damsté, J.S., Dickens, G.R.,
 1020 Huber, M., Reichart, G., Stein, R., Matthiessen, J., Lourens, L.J., Pedentchouk, N., Backman, J.,
 1021 Moran, K. and the Expedition 320 Scientists: Subtropical arctic ocean temperatures during the
 1022 Palaeocene/Eocene thermal maximum. *Nature* 441, 610–613, 2006.
- 1023 Sluijs, A., Schouten, S., Donders, T. H., Schoon, P. L., Rohl, U., Reichart, G.-J., Sangiorgi, F., Kim, J.-
 1024 H., Sinninghe Damsté, J. S., Brinkhuis, H.: Warm and wet conditions in the Arctic region during
 1025 Eocene Thermal Maximum 2, *Nature Geosci*, 2, 777-780, 2009.
- 1026 Taylor, K. W. R., Willumsen, P. S., Hollis, C. J., Pancost, R. D.: South Pacific evidence for the long-term
 1027 climate impact of the Cretaceous/Paleogene boundary event, *Earth-Science Reviews*, 179, 287-302,
 1028 2018.
- 1029 Taylor, K. W., Huber, M., Hollis, C. J., Hernandez-Sanchez, M. T., Pancost, R. D.: Re-evaluating modern
 1030 and Palaeogene GDGT distributions: Implications for SST reconstructions, *Global and Planetary*
 1031 *Change*, 108, 158-174, 2013.
- 1032 Thomas, E. K., S. C. Clemens, W. L. Prell, T. D. Herbert, Y. Huang, Z. Liu, J. S. S. Damsté, Y. Sun, and
 1033 X. Wen.: Temperature and leaf wax $\delta^2\text{H}$ records demonstrate seasonal and regional controls on
 1034 Asian monsoon proxies, *Geology*, 42(12), 1075–1078, doi:10.1130/G36289.1, 2014.
- 1035 Tierney, J. E.: GDGT Thermometry: Lipid Tools for Reconstructing Paleotemperatures. Retrieved from
 1036 https://www.geo.arizona.edu/~jesst/resources/TierneyPSP_GDGTs.pdf, 2012
- 1037 Tierney, J. E., and Tingley, M. P.: A Bayesian, spatially-varying calibration model for the TEX₈₆ proxy.
 1038 *Geochimica et Cosmochimica Acta*, 127, 83-106, 2014.
- 1039 Tierney, J. E., and Tingley, M. P.: A TEX₈₆ surface sediment database and extended Bayesian calibration,
 1040 *Scientific data*, 2, 150029, 2015.
- 1041 Williams, C.K.I., and Rasmussen, C.E.: Gaussian processes for machine learning. MIT Press Cambridge,
 1042 MA, 2006.

1043 Wuchter, C., Schouten, S., Coolen, M. J. L., and Sinninghe Damsté, J. S.: Temperature-dependent
1044 variation in the distribution 30 of tetraether membrane lipids of marine Crenarchaeota: Implications
1045 for TEX₈₆ paleothermometry, *Paleoceanography and Paleoclimatology*.
1046 doi:10.1029/2004PA001041, 2004.

1047 Zhang, Y. G., and Liu, X.: Export Depth of the TEX₈₆ Signal. *Paleoceanography and Paleoclimatology*.
1048 doi.org/10.1029/2018PA003337, 2018.

1049 Zhang, Y. G., Pagani, M., Wang, Z.: Ring Index: A new strategy to evaluate the integrity of TEX₈₆
1050 paleothermometry, *Paleoceanography*, 31, 220-232, 2016.

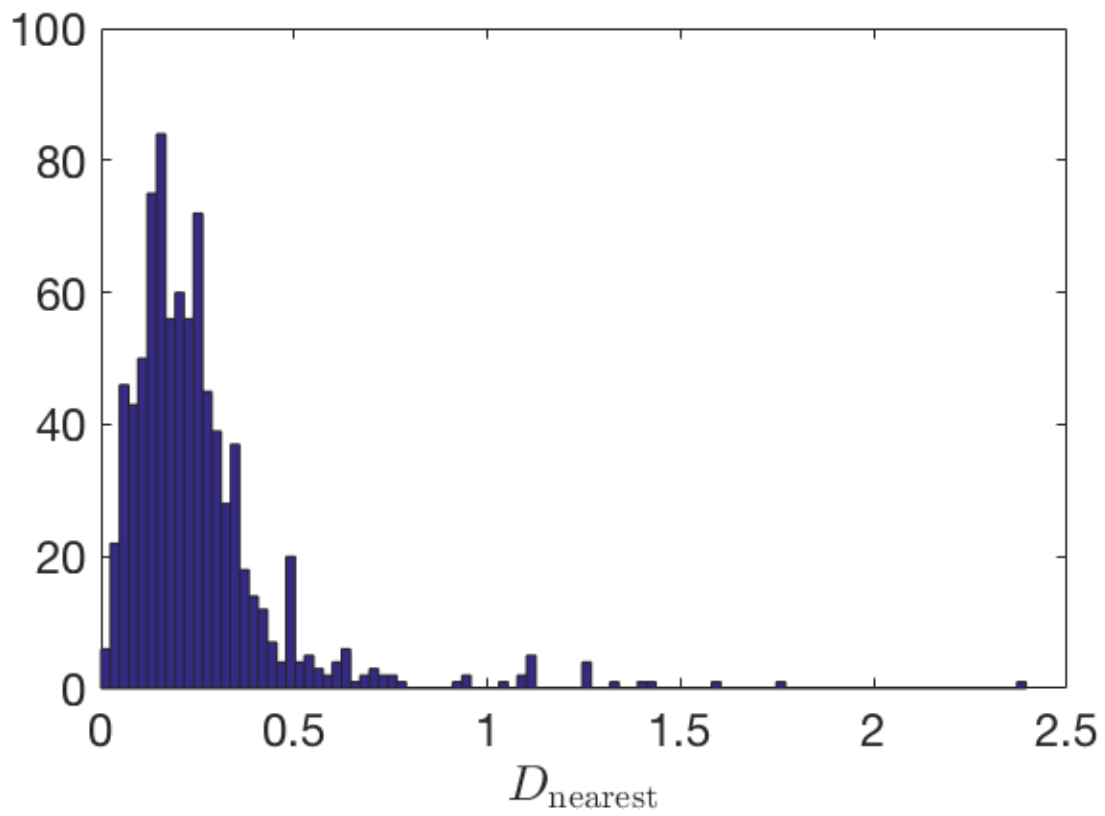
1051 Zhang, Y. G., Zhang, C. L., Liu, X.-L., Li, L., Hinrichs, K.-U., Noakes, J. E.: Methane Index: a tetraether
1052 archaeal lipid 15 biomarker indicator for detecting the instability of marine gas hydrates, *Earth and*
1053 *Planetary Science Letters*, 307, 525- 534, 2011.

1054 Zhang, Y.G., Pagani, M., Liu, Z.: A 12-million-year temperature history of the tropical Pacific
1055 Ocean: *Science*, 343, 84-86, 2014.

1056 Zhu, J., Poulsen, C. J., Tierney, J.: Simulation of Eocene extreme warmth and high climate sensitivity
1057 through cloud feedbacks, *Science Advances*, 5(9), eaax1874, 2019.

1058
1059
1060
1061
1062
1063
1064
1065

Figure 1

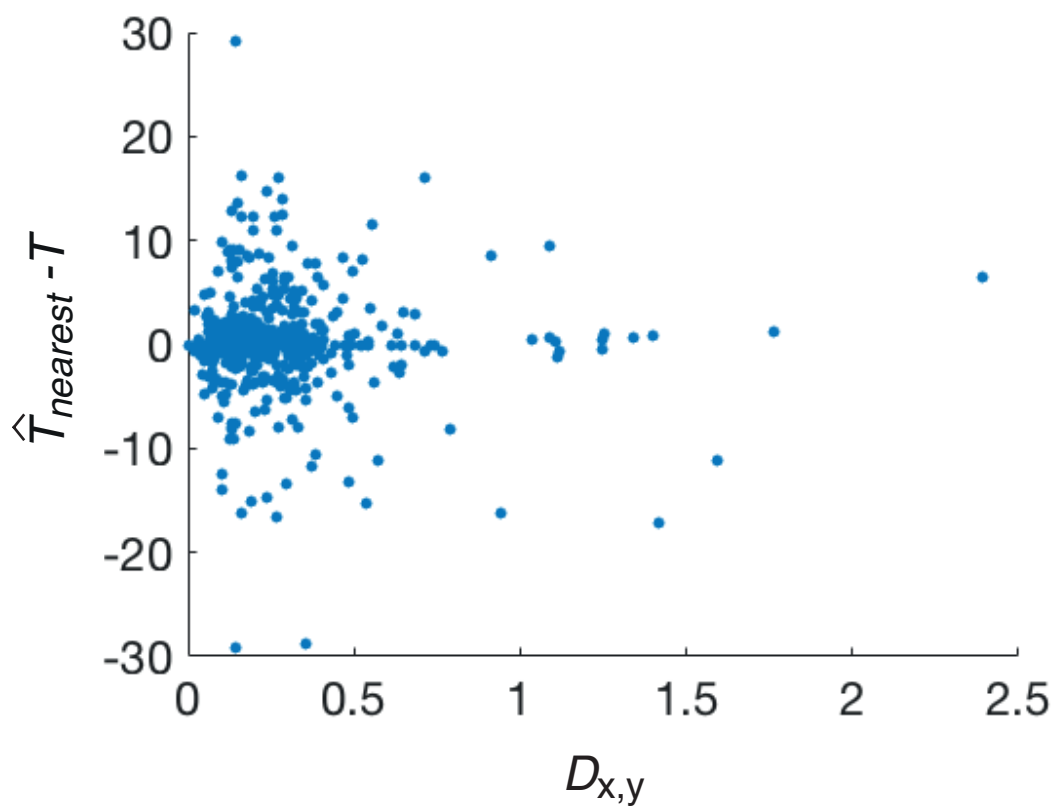


1066

1067

1068

Figure 2



1069

1070

1071

Figure 3

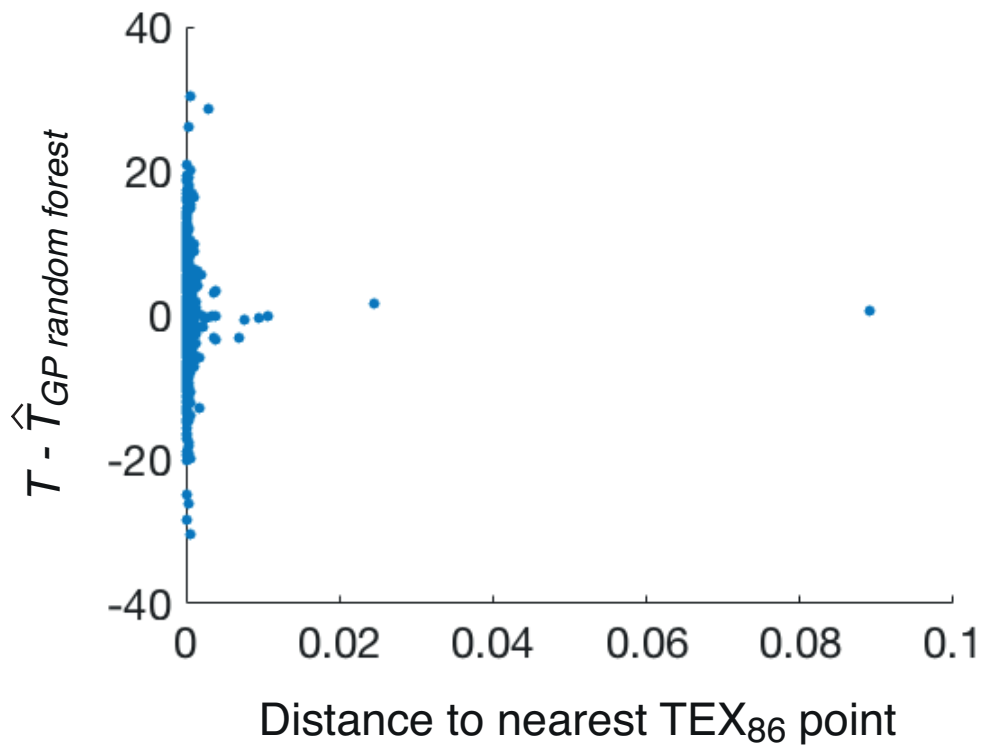
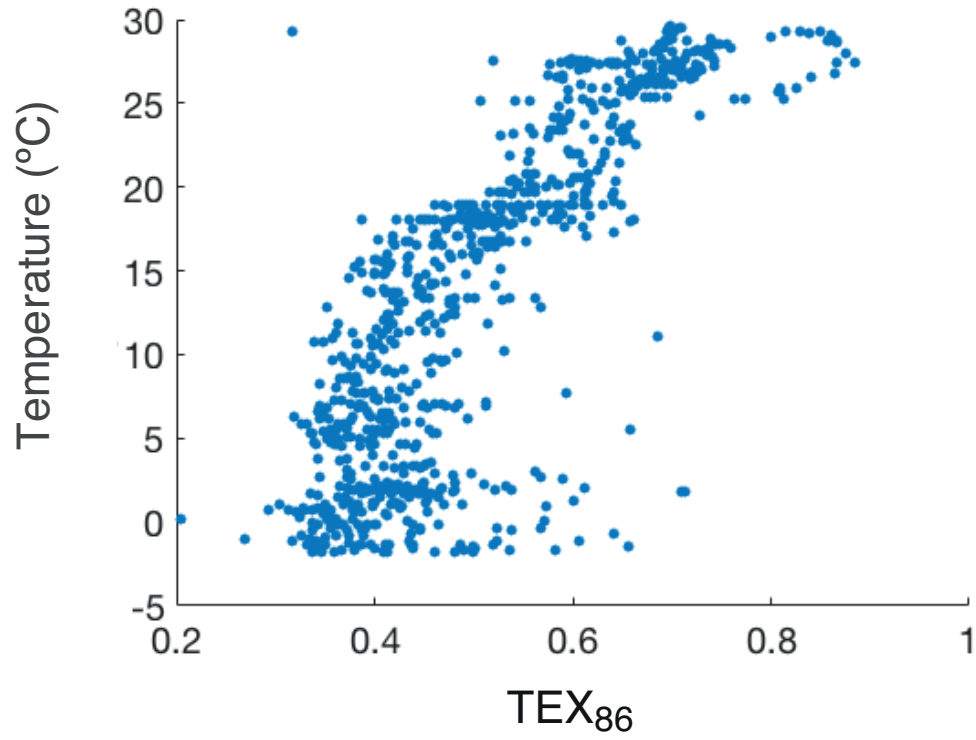
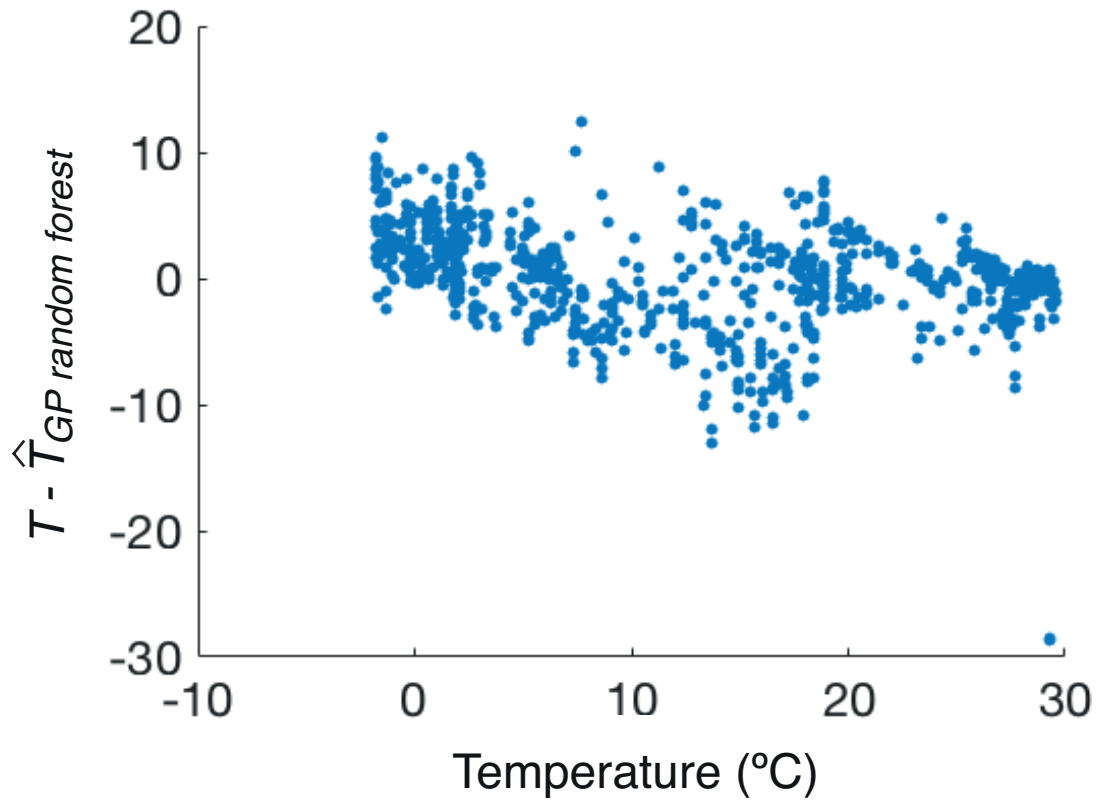


Figure 4



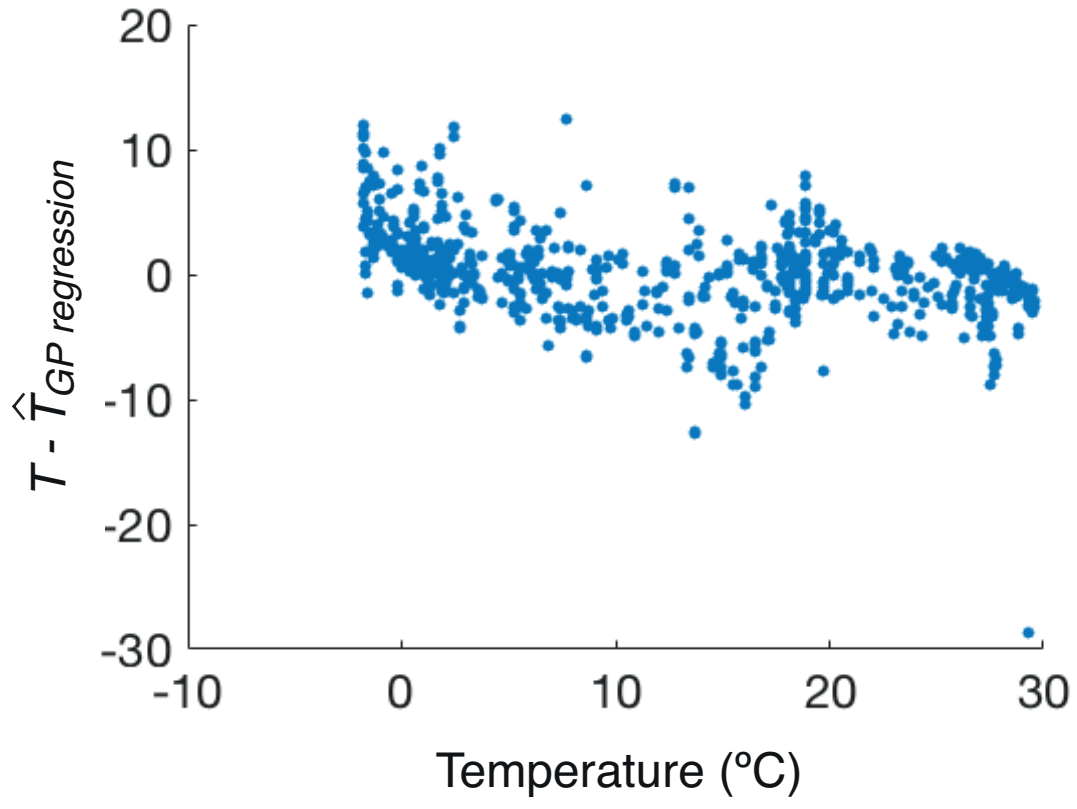
1073

1074

1075

1076

Figure 5

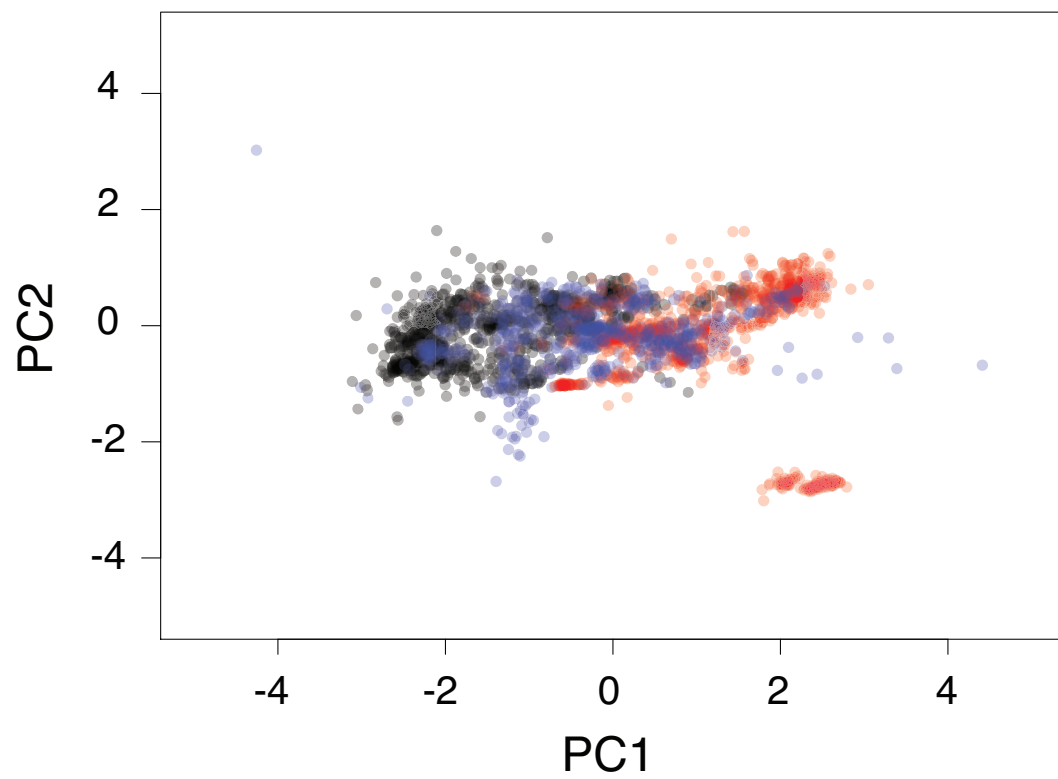


1077

1078

1079

Figure 6

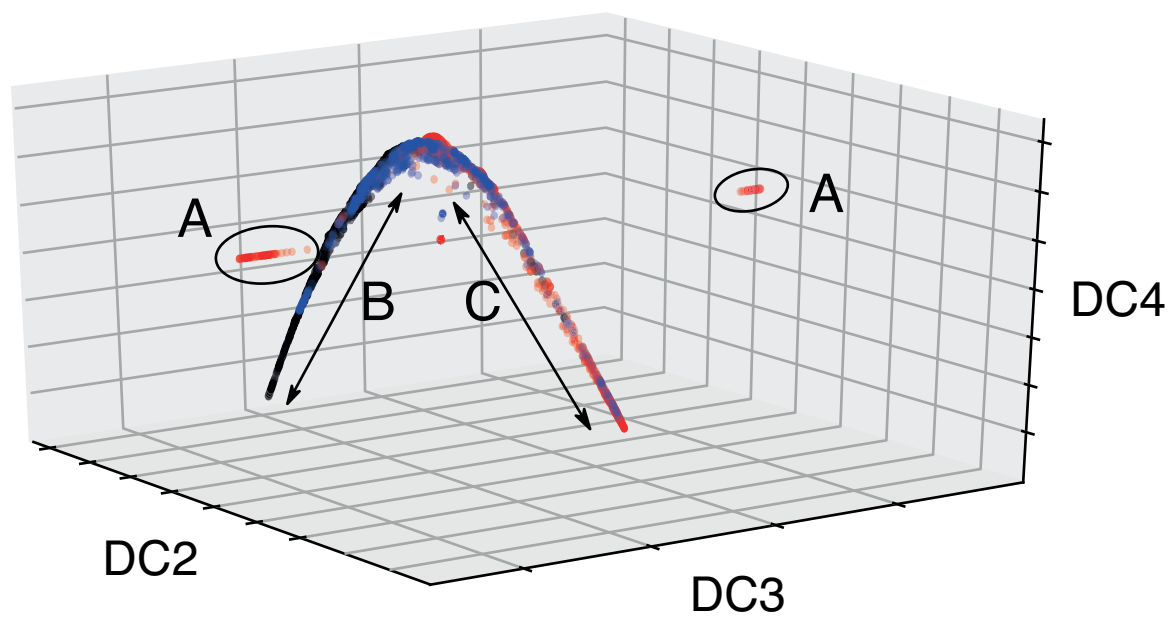


1080

1081

1082

Figure 7



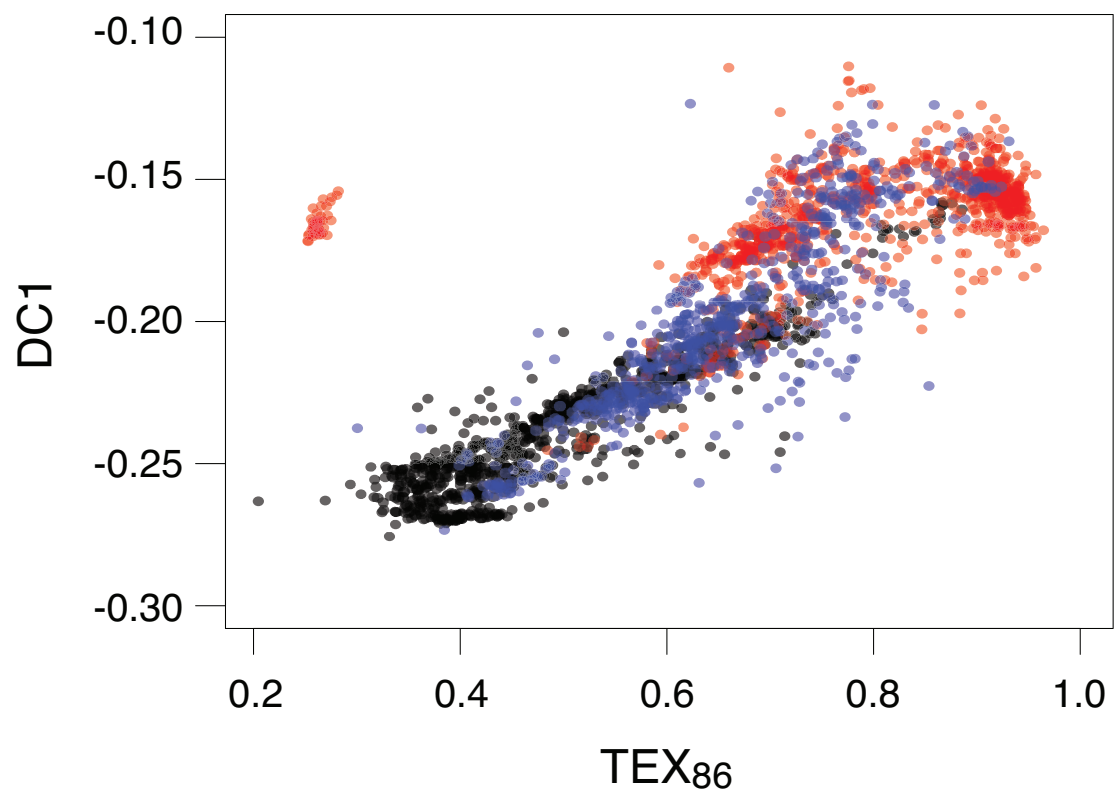
1083

1084

1085

1086

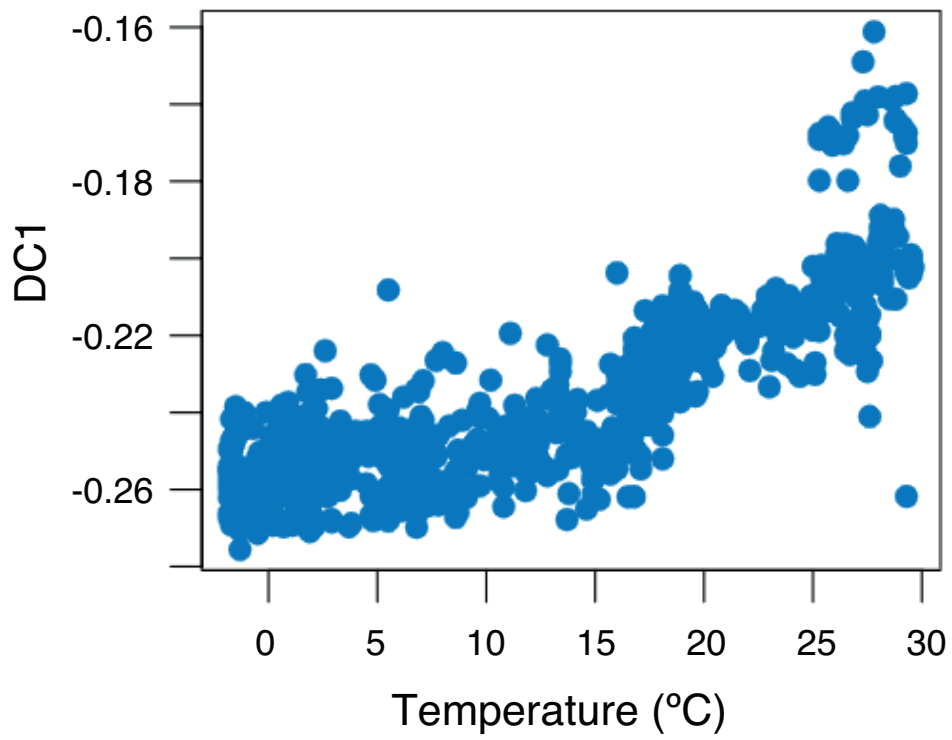
Figure 8



1087

1088

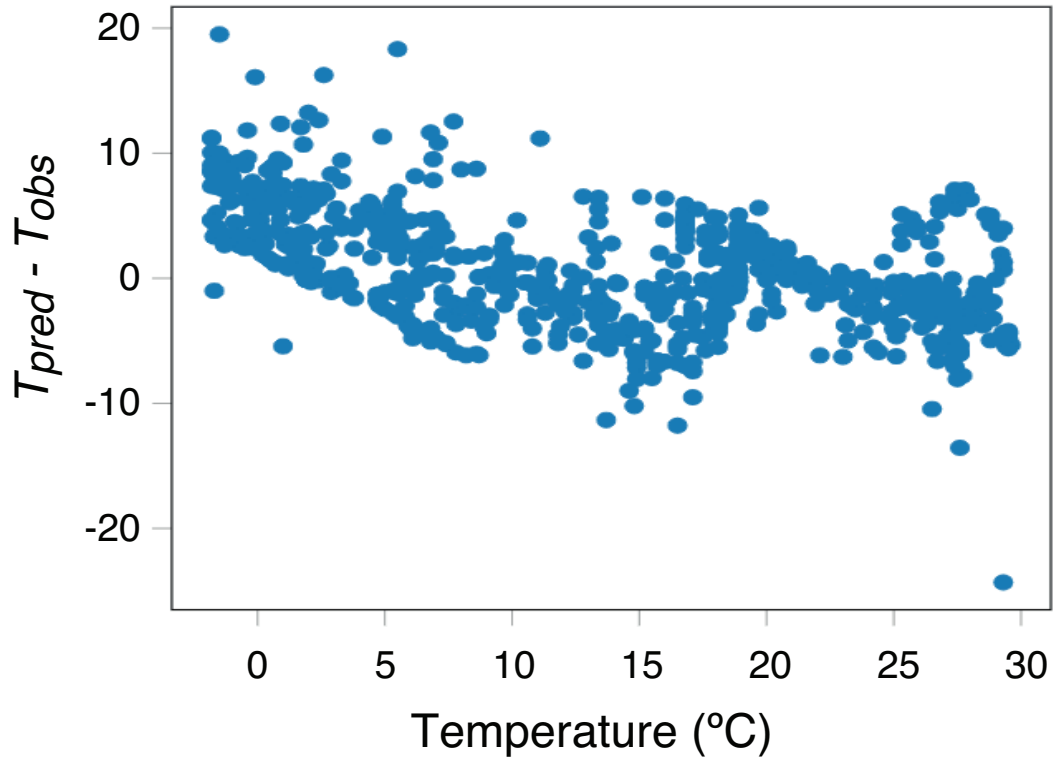
Figure 9



1089

1090

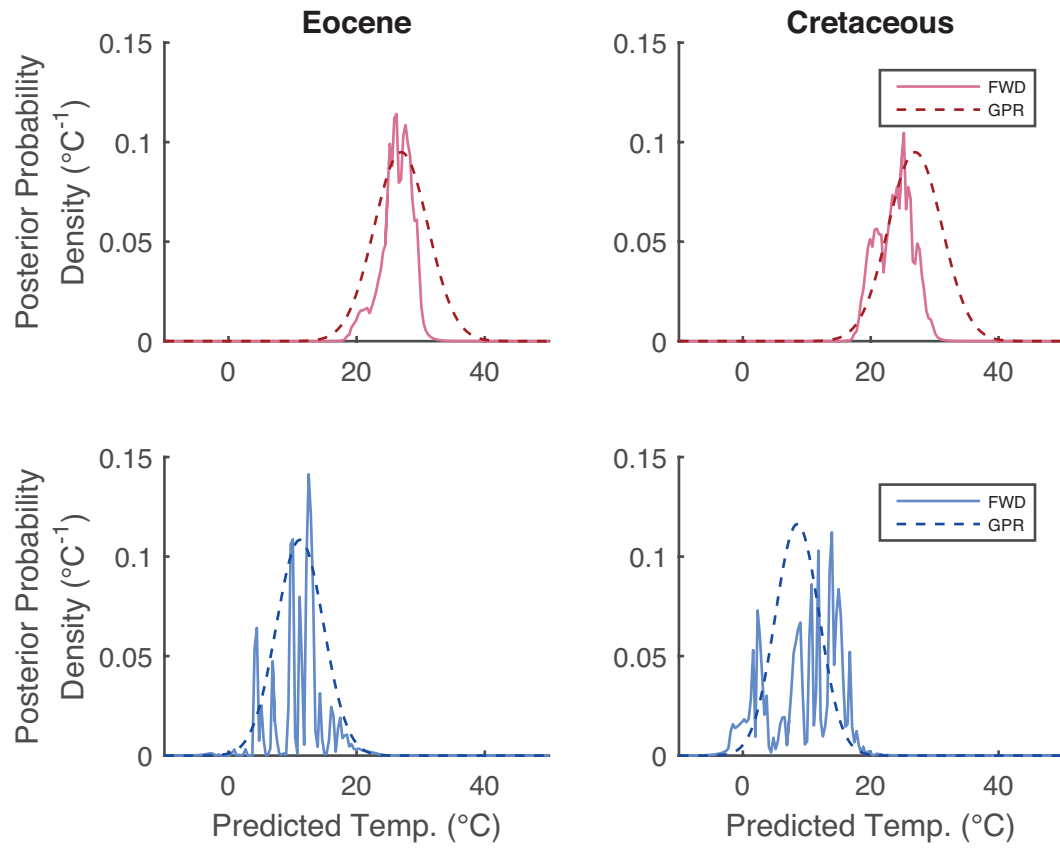
Figure 10



1091

1092

Figure 11

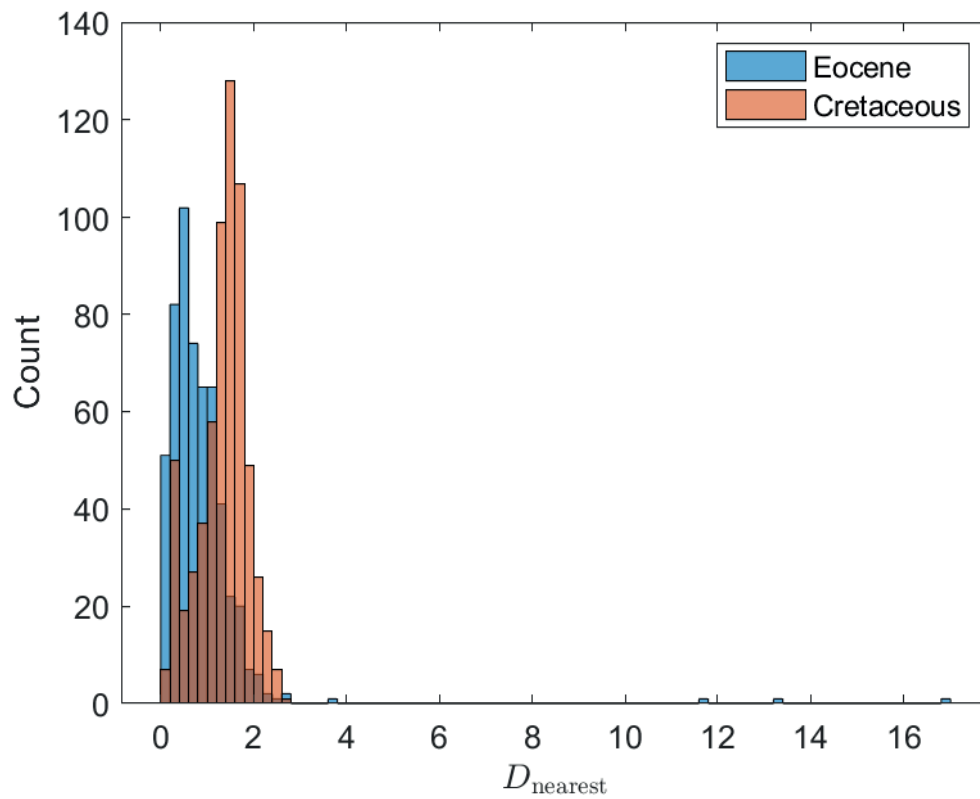


1093

1094

1095

Figure 11



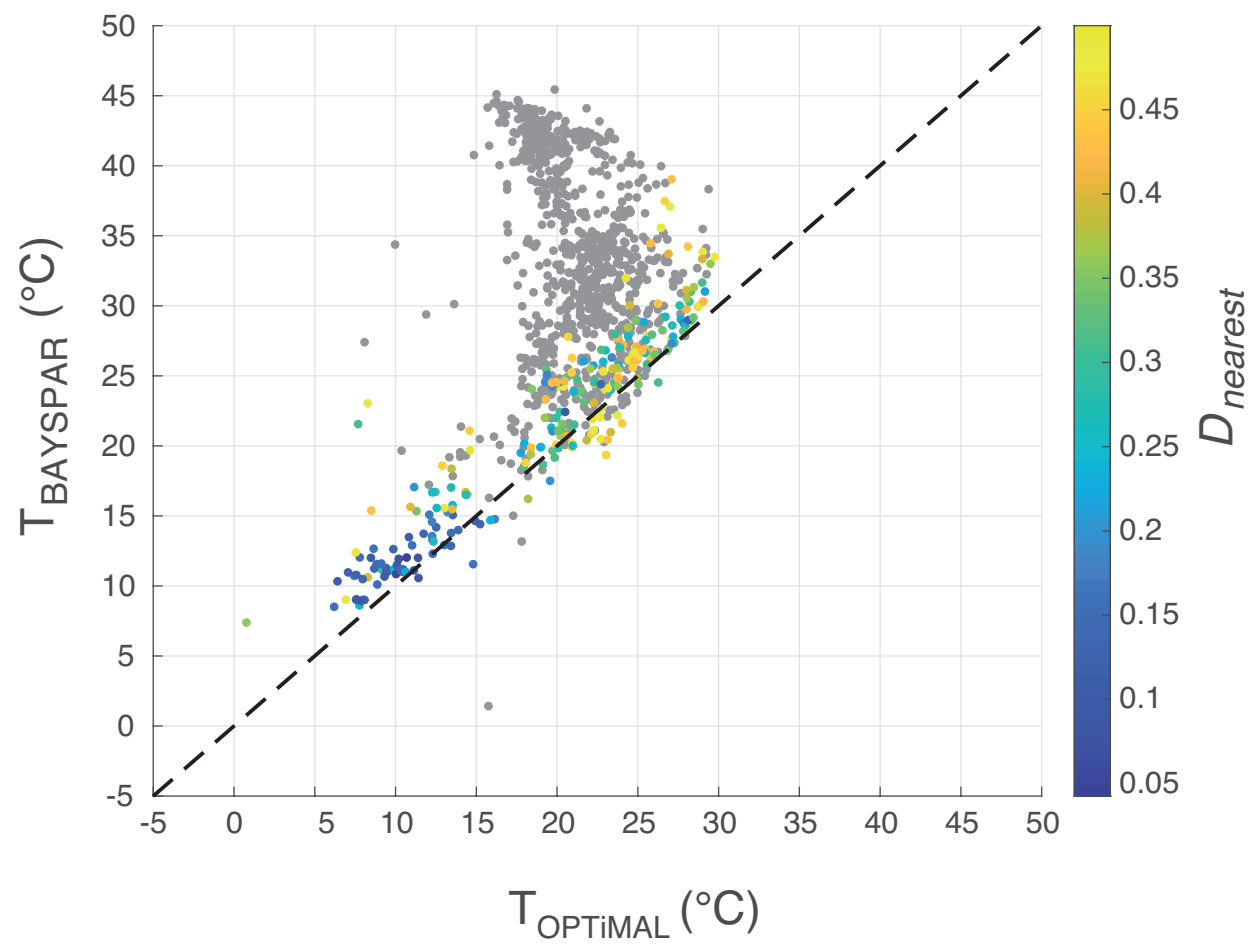
1096

1097

1098

1099

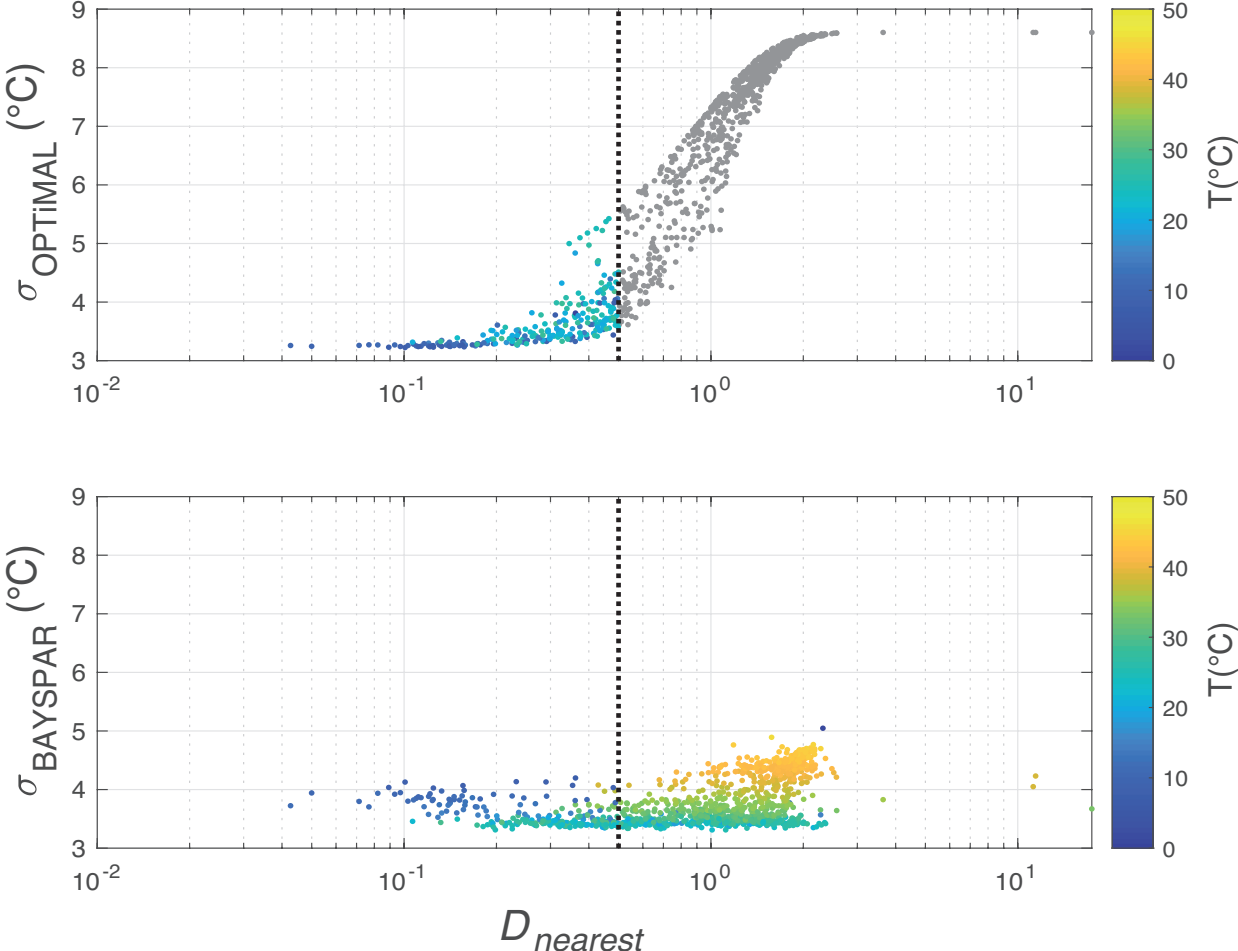
Figure 13



1100

1101

Figure 14



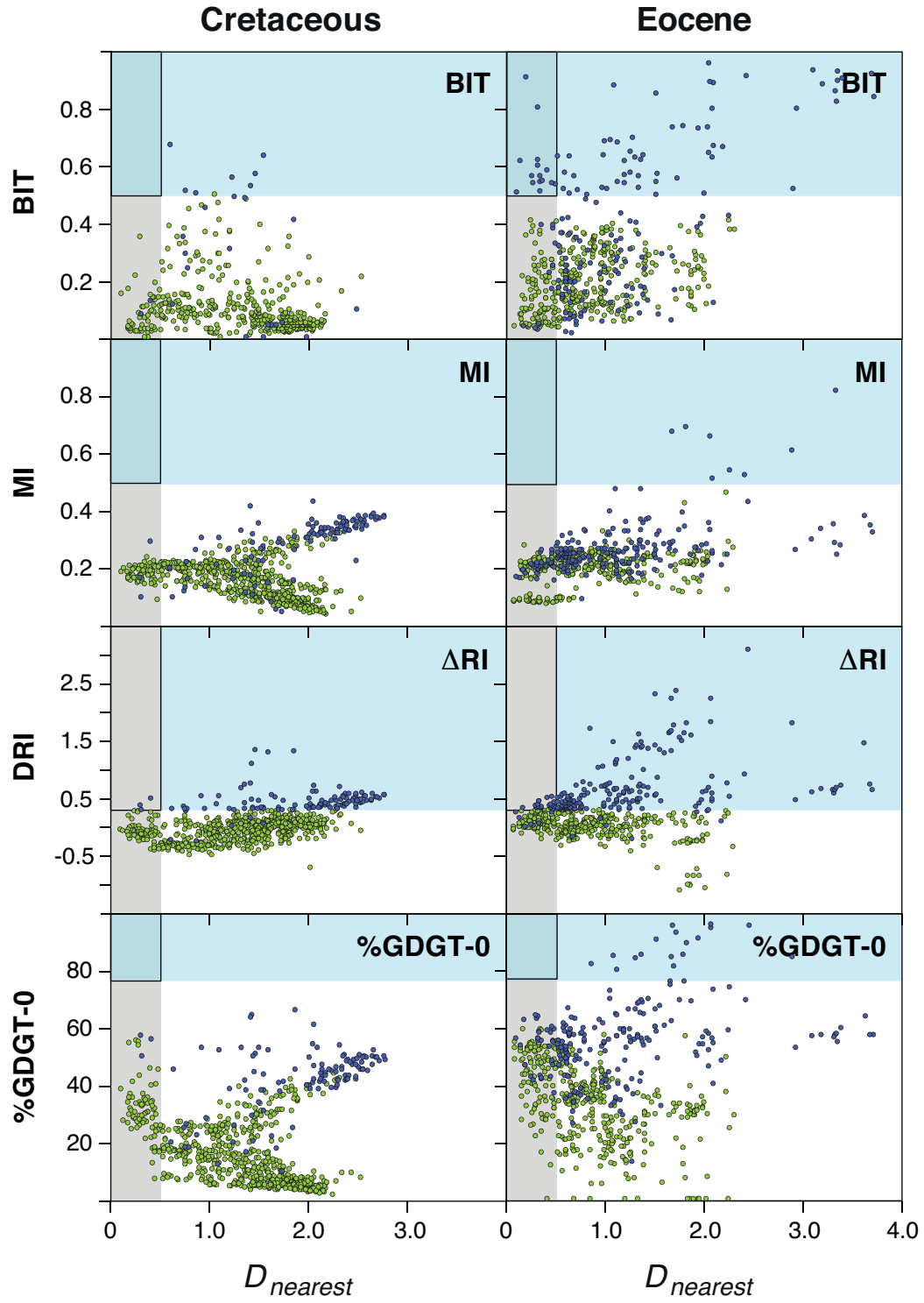
1102

1103

1104

1105

1106 Figure 15

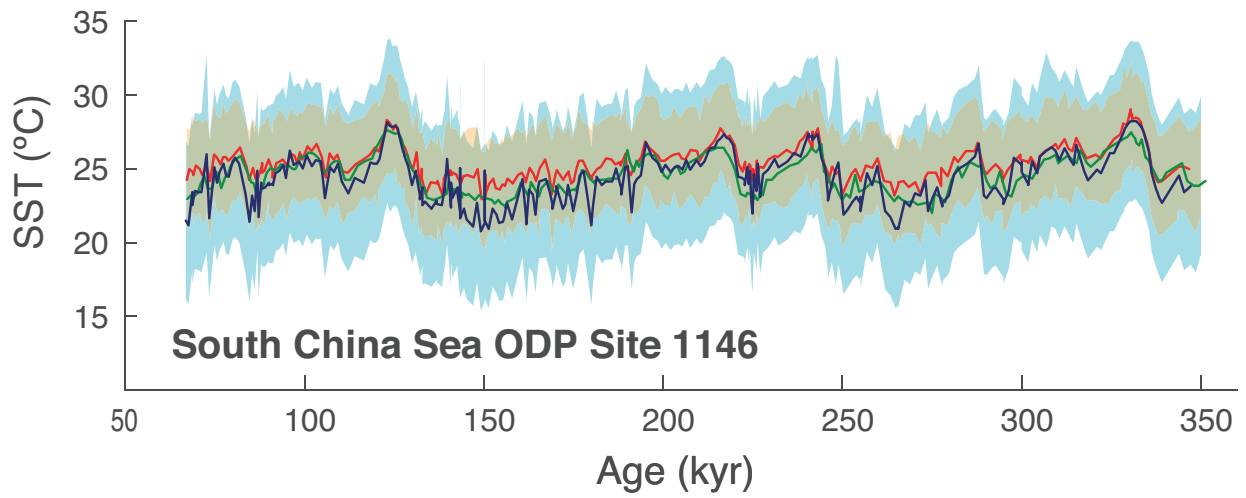
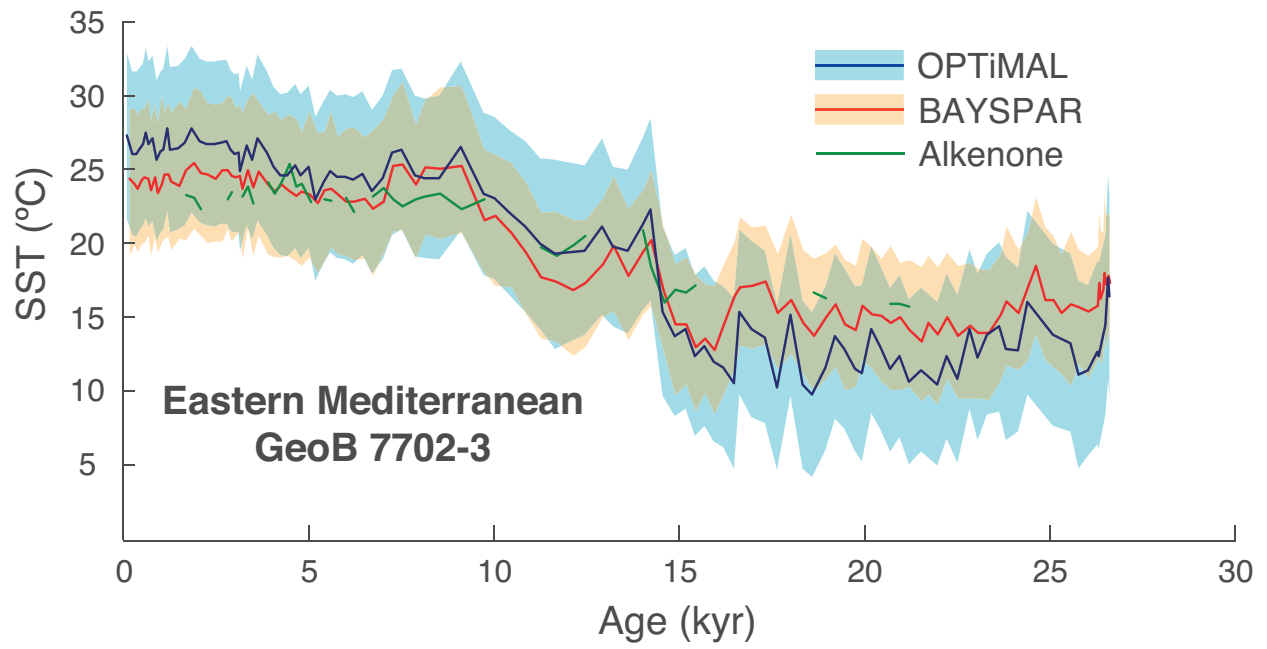


1107

1108

1109 **Figure 16**

1110

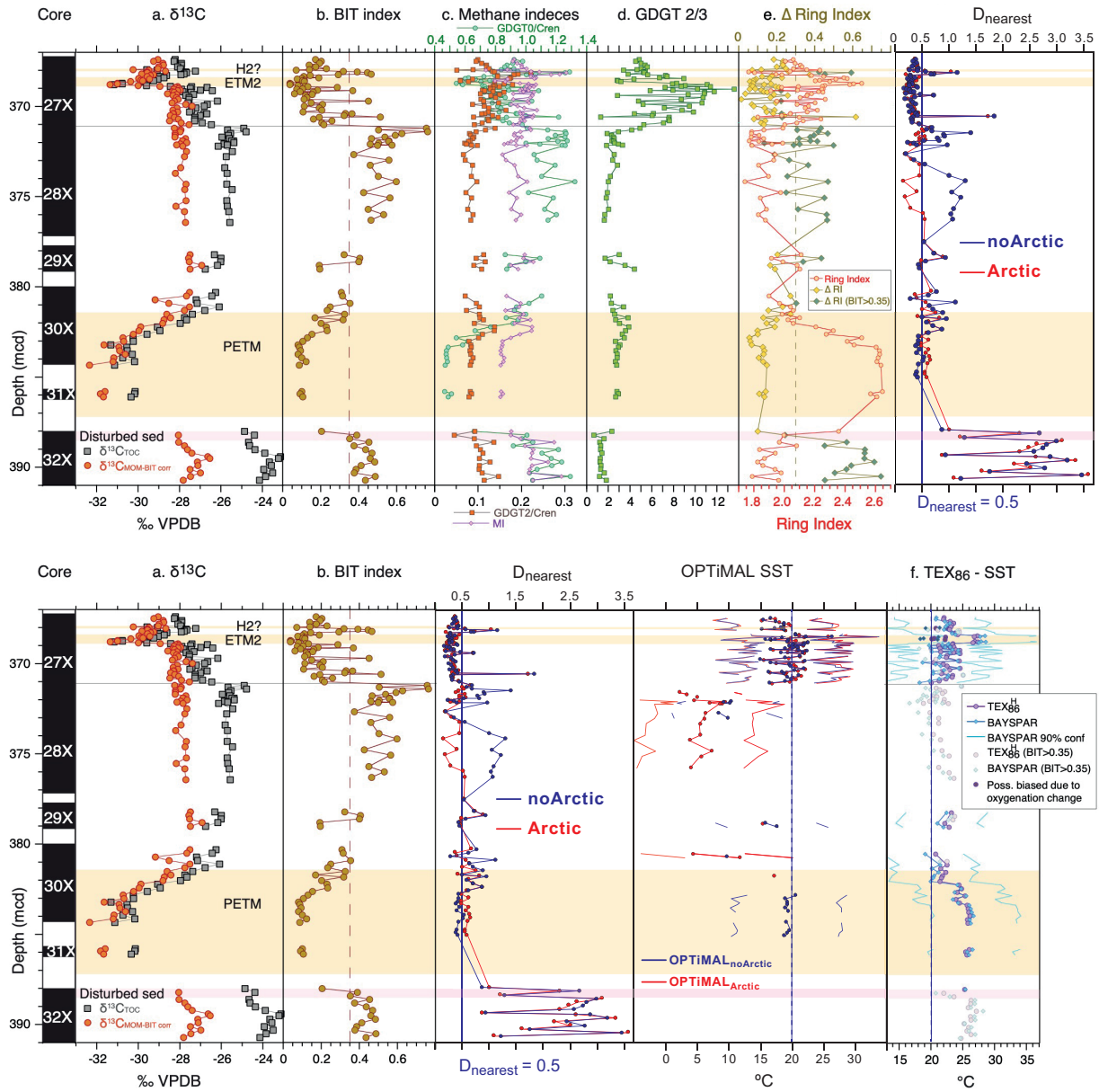


1111

1112

1113 **Figure 17**
 1114

Late Paleocene to early Eocene Arctic Ocean IODP Expedition 302 Hole 4A



1115