

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

2

3 Yvette L. Eley,¹ William Thomson², Sarah E. Greene¹, Ilya Mandel^{3,4}, Kirsty Edgar¹, James A. Bendle¹,
4 and Tom Dunkley Jones¹

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

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

14

15

16 **Abstract**

17

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

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

40

41 **1. Introduction**

42

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

54

$$55 \text{TEX}_{86} = (\text{GDGT-2} + \text{GDGT-3} + \text{cren}') / (\text{GDGT-1} + \text{GDGT-2} + \text{GDGT-3} + \text{cren}') \text{ (Eq. 1)}$$

56

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

67

$$68 \text{TEX}_{86} = 0.015 * \text{SST} + 0.287 \text{ (Eq. 2)}$$

69

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

$$81 \quad TEX_{86}^L = \log \left(\frac{[GDGT2]}{[GDGT1]+[GDGT2]+[GDGT3]} \right) \quad \text{Eq. 3}$$

$$84 \quad TEX_{86}^H = \log \left(\frac{[GDGT2]+[GDGT3]+[Cren']}{[GDGT1]+[GDGT2]+[GDGT3]+[Cren']} \right) \quad \text{Eq. 4}$$

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

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

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

113
114 Second, from the earliest applications of the TEX₈₆ proxy to deep-time warm climate states (Schouten et
115 al., 2003) it was recognized that reconstructed temperatures beyond the range of the modern calibration
116 (>30 °C), were highly sensitive to model choice within the modern calibration range. Thus, Schouten et al.
117 (2003) restricted their calibration data for deep-time temperature estimates to core-top data in the modern
118 with mean annual SSTs over 20 °C. However, this problem of model choice, and its impact on temperature
119 estimation beyond the modern calibration range, persists (Hollis et al. 2019), with current arguments
120 focused on whether there is an exponential (e.g. Cramwinckel et al., 2018) or linear (Tierney & Tingley,
121 2015) dependency of TEX₈₆ on SSTs, and the effect of these models on temperature estimates over 30 °C.

122
123 Culture and mesocosm studies are sometimes cited in support of extrapolations beyond the modern
124 calibration range when reconstructing ancient SSTs (Kim et al., 2010, Hollis et al., 2019). While there is a
125 basic underlying trend for more rings within GDGT structures at higher temperatures (Zhang et al. 2015;
126 Qin et al., 2015), the lack of a uniform response to archaeal GDGT production in response to increasing
127 growth temperatures (e.g., Elling et al., 2015; Qin et al., 2015) suggests that this does not easily translate
128 into a simple linear model at the community scale (i.e. the core top calibration dataset). Wuchter et al.
129 (2004) and Schouten et al. (2007) show a compiled linear calibration of TEX₈₆ against incubation
130 temperature (up to 40°C in the case of Schouten et al., 2007) based on strains that were enriched from
131 surface seawater collected from the North Sea and Indian Ocean respectively. Like Qin et al. (2015), we
132 note the *non*-linear nature of the individual experiments in Wuchter et al. (2004; see Fig. 5). Moreover, the
133 relatively lower Cren’ in these studies yield a very different intercept and slope compared to core-top
134 calibrations (e.g. Kim et al. 2010) making direct comparisons problematic.

135
136

137 More recently, Elling et al. (2015) studied three different strains (*N. maritimus*, NAOA6, NAOA2) isolated
138 from open ocean surface waters (South Atlantic) whilst Qin et al., (2015) studied a culture of *N. maritimus*
139 and three *N. maritimus*-like strains isolated from Puget Sound. All strains are of marine, mesophilic,
140 Thaumarchaeota within Marine Group 1 (equivalent to Crenarchaeota Group 1). Both of these papers
141 clearly demonstrate distinctly different responses of membrane lipid composition to temperature in these
142 strains, whilst Qin et al. (2015) additionally show that oxygen concentration is at least as important as
143 temperature in controlling TEX₈₆ values in culture. The impact of Thaumarchaeota community change on
144 TEX₈₆ in palaeoclimate studies is further suggested by the downcore study of Polik et al (2019). All of these
145 culture studies, made on marine, mesophilic archaea demonstrate how community composition may have
146 a significant impact on measured environmental TEX₈₆ signatures. In these cases (e.g., Zhang et al. 2015;
147 Qin et al., 2015; Elling et al., 2015) cultured strains of Thaumarchaeota were obtained from surface waters
148 which overlie the epi-continental or continental shelf regions of the North Sea, Indian Ocean, South Atlantic
149 and North Pacific - in addition to the pure culture strain *N. maritimus* in Qin et al. (2015) and Elling et al.
150 (2015). As such, these are collectively more representative of the community production contributing to
151 samples in the global core-top TEX₈₆ calibrations of Kim et al., (2010) and BAYSPAR (Tierney & Tingley,
152 2014), which predominantly sample continental margin environments, rather than deep ocean / pelagic
153 environments.

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

171
172 Finally, the original uses of the TEX₈₆ proxy had a relatively poor representation of the true uncertainty
173 associated with palaeotemperature estimates, as they included no assessment of non-analogue behavior
174 relative to the modern core-top data. Instead, uncertainty was typically based on the residuals on the modern
175 calibration, with no reference to the relationship between GDGT distributions of an ancient sample and the
176 modern calibration data. An improved Bayesian uncertainty model “BAYSPAR” is now in widespread use
177 for SST estimation, which models TEX₈₆ to SSTs regression parameters, and associated uncertainty, as
178 spatially varying functions (Tierney and Tingley, 2015). The Bayesian approach, as with all approaches
179 based on the TEX₈₆ index, however, still does not include an uncertainty that reflects how well modelled
180 ancient GDGT assemblages are by the modern calibration – i.e. the degree to which they are non-analogue
181 - as it still functions on one-dimensional TEX₈₆ index values.

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

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

204 datasets, which outperforms previous GDGT palaeothermometers and includes robust error estimates that,
205 for the first time, accounts for model uncertainty.

206

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

208

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

217

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

224

225

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

226

(Eq. 5)

227

228 where the sum is taken over each of $N_v = 85$ validation points, T is the known measured temperature (which
229 we refer to as the true temperature) and \hat{T} is the predicted temperature. For conciseness, we refer to δT as
230 the predictor standard error. It is useful to compare the accuracy of the predictor to the standard deviation
231 of all temperatures in the data set σT , which corresponds to using the mean temperature as the predictor in
232 Equation 1; for the modern data set, $\sigma T = 10.0$ °C. The coefficient of determination, R^2 , provides a measure
233 of the fraction of the fluctuation in the temperature explained by the predictor. To facilitate performance
234 comparisons between different methods of predicting temperature, we use the same subset of validation

235 points for all analyses. To avoid sensitivity to the choice of validation points, we repeat the calibration-
236 validation procedure for 10 random choices from the validation dataset.

237

238 *2.1 Nearest neighbours*

239

240 We begin with an agnostic approach to using some combination of the proportions of each of the six
241 observables - GDGT-0, GDGT-1, GDGT-2, GDGT-3, crenarchaeol and cren', which we will jointly refer
242 to as GDGTs - to predict sea surface temperatures. Whatever functional form the predictor might take, it
243 can only provide accurate temperature predictions if nearby points in the six-dimensional observable space
244 - i.e. the distribution of all of the six commonly reported GDGTs - can be translated to nearby points in
245 temperature space. Conversely, if nearby points in the observable space correspond to vastly different
246 temperatures, then no predictor, regardless of which combination of GDGTs are used, will be able to
247 provide a useful temperature estimate. In other words, the structuring of GDGT distributions within multi-
248 dimensional space, must have some correspondence to the temperatures of formation (or rather the mean
249 annual SSTs used for standard calibrations).

250

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

260

$$261 \quad D_{x,y}^2 \equiv \sum_{i=0}^6 \frac{GDGT_i(x) - (GDGT_i(y))^2}{var(GDGT_i)} \quad (Eq. 7)$$

262

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

265

266 The nearest-sample temperature predictor is $\hat{T}_{nearest}(x) = T(y)$ where y is the nearest point to x over the
267 calibration data set, i.e., one that minimises $D_{x,y}$. Fig. 2 shows the scatter in the predicted temperature when
268

269 using the temperature of the nearest data point to make the prediction. Overall, the failure of the nearest-
270 neighbour predictor to provide accurate temperature estimates even when the normalised distance to the
271 nearest point is small, $D_{x,y} \leq 0.5$, casts doubt on the possibility of designing an accurate predictor for
272 temperature based on GDGT observations. This is most likely due to additional environmental controls on
273 GDGT abundance distributions in natural systems, in particular the water depth (Zhang and Liu, 2018),
274 nutrient availability (Hurley et al., 2016; Polik et al., 2018; Park et al., 2018), seasonality, growth rate
275 (Elling et al., 2014; Hurley et al., 2016) and ecosystem composition (Polik et al., 2018), that obscure a
276 predominant relationship to mean annual SSTs.

277

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

286

287 *2.2 TEX₈₆ and Bayesian applications*

288

289 The TEX₈₆ index reduces the six-dimensional observable GDGT space to a single number. While this has
290 the advantage of convenience for manipulation and the derivation of simple analytic formulae for
291 predictors, as illustrated below, this approach has one critical disadvantage: it wastes significant information
292 embedded in the hard-earned GDGT distribution data. Fig. 3 illustrates both the advantage and
293 disadvantage of TEX₈₆. On the one hand, there is a clear correlation between TEX₈₆ and temperature (top
294 panel of Fig. 3), with a correlation coefficient of 0.81 corresponding to an overwhelming statistical
295 significance of 10^{-198} . On the other hand, very similar TEX₈₆ values can correspond to very different
296 temperatures. We can apply the nearest-neighbour temperature prediction approach to the TEX₈₆ value
297 alone rather than the full GDGT parameter space; this predictor yields a large standard error of $\delta T_{\text{nearestTEX86}}$
298 $= 8.0 \text{ }^\circ\text{C}$ (bottom panel of Fig. 3). While smaller than σT , this is significantly larger than $\delta T_{\text{nearest}}$ (Fig. 2),
299 consistent with the loss of information in TEX₈₆. We therefore do not expect other predictors based on
300 TEX₈₆ to perform as well as those based on the full available data set.

301

302 Indeed, this is what we find when we consider predictors of the form $\hat{T}_{1/\text{TEX}} = a + b/\text{TEX}_{86}$ and $\hat{T}_{\text{TEXH}} = c$
303 $+ d \log_{\text{TEX}_{86}}$ (Liu et al., 2009; Kim et al., 2010), i.e., the established relationships between GDGT
304 distributions and SST. We fit the free parameters a , b , c , and d by minimising the sum of squares of the
305 residuals over the calibration data sets (least squares regression). We find that $\delta T_{1/\text{TEX}} = 6.1$ °C (note that
306 this is slightly better than using the fixed values of a and b from (Kim et al., 2010), which yield $\delta T_{1/\text{TEX}} =$
307 6.2 °C). We note that the corresponding R^2 value associated with these TEX_{86} based predictors is 0.64,
308 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
309 dataset based on Tierney and Tingley (2015), including data from the Red Sea (Kim et al. 2010).

310

311 Tierney and Tingley (2014) proposed a more sophisticated approach to obtaining the transfer function from
312 TEX_{86} to temperature, continuing to use simple linear regression, but with the addition of Gaussian
313 processes to model spatial variability in the temperature- TEX_{86} relationship and working with a forward
314 model which is subsequently inverted to produce temperature predictions. This forward model
315 ‘BAYSPAR’ is capable of generating an infinite number of calibration curves relating TEX_{86} to sea surface
316 temperatures (Tierney and Tingley, 2014). In order to derive a calibration for a specific dataset, the user
317 edits a range of parameters which vary depending on whether the dataset in question is from the relatively
318 recent past or deep time (Tierney and Tingley, 2014). For deep time applications, the authors propose a
319 modern analogue-type approach, in which they search the modern data for $20^\circ \times 20^\circ$ grid boxes containing
320 ‘nearby’ TEX_{86} measurements and subsequently apply linear regression models calibrated on the analogous
321 samples for making predictions.

322

323 However, along with the simpler TEX_{86} -based models described above, this approach still suffers from the
324 reduction of a six-dimensional data set to a single number. Therefore, it is not surprising that even the
325 simplest nearest-neighbour predictor (such as the one described above) that makes use of the full six-
326 dimensional dataset outperforms single-dimensional forward modelling approaches. Additionally,
327 uncertainty estimates do not account for the fact that TEX_{86} is, fundamentally, an empirical proxy, and so
328 its validity outside the range of the modern calibration is not guaranteed. This is a fundamental issue for
329 attempts to reconstruct surface temperatures during Greenhouse climate states, when tropical and sub-
330 tropical SSTs were likely hotter than those observed in the modern oceans.

331

332 *2.3 Machine learning Approaches – Random Forests*

333

334 There are a number of options to improve on nearest-neighbour predictions using machine learning
335 techniques such as artificial neural networks and random forests. These flexible, non-parametric models

336 would ideally be based on the underlying processes driving the GDGT response to temperature, but since
337 these processes remain unconstrained at present, we choose to deploy models which can reasonably reflect
338 predictive uncertainty and will be sufficiently adaptable in future (as new information regarding controls
339 on GDGTs emerge). These machine learning approaches are all based on the idea of training a predictor by
340 fitting a set of coefficients in a sufficiently complex multi-layer model in order to minimise residuals on
341 the calibration data set. As an example of the power of this approach, we train a random forest of decision
342 trees with 100 learning cycles using a least-squares boosting to fit the regression ensemble. Figure 4 shows
343 the prediction accuracy for this random forest implementation. This machine learning predictor yields δT
344 = 4.1 °C degrees, outperforming the naive nearest-neighbour predictor by effectively applying a suitable
345 weighted average over multiple near neighbours. This corresponds to a very respectable $R^2 = 0.83$, meaning
346 that 83% of the variation in the observed temperature is successfully explained by our GDGT-based model.

347

348 *2.4 Gaussian Process Regression*

349

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

361

362 The weighting coefficients learned by the GP emulator represent a covariance matrix on the GDGT
363 parameter space. We can use this as a distance metric to provide meaningfully normalised distances
364 between points, removing the arbitrariness from the nearest neighbour distance ($D_{x,y}$) definition used
365 earlier. If the temperature is insensitive to a particular GDGT input coordinate (i.e., the value of that input
366 has a minimal effect on the temperature) then points within GDGT space that have large differences in
367 absolute input values in that coordinate are still near. We find that Cren has very limited predictive power,
368 and so points with large Cren differences are close in term of the normalised distance. Conversely, if the
369 temperature is sensitive to small changes in a particular GDGT variant, then points with relatively nearby

370 absolute input values in that coordinate are still distant. We find that most GDGT parameters other than
371 Cren are comparably useful in predicting temperature, with GDGT-0 and GDGT-3 marginally the most
372 informative. We considered whether interdependency of percentage GDGT data could influence our
373 calculations. Our analysis suggests that there are only five free parameters. Machine learning tools should
374 be able to pick up this correlation and effectively ignore one of the parameters (or one parameter
375 combination). For example, we do find that the GP emulator has a very broad kernel in at least one
376 dimension, signaling this. In principle, we could have considered only five of six parameters. The smaller
377 scale of some of the parameters is automatically accounted for by the trained kernel size in GP regression,
378 or by normalising to the appropriate dynamical range in our initial investigation.

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

391 392 *2.5 Data Structure*

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

404 approximately on a one-dimensional curve in GDGT space, and the arclength along this curve should
405 correspond to temperature; we will revisit this point below.

406
407 We first employed a version of principal component analysis (PCA) tailored to compositional data
408 (Aitchison, 1982, 1983; Aitchison and Greenacre, 2002; Filzmoser et al., 2009a; Filzmoser et al., 2009b;
409 Filzmoser et al., 2012). Taking into account the compositional nature of the data is important because the
410 sum-to-one constraint induces correlations between variables which are not accounted for by classical PCA.
411 Furthermore, apparently nonlinear structure in Euclidean space often corresponds to linearity in the simplex
412 (i.e. the restricted space in which all elements sum to one) (Egozcue et al., 2003). Figure 6 shows the
413 modern, Eocene and Cretaceous data projected onto the first two principal components. Aside from the
414 obvious outlying cluster of Cretaceous data, characterised by GDGT-3 fractions above 0.6, the bulk of the
415 data occupy a two-dimensional point cloud with a small amount of curvature. The large majority of the
416 Cretaceous data has more positive PC1 values relative to the modern data.

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

430
431 The relationship between the first diffusion component and TEX_{86} for all data is shown in Figure 8. There
432 is a clear correlation, despite the presence of some outlying Cretaceous points, some of which are not shown
433 because they lie so far outside the majority data range within this projection. This suggests that TEX_{86} is,
434 in one sense, a natural one-dimensional representation of the data. We also plot the first diffusion
435 component for the modern data as a function of temperature (Figure 9). We see a similar pattern emerging
436 to that displayed by TEX_{86} - there is little sensitivity to temperature below 15 °C, and between ~20 and 25
437 °C. An interesting avenue for future research might be to explore the temperature-GDGT system from a

438 dynamical systems perspective, i.e. use simple mechanistic mathematical models to explore the
439 temperature-dependence of steady-state GDGT distributions. It may be that such models suggest that only
440 a few steady-states exist, and that temperature is a bifurcation parameter, i.e. it controls the switch between
441 the steady states. Note also the downward slope in the residual pattern in Figure 4 between 0 and 15-17
442 degrees celsius, and again at higher temperatures. This pattern is consistent with predictions that are biased
443 towards the centre of each 'cluster', i.e. a system which is not very sensitive to temperature, but can
444 distinguish between high and low temperatures reasonably well. This observation also links to recent culture
445 studies (Elling et al., 2015) and Pliocene-Pleistocene sapropel data (Polik et al., 2018), which support the
446 existence of discrete populations with unique GDGT-temperature relationships and that temporal changes
447 in population over time can drive changes in TEX_{86} .

448

449 *2.6 Forward Modelling*

450

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

470

471 The forward modelling approach is similar to that of Haslett et al. (2006), who argue that it is preferable to
472 model measured compositions as functions of climate, before probabilistically inverting the model to infer
473 plausible climates given a composition. The cost of modelling the data in this more natural way is the loss
474 of degrees of freedom -- we are now attempting to fit a one-dimensional line through a multidimensional
475 point cloud rather than fit a multidimensional surface to the GDGT data, which means that the predictive
476 power of the model suffers, at least on the modern data. The existing BAYSPAR calibration also specifies
477 the model in the forward direction, however while BAYSPAR does model spatial variability, it does not
478 account for the systematic uncertainty in the model when extrapolated beyond the calibration range. As
479 with all GP models, the choice of kernel has a substantial impact on predictions (and their associated
480 uncertainty) outside the range of the modern data, where predictions revert to the prior implied by the
481 kernel. Given that we have no mechanistic model for the data generating process, we recommend the use
482 of kernels which do not impose strong prior assumptions on the form of the GDGT-temperature relationship
483 (e.g. kernels with a linear component) and thus reasonably represent model uncertainty outside the range
484 of the modern data. We choose a zero-mean Matern 3/2 kernel for the applications below. Note, however,
485 that since we are working in ilr-transformed coordinates, this corresponds to a prior assumption of uniform
486 compositions at all temperatures, i.e. all components are equally abundant.

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

495
496 An advantage of the forward modelling approach is that the inversion can incorporate substantive prior
497 information about temperatures for individual data points. In particular, other proxy systems can be used to
498 elicit prior distributions over temperatures to constrain GDGT-based predictions, particularly when
499 attempting to reconstruct ancient climates with no modern analogue in GDGT-space. We emphasise that
500 outside the range of the modern data, the utility of the models is almost solely due to the prior information
501 included in the reconstruction. At present, the only priors being used in the forward model prescribe a
502 reasonable upper limit and lower limit on temperatures (see Supplementary Information). The only way to
503 improve these reconstructions will be for future iterations to incorporate prior information from other
504 proxies. It is worth noting that the predictive uncertainty, while reasonably well-described by the standard

505 deviation in cases where ancient data lie quite close to the modern data in GDGT space, can be highly
506 multimodal (Fig. 11). This is the case when estimates are significantly outside of the modern calibration
507 dataset, such as low latitude data in the Cretaceous, or where there is considerable scatter in the modern
508 calibration data, for example in the low temperature range (<5 °C).

509

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

511

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

535

536 Where ancient GDGT distributions lie far from the modern calibration data set ($D_{\text{nearest}} > 0.5$), we argue that
537 there is no suitable set of modern analogue GDGT distributions from which to infer growth temperatures
538 for this ancient GDGT distribution. Both the GPR and Fwd models revert to imposed priors once the

539 distance from the modern calibration dataset increases. We propose that this is more rigorous and justified
540 model behavior than extrapolation of TEX₈₆ or BAYSPAR predictors to non-analogue samples far from
541 the modern calibration data. As a result, the predictive models can only be applied to a subset of the Eocene
542 and Cretaceous data. We also note that there are two broad, non-mutually-exclusive categories of samples
543 that lie far from the modern calibration dataset ($D_{\text{nearest}} > 0.5$), the first are samples that seem to lie ‘beyond’
544 the temperature-GDGT calibration relationship, likely with (unconstrained) GDGT formation temperatures
545 higher than the modern core-top calibrations; the second are samples with anomalous GDGT distributions
546 lying on the margins of, or far away from the main GDGT clustering in 6-dimensional space (see outliers
547 in Fig. 8).

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

573 than temperature. In view of these existing data, we see no robust justification at present for the
574 extrapolation of modern core-top calibration data sets into the unknown above 30 °C, although the
575 coherent patterns apparent across GDGT space, between modern, Eocene and Cretaceous data (Figures
576 7), do provide some grounds for hope that the extension of GDGT palaeothermometry beyond 30°C might
577 be possible in future.

578

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

580

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

589

590 To investigate the behaviour of the new OPTiMAL model, we compare temperature predictions including
591 uncertainties for the Eocene and Cretaceous datasets, made by OPTiMAL and the BAYSPAR methodology
592 of Tierney and Tingley (2014) (Figures 13 and 14), using the default priors specified in the model code for
593 the BAYSPAR estimation. The OPTiMAL model systematically estimates slightly cooler temperatures
594 than BAYSPAR, with the biggest offsets below ~15 °C (Figure 13). We suggest that this is driven by our
595 inclusion of all data present in our modern calibration dataset, including data north of 70°N, which is
596 excluded from the BAYSPAR model (Tierney and Tingley, 2015). Fossil GDGT assemblages that fail the
597 D_{nearest} test are shown in grey, which clearly illustrate the regression to the mean in the OPTiMAL model,
598 whereas BAYSPAR continues to make SST predictions up to and exceeding 40 °C for these “non-analogue”
599 samples due to the fact that BAYSPAR assumes that higher TEX86 values equate to higher temperatures
600 as part of the functional form of the model, whereas the GPR model is agnostic on this. A comparison of
601 error estimation between OPTiMAL and BAYSPAR is shown in Figure 14. For most of the predictive
602 range below the D_{nearest} cut-off of 0.5, OPTiMAL has smaller errors than BAYSPAR, especially in the lower
603 temperature range. As D_{nearest} increases, i.e. as the fossil GDGT assemblage moves further from the
604 constraints of the modern calibration dataset, the error on OPTiMAL increases, until it reaches the standard
605 deviation of the modern calibration dataset (i.e., is completely unconstrained). In other words, OPTiMAL
606 generates maximum likelihood SSTs with robust confidence intervals, which appropriately reflect the

607 relative position of an ancient sample used for SST estimation and the structure of the modern calibration
608 data set. Where there are strong constraints from near analogues in the modern data, uncertainties will be
609 small, where there are weak constraints, uncertainty increases. In contrast, BAYSPAR, because it is
610 fundamentally based on a *parametric* linear model and therefore does not account for model uncertainty,
611 assigns similar uncertainty intervals as to the rest of the data, despite there being no way of reasonably
612 testing whether the linear model is an appropriate description of the data far from the modern dataset.

613
614 We further provide a first assessment of the inter-relationship between standard screening indices and
615 D_{nearest} , with an additional figure (Figure 15) showing the relationship between D_{nearest} , BIT, and MI for the
616 Eocene and Cretaceous compilations (where these data are available). These plots show little relationship
617 between the BIT and MI screening indices and D_{nearest} values. Whilst in the Eocene, samples with the highest
618 D_{nearest} values (>3) also show very elevated BIT values (>0.8), in the Cretaceous the exceptionally
619 anomalous assemblages (D_{nearest} values >100) are not anomalous in either BIT or MI. Conversely, in the
620 Eocene there are many samples with relatively high BIT (>0.3) that are below the D_{nearest} threshold of 0.5.
621 The behaviour of these systems needs to be examined in detail in future studies, but a conservative approach
622 would be to apply all three screening indices (BIT, MI and D_{nearest}) to have the most confidence in resulting
623 temperature estimates. To investigate these behaviours requires the publication of the full range GDGT
624 abundance data. Whilst key compilations of Eocene and Cretaceous GDGT data have strongly encouraged
625 the release of such datasets (Lunt et al. 2012; Dunkley Jones et al. 2013; Inglis et al. 2015; O'Brien et al.
626 2017), most Neogene studies only publish TEX_{86} values. Without full GDGT assemblage data neither
627 OPTiMAL nor other detailed assessments of GDGT behaviour and type can be made, and we would
628 strongly encourage authors, reviewers and editors to ensure the publication of full GDGT assemblages in
629 future.

630 Finally, we provide two example time series from the Neogene to modern, where full GDGT assemblage
631 data were made available, and there are comparison alkenone-based U^k_{37} data from the same sampling
632 location – ODP 806 and ODP 850, respectively in the West and Eastern Equatorial Pacific (Figure 16;
633 Zhang et al. 2014). Where U^k_{37} temperatures are not at the limit of alkenone saturation in ODP 806,
634 OPTiMAL and U^k_{37} agree well in the Plio-Pleistocene. In ODP 850, there is a strong agreement in the
635 reproduction of a long-term late Miocene to Recent cooling trend in both U^k_{37} and OPTiMAL of $\sim 5^\circ\text{C}$.
636 There is, however a consistent $\sim 2^\circ$ offset between cooler OPTiMAL and warmer U^k_{37} temperatures at this
637 location. This offset is very similar in magnitude and direction as that between $\text{TEX}_{86}^{\text{H}}$ and U^k_{37} used in
638 the original study. and is more likely due to an inherent feature (seasonality or depth) of archaeal versus
639 eukaryotic production at this site (Zhang et al. 2014).

640

641 **5. Conclusions**

642

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

657

658 In this study, we apply modern machine-learning tools, including Gaussian Process Emulators and forward
659 modelling, to improve temperature estimation and the representation of uncertainty in GDGT-based SST
660 reconstructions. Using our new nearest neighbour test, we demonstrate that >60% of Eocene, and >90% of
661 Cretaceous, fossil GDGT distribution patterns differ so significantly from modern as to call into question
662 SSTs derived from these assemblages. For data that does show sufficient similarity to modern, we present
663 OPTiMAL, a new multi-dimensional Gaussian Process Regression tool which uses all six GDGTs (GDGT-
664 0, -1, -2, -3, Cren and Cren') to generate an SST estimate with associated uncertainty. The key advantages
665 of the OPTiMAL approach are: 1) that these uncertainty estimates are intrinsically linked to the strength of
666 the relationship between the fossil GDGT distributions and the modern calibration data set, and 2) by
667 considering all GDGT compounds in a multi-dimensional regression model it avoids the dimensionality
668 reduction and loss of information that takes place when calibrating single parameters (TEX_{86}) to
669 temperature. The methods presented above make very few assumptions about the data. We argue that such
670 methods are appropriate with the current absence of any reasonable mechanistic model for the data
671 generating process, in that they reflect model uncertainty in a natural way. Finally, we note the potential
672 for multi-proxy machine learning approaches, synthesising data from other palaeothermometers with
673 independent uncertainties and biases, to improve calibration of ancient GDGT-derived SST reconstructions.

674

675

676 **Acknowledgements:**

677 TDJ, JAB, IM, KME and YE acknowledge NERC grant NE/P013112/1. SEG was supported by NERC
678 Independent Research Fellowship NE/L011050/1 and NERC large grant NE/P01903X/1. WT
679 acknowledges the Wellcome Trust (grant code: 1516ISSFFEL9, www.wellcome.ac.uk/) for funding a
680 parameterisation workshop at the University of Birmingham (UK). WT, TDJ and IM would like to thank
681 the BBSRC UK Multi-Scale Biology Network Grant No. BB/M025888/1.

682

683

684

685 **Figure Captions:**

686

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

689

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

692

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

696

697 **Figure 4.** The error of a random forest predictor as a function of the true temperature (see Section
698 2.3). This machine learning based predictor yields δT 4.1 °C degrees, by applying a suitable weighted
699 average over multiple near neighbours. $R^2 = 0.83$, meaning that 83% of the variation in the observed
700 temperature is successfully explained by our GDGT-based model.

701

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

704

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

707

708 **Figure 7.** Diffusion map projection of the modern and ancient data. Black: Modern; Blue: Eocene (Inglis
709 et al., 2015); Red: Cretaceous (O'Brien et al., 2017). separate clusters marked 'A' are the outlying

710 Cretaceous points with high GDGT-3 values. Branch ‘B’ is dominated by modern data points; branch ‘C’
711 by Cretaceous data.

712

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

716

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

718

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

720

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

724

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

728

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

734

735 **Figure 14.** Inter-comparison of temperature estimates and errors (y-axis) for compiled Eocene and
736 Cretaceous data calculated using OPTiMAL (top) and BAYSPAR (bottom). Greyed out data fails the
737 $D_{nearest}$ test (>0.5), and the colour scaling reflects $D_{nearest}$ values for those datapoints that pass. The black
738 dashed line shows the $D_{nearest}$ threshold (>0.5).

739

740 **Figure 15.** Comparison of $D_{nearest}$ threshold (>0.5), BIT and MI index for the Eocene (Inglis et al., 2015)
741 and Cretaceous (O’Brien et al., 2017) datasets.

742

743 **Figure 16.** GDGT-derived OPTiMAL palaeotemperatures, for the late Miocene to Recent, from two sites
744 in the Eastern (ODP Site 850) and Western (ODP Site 806) Equatorial Pacific; uncertainty envelopes are

745 one standard deviation each side of the maximum likelihood temperature estimator. Also shown for
746 comparison are U^k_{37} data from the same sites and modern mean annual SSTs.

747

748

749

750

751 **References:**

752 Aitchison, J.: The Statistical Analysis of Compositional Data. J. R. Stat. Soc. Series B Stat. Methodol. 44,
753 139–160, 1982.

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

755 Aitchison, J., Greenacre, M.: Biplots of compositional data. J. R. Stat. Soc. Ser. C Appl. Stat. 51, 375–392,
756 2002.

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

759 Bale, N. J., Palatinszky, M., Rijpstra, I. C., Herbold, C. W., Wagner, M., Sinnighe Damste, J. S.: Membrane
760 lipid composition of the moderately thermophilic ammonia-oxidizing Archaeon “*Candidatus*
761 *Nitrosotenus uzonensis*” at different growth temperatures, Applied and Environmental
762 Microbiolpogy, DOI: 10.1128/AEM.01332-19, 2019.

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

765 Bijl, P. K., Bendle, J.A.P., Bohaty, S.M., Pross, J., Schouten, S., Tauxe, L., Stickley, C., McKay, R.M.,
766 Röhl, U., Olney, M., Sluijs, A., Escutia, C., Brinkhuis, H. and Expedition 318 Scientists.: Eocene
767 cooling linked to early flow across the Tasmanian Gateway. Proc. Natl. Acad. Sci. U.S.A., 110,
768 9645–9650, 2013.

769 Brassell, S. C.: Climatic influences on the Paleogene evolution of alkenones, Paleoceanography, 29, 255-
770 272, doi:10.1002/2013PA002576, 2014.

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

776 Cadillo-Quiroz, H., Didelot, X., Held, N. L., Herrera, A., Darling, A., Reno, M. L., Krause, D. J., Whitaker,
777 R. J.: Patterns of gene flow define species of thermophilic Archaea, PLOS, Cadillo-Quiroz, H. et al.
778 (2012) PLOS Biology, <https://doi.org/10.1371/journal.pbio.1001265>, 2012.
779

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

- 783 Coulston, J.W., Blinn, C.E., Thomas, V.A.: Approximating prediction uncertainty for random forest
784 regression models. *Photogrammetric Engineering & Remote Sensing*, Volume 82, 189-197,
785 <https://doi.org/10.14358/PERS.82.3.189>, 2016.
- 786 Cramwinckel, M. J., Huber, M., Kocken, I. J., Agnini, C., Bijl, P. K., Bohaty, S. M., Frieling, J., Goldner,
787 A., Hilgen, F. J., Kip, E. L., Peterse, F., van der Ploeg, R., Röhl, U., Schouten, S., and Sluijs, A.:
788 Synchronous tropical and polar temperature evolution in the Eocene, *Nature*, 559, 382-386, 2018.
- 789 De Rosa, M., Esposito, E., Gambacorta, A., Nicolaus, B., Bu'Lock, J. D.: Effects of temperatures on ether
790 lipid composition of *Caldariella acidophila*, 19, 827 – 831 1980.
- 791 Douglas, P. M. J., Affek, H. P., Ivany, L. C., Houben, A. J. P., Sijp, W. P., Sluijs, A., Schouten, S., Pagani,
792 M.: Pronounced zonal heterogeneity in Eocene southern high-latitude sea surface temperatures.
793 *Proceedings of the National Academy of Sciences*, 111, 6582–6587, 2014.
- 794 Dunkley Jones, T., Lunt, D. J., Schmidt, D. N., Ridgwell, A., Sluijs, A., Valdes, P. J., and Maslin, M.:
795 Climate model and proxy data constraints on ocean warming across the Paleocene–Eocene Thermal
796 Maximum, *Earth-Science Reviews*, 125, 123-145, 2013.
- 797 Egozcue, J.J., Pawłowsky-Glahn, V., Mateu-Figueras, G., Barceló-Vidal, C.: Isometric Logratio
798 Transformations for Compositional Data Analysis. *Math. Geol.* 35, 279–300, 2003.
- 799 Elling, F. J., Könneke, M., Lipp, J. S., Becker, K. W., Gagen, E. J., and Hinrichs, K.-U.: Effects of growth
800 phase on the membrane lipid composition of the thaumarchaeon *Nitrosopumilus maritimus* and
801 their implications for archaeal lipid 20 distributions in the marine environment, *Geochimica et*
802 *Cosmochimica Acta*, 141, 579-597, 2014.
- 803 Elling, F. J., Könneke, M., Mußmann, M., Greve, A., and Hinrichs, K.-U.: Influence of temperature, pH,
804 and salinity on membrane lipid composition and TEX₈₆ of marine planktonic thaumarchaeal
805 isolates, *Geochimica et Cosmochimica Acta*, 171, 238-255, 2015.
- 806 Elling, F.J., Konnecke, M., Nicol, G. W., Stieglmeier, M., Bayer, B., Spieck, E., de la Torre, J. R., Becker,
807 K. W., Thomm, M. Prosser, J. I., Herndl, G., Schleper, C., Hinrichs, K-U.: Chemotaxonomic
808 characterisation of the thaumarchaeal lipidome, *Environmental Microbiology* 19, 2681–2700 ,
809 2017.
810
- 811 Filzmoser, P., Hron, K., Reimann, C.: Principal component analysis for compositional data with outliers.
812 *Environmetrics* 20, 621–632, 2009a.
- 813 Filzmoser, P., Hron, K., Reimann, C., Garrett, R.: Robust factor analysis for compositional data. *Comput.*
814 *Geosci.* 35, 1854–1861, 2009b.
- 815 Filzmoser, P., Hron, K., Reimann, C.: Interpretation of multivariate outliers for compositional data.
816 *Comput. Geosci.* 39, 77–85, 2012.
- 817 Gliozzi, A., Paoli, G., De Rosa, M., Gambacorta, A.: Effect of isoprenoid cyclization on the transition
818 temperature of lipids in thermophilic archaeobacteria, *Biochimica et Biophysica Acta (BBA)*
819 *Biomembranes*, 735, 234 – 242, 1983.
- 820 Haghverdi, L., Buettner, F., Theis, F.J: Diffusion maps for high-dimensional single-cell analysis of
821 differentiation data. *Bioinformatics* 31, 2989–2998, 2015.

- 822 Haghverdi, L., Büttner, M., Wolf, F.A., Buettner, F., Theis, F.J.: Diffusion pseudotime robustly
823 reconstructs lineage branching. *Nat. Methods* 13, 845–848, 2016.
- 824 Haslett, J., Whitley, M., Bhattacharya, S., Salter-Townshend, M., Wilson, S.P., Allen, J.R.M., Huntley, B.,
825 Mitchell, F.J.G.: Bayesian palaeoclimate reconstruction. *J Royal Statistical Soc A* 169, 395–438,
826 2006.
- 827 Herfort, L., Schouten, S., Boon, J. P., and Sinninghe Damsté, J. S.: Application of the TEX₈₆ temperature
828 proxy to the southern North Sea, *Organic Geochemistry*, 37, 1715-1726, 2006.
- 829 Hertzberg, J. E., Schmidt, M. W., Bianchi, T. S., Smith, R. K., Shields, M. R., & Marcantonio, F.:
830 Comparison of eastern tropical Pacific TEX₈₆ and Globigerinoides ruber Mg/Ca derived sea surface
831 temperatures: Insights from the Holocene and Last Glacial Maximum. *Earth and Planetary Science*
832 *Letters*, 434, 320–332, 2016.
- 833 Hollis, C. J., Taylor, K. W. R., Handley, L., Pancost, R. D., Huber, M., Creech, J. B., Hines, B. R., Crouch,
834 E. M., Morgans, H. E. G., Crampton, J. S., Gibbs, S., Pearson, P. N., and Zachos, J. C.: Early
835 Paleogene temperature history of the Southwest Pacific Ocean: Reconciling proxies and models,
836 *Earth and Planetary Science Letters*, 349–350, 53-66, 2012.
- 837 Hollis, C. J., Dunkley Jones, T., Anagnostou, E., Bijl, P. K., Cramwinckel, M. J., Cui, Y., Dickens, G. R.,
838 Edgar, K. M., Eley, Y., Evans, D., Foster, G. L., Frieling, J., Inglis, G. N., Kennedy, E. M.,
839 Kozdon, R., Lauretano, V., Lear, C. H., Littler, K., Meckler, N., Naafs, B. D. A., Pälike, H.,
840 Pancost, R. D., Pearson, P., Royer, D. L., Salzmann, U., Schubert, B., Seebeck, H., Sluijs, A.,
841 Speijer, R., Stassen, P., Tierney, J., Tripathi, A., Wade, B., Westerhold, T., Witkowski, C., Zachos,
842 J. C., Zhang, Y. G., Huber, M., and Lunt, D. J.: The DeepMIP contribution to PMIP4:
843 methodologies for selection, compilation and analysis of latest Paleocene and early Eocene climate
844 proxy data, incorporating version 0.1 of the DeepMIP database, *Geosci. Model Dev. Discuss.*,
845 <https://doi.org/10.5194/gmd-2018-309>, in review, 2019.
- 846 Hollis, C. J., Handley, L., Crouch, E. M., Morgans, H. E., Baker, J. A., Creech, J., Collins, K. S., Gibbs, S.
847 J., Huber, M., Schouten, S.: Tropical sea temperatures in the high-latitude South Pacific during the
848 Eocene. *Geology*, 37, 99–102, 2009.
- 849 Hopmans, E. C., Weijers, J. W. H., Schefuss, E., Herfort, L., Sinninghe Damsté, J. S., Schouten, S.: A
850 novel proxy for terrestrial organic matter in sediments based on branched and isoprenoid tetraether
851 lipids, *Earth and Planetary Science Letters*, 224, 107-116, 2004.
- 852 Huguet C, Kim J-H, Sinninghe Damsté J.S., Schouten S: Reconstruction of sea surface temperature
853 variations in the Arabian Sea over the last 23 kyr using organic proxies (TEX₈₆ and UK 0 37 .
854 *Paleoceanography* 21(3): PA3003, 2006.
- 855 Hurley, S. J., Elling, F. J., Könneke, M., Buchwald, C., Wankel, S. D., Santoro, A. E., Lipp, J.S., Hinrichs,
856 K., Pearson, A.: Influence of ammonia oxidation rate on thaumarchaeal lipid composition and the
857 TEX₈₆ temperature proxy. *Proceedings of the National Academy of Sciences*, 113, 7762–7767,
858 2016.
- 859 Inglis, G. N., Farnsworth, A., Lunt, D., Foster, G. L., Hollis, C. J., Pagani, M., Jardine, P. E., Pearson, P.
860 N., Markwick, P., Galsworthy, A. M. J., Raynham, L., Taylor, K. W. R., and Pancost, R. D.:

- 861 Descent toward the Icehouse: Eocene sea surface cooling inferred from GDGT distributions,
862 *Paleoceanography*, 30, 1000-1020, 2015.
- 863 Jenkyns H.C., Schouten-Huibers L., Schouten S., Damsté J.S.S.: Warm Middle Jurassic-Early Cretaceous
864 high-latitude sea-surface temperatures from the Southern Ocean. *Clim Past* 8 (1):215–226, 2012.
- 865 Kim, J.-H., Schouten, S., Hopmans, E. C., Donner, B., and Sinninghe Damsté, J. S.: Global sediment core-
866 top calibration of the TEX₈₆ paleothermometer in the ocean, *Geochimica et Cosmochimica Acta*,
867 72, 1154-1173, 2008.
- 868 Kim, J.-H., van der Meer, J., Schouten, S., Helmke, P., Willmott, V., Sangiorgi, F., Koç, N., Hopmans, E.
869 C., and Sinninghe Damsté, J. S.: New indices and calibrations derived from the distribution of
870 crenarchaeal isoprenoid tetraether lipids: Implications for past sea surface temperature
871 reconstructions, *Geochimica et Cosmochimica Acta*, 74, 4639-4654, 2010.
- 872 Linnert, C., Robinson, S. A., Lees, J. A., Bown, P. R., Perez-Rodriguez, I., Petrizzo, M. R., Falzoni, F.,
873 Littler, K., Antonio Arz, J., Russell, E. E. : Evidence for global cooling in the Late Cretaceous.
874 *Nature Communications*, 5, 1–7, 2014.
- 875 Liu, X-L., Zhu, C., Wakeham, S.G., Hinrichs, K-U.: In situ production of branched glycerol dialkyl
876 glycerol tetraethers in anoxic marine water columns, *Marine Chemistry*, 166, 1 – 8, 2014.
- 877 Lunt, D. J., Dunkley Jones, T., Heinemann, M., Huber, M., LeGrande, A., Winguth, A., Loptson, C.,
878 Marotzke, J., Tindall, J., 15 Valdes, P., Winguth, C.: A model-data comparison for a multi-model
879 ensemble of early Eocene atmosphere-ocean simulations: EoMIP, *Clim. Past Discuss.*, 8, 1229-
880 1273, 2012.
- 881 Mentch, L., Hooker, G.: Quantifying Uncertainty in Random Forests via Confidence Intervals and
882 Hypothesis Tests. *Journal of Machine Learning Research*, 17, 1-41, 2016.
- 883 O'Brien, C. L., Robinson, S. A., Pancost, R. D., Sinninghe Damsté, J. S., Schouten, S., Lunt, D. J., Alsenz,
884 H., Bornemann, 20 A., Bottini, C., Brassell, S. C., Farnsworth, A., Forster, A., Huber, B. T., Inglis,
885 G. N., Jenkyns, H. C., Linnert, C., Littler, K., Markwick, P., McAnena, A., Mutterlose, J., Naafs, B.
886 D. A., Püttmann, W., Sluijs, A., van Helmond, N. A. G. M., Vellekoop, J., Wagner, T., Wrobel, N.
887 E.: Cretaceous sea-surface temperature evolution: Constraints from TEX₈₆ and planktonic
888 foraminiferal oxygen isotopes, *Earth-Science Reviews*, 172, 224-247, 2017.
- 889 Park, E., Hefter, J., Fischer, G., Mollenhauer, G.: TEX₈₆ in sinking particles in three eastern Atlantic
890 upwelling regimes. *Organic Geochemistry*, 124, 151–163, 2018.
- 891 Pearson, P. N., van Dongen, B. E., Nicholas, C. J., Pancost, R. D., Schouten, S., Singano, J. M., Wade, B.
892 S.: Stable warm tropical climate through the Eocene Epoch. *Geology*, 35, 211-214, 2007.
- 893 Polik, C. A., Elling, F. J., Pearson, A.: Impacts of Paleoeology on the TEX₈₆ Sea Surface Temperature
894 Proxy in the Pliocene-Pleistocene Mediterranean Sea. *Paleoceanography and Paleoclimatology*, 33,
895 1472–1489, 2018.
- 896 Qin, W., Amin, S. A., Martens-Habbena, W., Walker, C. B., Urakawa, H., Devol, A. H., Ingalls, A.E.,
897 Moffett, J.W., Ambrust, E.V., Stahl, D. A.: Marine ammonia-oxidizing archaeal isolates display
898 obligate mixotrophy and wide ecotypic variation. *Proceedings of the National Academy of*
899 *Sciences of the United States of America*, 111, 12504–12509, 2014.

- 900 Qin, W., Carlson, L. T., Armbrust, E. V., Stahl, D. A., Devol, A. H., Moffett, J. W., Ingalls, A. E.:
901 Confounding effects of oxygen and temperature on the TEX₈₆ signature of marine
902 Thaumarchaeota. *Proceedings of the National Academy of Sciences*, 112, 10979–10984, 2015.
- 903 Rasmussen, C.E., Nickisch, H.: Gaussian Processes for Machine Learning (GPML) Toolbox. *J. Mach.*
904 *Learn. Res.* 11, 3011–3015, 2010.
- 905 Sangiorgi, F., van Soelen Els, E., Spofforth David, J. A., Pälike, H., Stickley Catherine, E., St. John, K.,
906 Koç, N., Schouten, S., Sinninghe Damsté Jaap, S., Brinkhuis, H.: Cyclicity in the middle Eocene
907 central Arctic Ocean sediment record: Orbital forcing and environmental response,
908 *Paleoceanography*, 23, 10.1029/2007PA001487, 2008.
- 909 Schouten, E., Hopmans, E.C., Forster, A., Van Breugel, Y., Kuypers, M.M.M., Sinninghe Damsté, J.S.:
910 Extremely high seasurface temperatures at low latitudes during the middle Cretaceous as revealed
911 by archaeal membrane lipids. *Geology*, 31, 1069–1072, 2003.
- 912 Schouten, S., Forster, A., Panoto, F. E., and Sinninghe Damsté, J. S.: Towards calibration of the TEX₈₆
913 palaeothermometer for 20 tropical sea surface temperatures in ancient greenhouse worlds, *Organic*
914 *Geochemistry*, 38, 1537-1546, 2007.
- 915 Schouten, S., Hopmans, E. C., Sinninghe Damsté, J. S.: The organic geochemistry of glycerol dialkyl
916 glycerol tetraether lipids: A review, *Organic Geochemistry*, 54, 19-61, 2013.
- 917 Schouten, S., Hopmans, E. C., Schefuß, E., Sinninghe Damsté, J. S.: Distributional variations in marine
918 crenarchaeotal membrane lipids: a new tool for reconstructing ancient sea water temperatures?
919 *Earth and Planetary Science Letters*, 204, 15 265-274, 2002.
- 920 Seki, O., Bendle, J. A., Harada, N., Kobayashi, M., Sawada, K., Moossen, H., Sakamoto, T.: Assessment
921 and calibration of TEX₈₆ paleothermometry in the Sea of Okhotsk and sub-polar North Pacific
922 region: Implications for paleoceanography. *Progress in Oceanography*, 126, 254–266, 2014.
- 923 Sluijs A, Schouten S, Pagani M, Woltering, M., Brinkhuis, H., Sinninghe Damsté, J.S., Dickens, G.R.,
924 Huber, M., Reichart, G., Stein, R., Matthiessen, J., Lourens, L.J., Pedentchouk, N., Backman, J.,
925 Moran, K. and the Expedition 320 Scientists: Subtropical arctic ocean temperatures during the
926 Palaeocene/Eocene thermal maximum. *Nature* 441, 610–613, 2006.
- 927 Sluijs, A., Schouten, S., Donders, T. H., Schoon, P. L., Rohl, U., Reichart, G.-J., Sangiorgi, F., Kim, J.-H.,
928 Sinninghe Damsté, J. S., Brinkhuis, H.: Warm and wet conditions in the Arctic region during
929 Eocene Thermal Maximum 2, *Nature Geosci*, 2, 777-780, 2009.
- 930 Taylor, K. W. R., Willumsen, P. S., Hollis, C. J., Pancost, R. D.: South Pacific evidence for the long-term
931 climate impact of the Cretaceous/Paleogene boundary event, *Earth-Science Reviews*, 179, 287-302,
932 2018.
- 933 Taylor, K. W., Huber, M., Hollis, C. J., Hernandez-Sanchez, M. T., Pancost, R. D.: Re-evaluating modern
934 and Palaeogene GDGT distributions: Implications for SST reconstructions, *Global and Planetary*
935 *Change*, 108, 158-174, 2013.
- 936 Tierney, J. E.: GDGT Thermometry: Lipid Tools for Reconstructing Paleotemperatures. Retrieved from
937 https://www.geo.arizona.edu/~jesst/resources/TierneyPSP_GDGTs.pdf, 2012

- 938 Tierney, J. E., and Tingley, M. P.: A Bayesian, spatially-varying calibration model for the TEX₈₆ proxy.
939 *Geochimica et Cosmochimica Acta*, 127, 83-106, 2014.
- 940 Tierney, J. E., and Tingley, M. P.: A TEX₈₆ surface sediment database and extended Bayesian calibration,
941 *Scientific data*, 2, 150029, 2015.
- 942 Williams, C.K.I., and Rasmussen, C.E.: Gaussian processes for machine learning. MIT Press Cambridge,
943 MA, 2006.
- 944 Wuchter, C., Schouten, S., Coolen, M. J. L., and Sinninghe Damsté, J. S.: Temperature-dependent variation
945 in the distribution 30 of tetraether membrane lipids of marine Crenarchaeota: Implications for
946 TEX₈₆ paleothermometry, *Paleoceanography and Paleoclimatology*. doi:10.1029/2004PA001041,
947 2004.
- 948 Zhang, Y. G., and Liu, X.: Export Depth of the TEX₈₆ Signal. *Paleoceanography and Paleoclimatology*.
949 doi.org/10.1029/2018PA003337, 2018.
- 950 Zhang, Y. G., Pagani, M., Wang, Z.: Ring Index: A new strategy to evaluate the integrity of TEX₈₆
951 paleothermometry, *Paleoceanography*, 31, 220-232, 2016.
- 952 Zhang, Y. G., Zhang, C. L., Liu, X.-L., Li, L., Hinrichs, K.-U., Noakes, J. E.: Methane Index: a tetraether
953 archaeal lipid 15 biomarker indicator for detecting the instability of marine gas hydrates, *Earth and*
954 *Planetary Science Letters*, 307, 525- 534, 2011.
- 955 Zhang, Y.G., Pagani, M., Liu, Z.: A 12-million-year temperature history of the tropical Pacific
956 *Ocean: Science*, 343, 84-86, 2014.
- 957 Zhu, J., Poulsen, C. J., Tierney, J.: Simulation of Eocene extreme warmth and high climate sensitivity
958 through cloud feedbacks, *Science Advances*, 5(9), eaax1874, 2019.
- 959
- 960

Figure 1

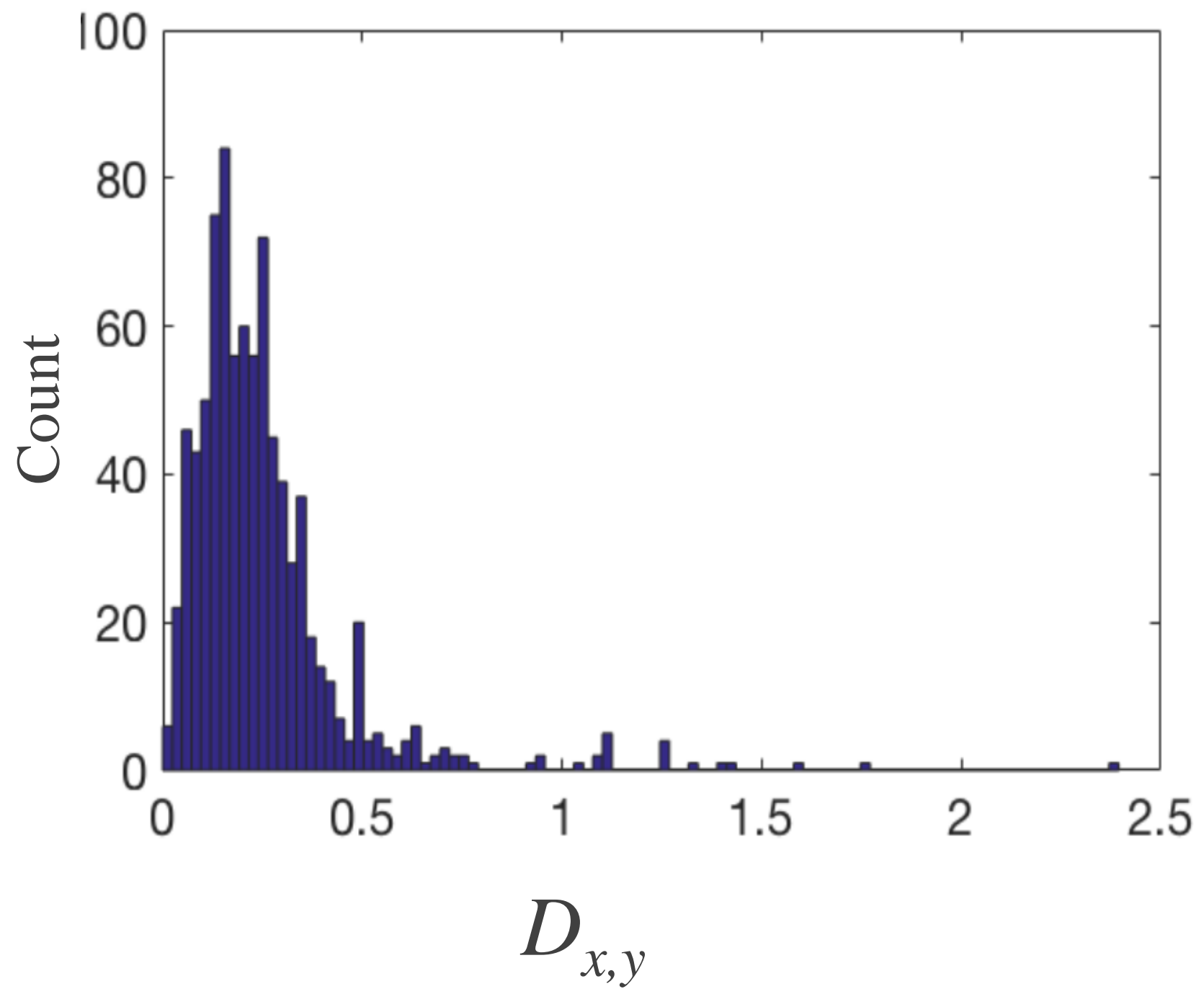


Figure 2

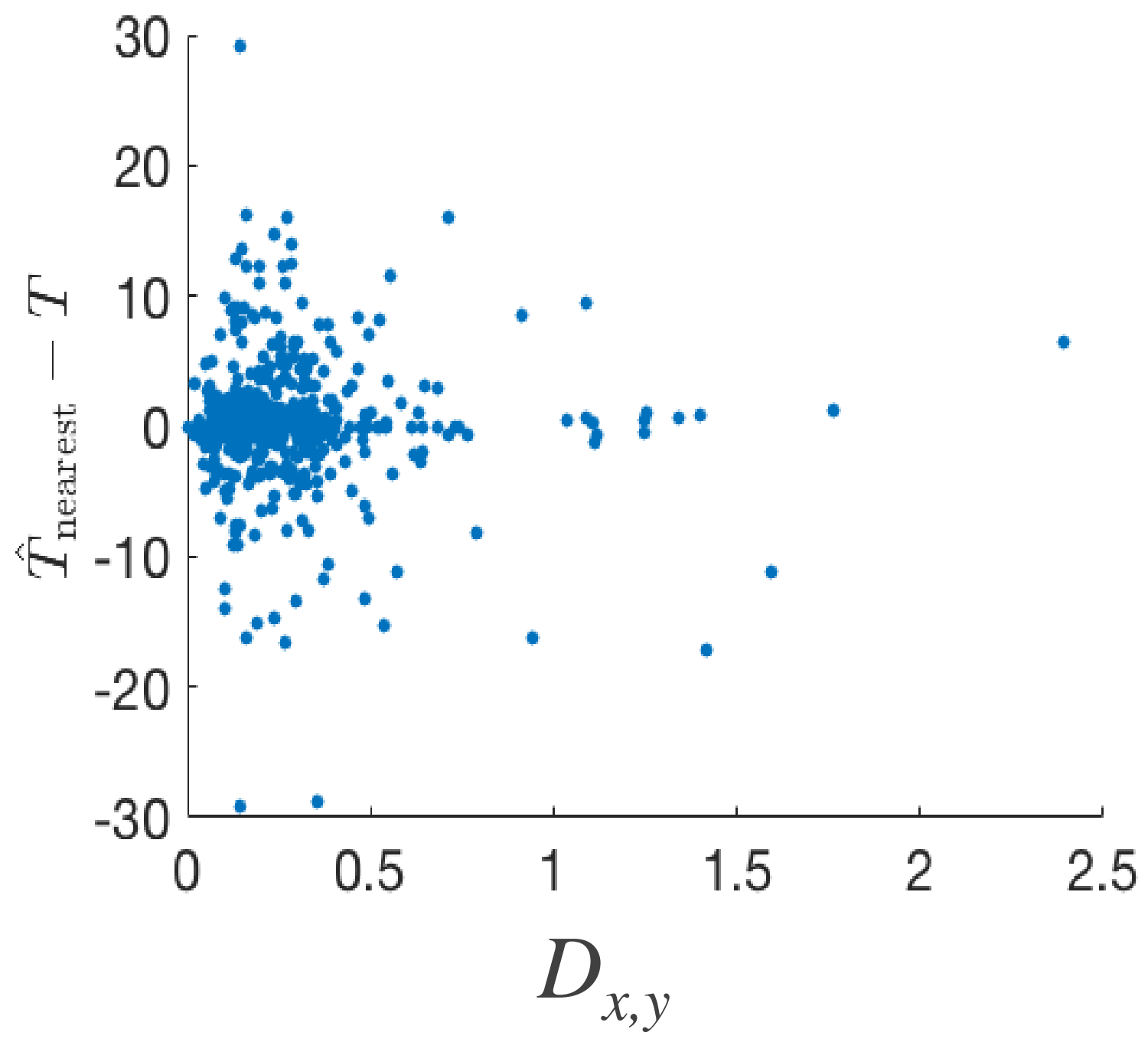


Figure 3

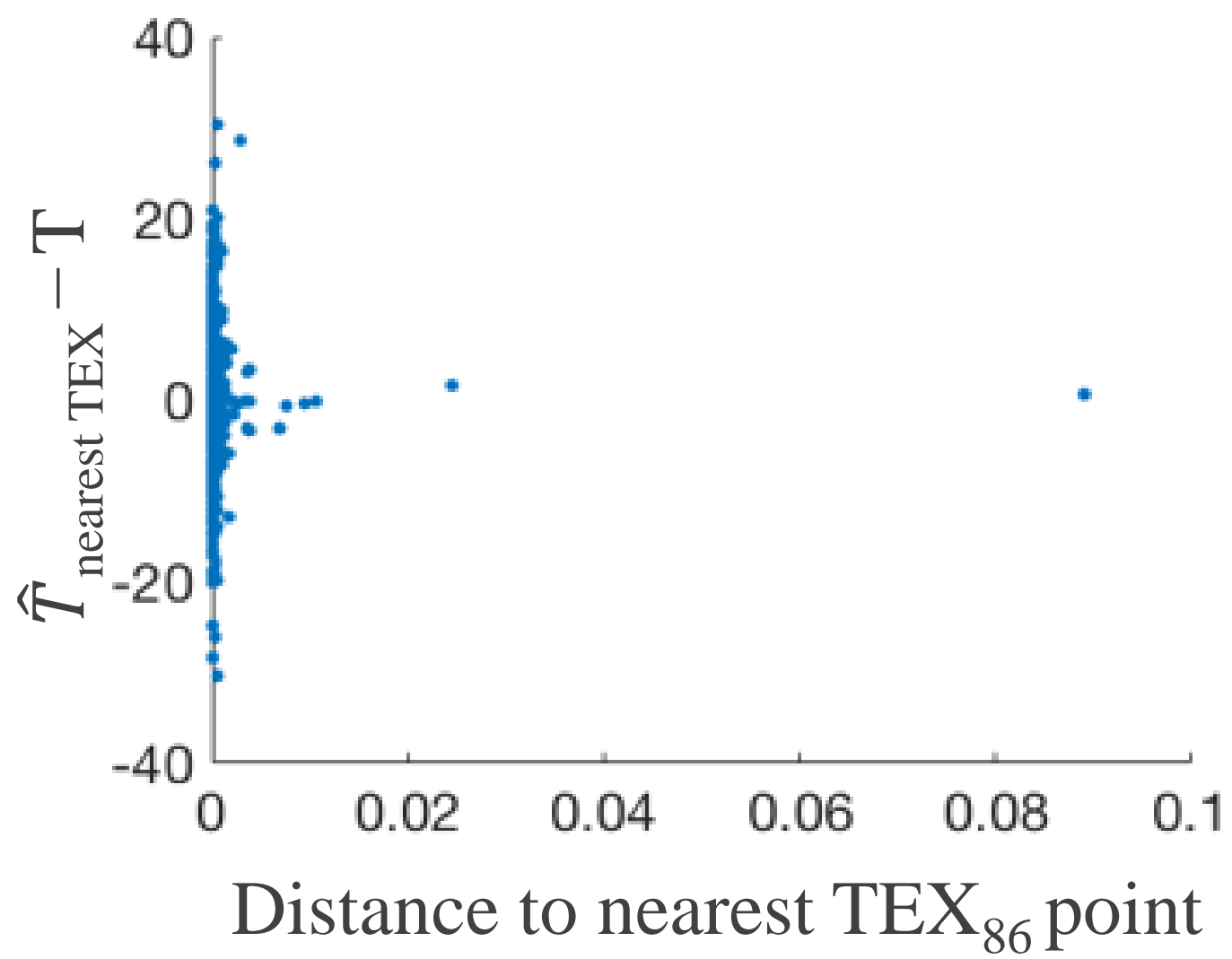
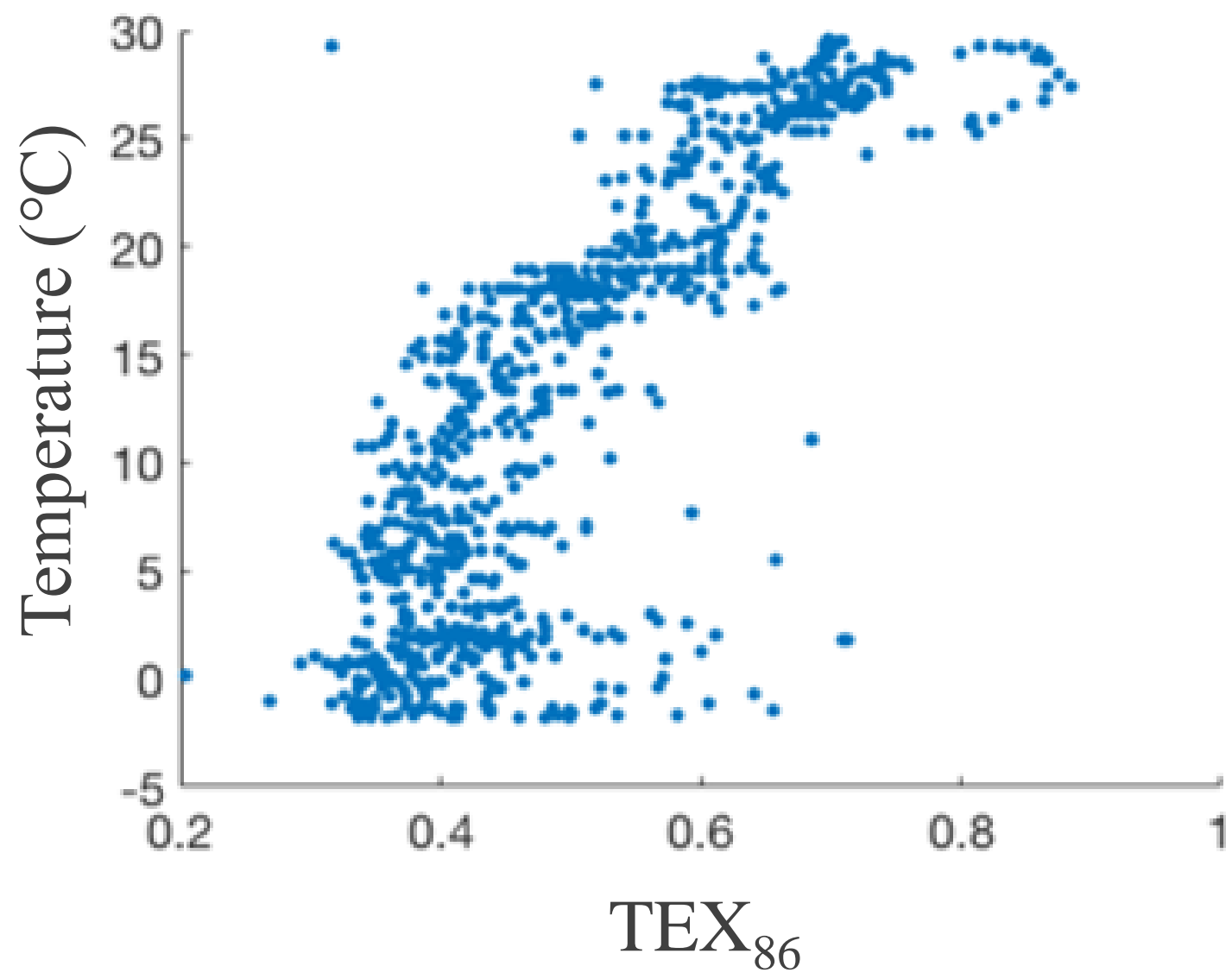


Figure 4

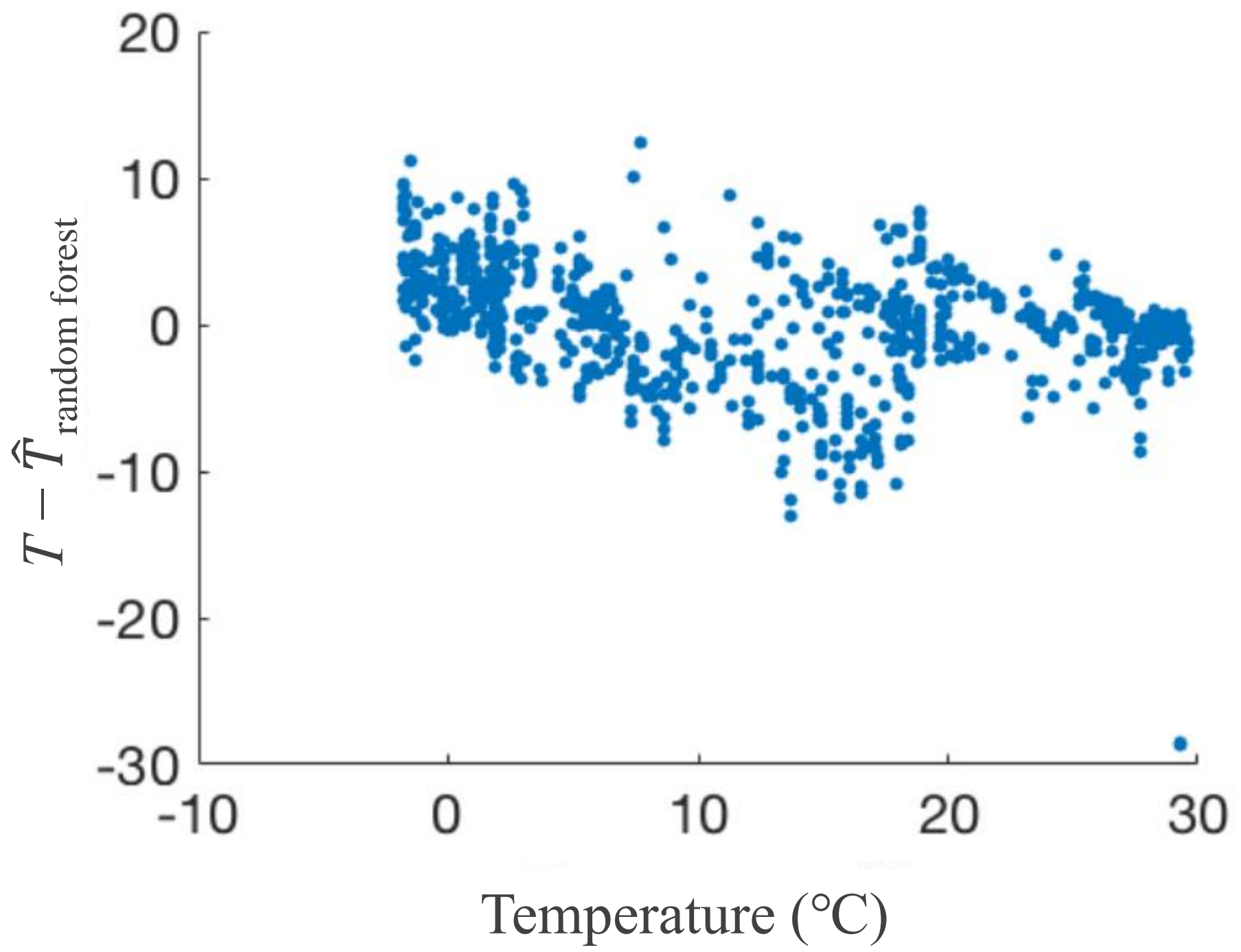


Figure 5

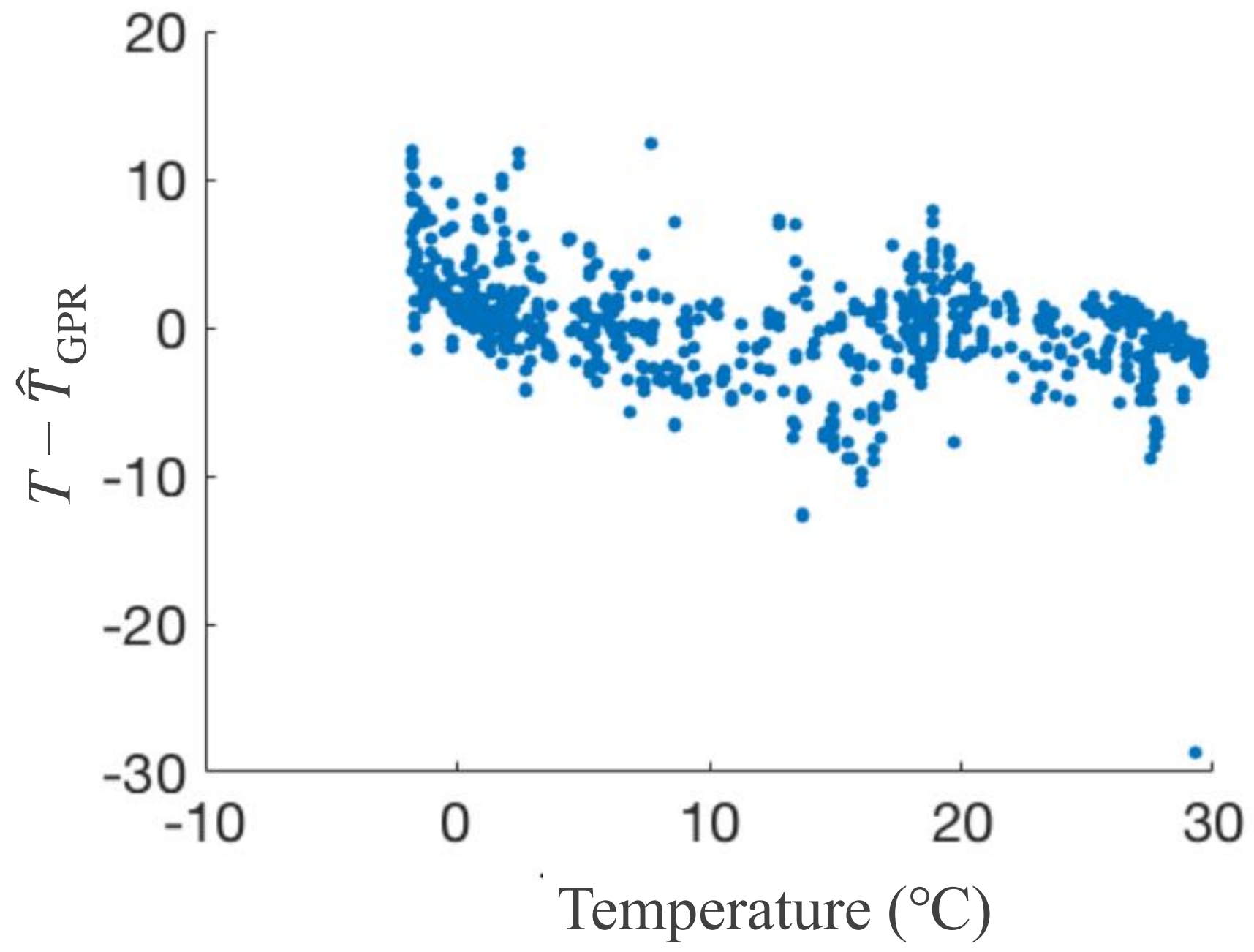


Figure 6

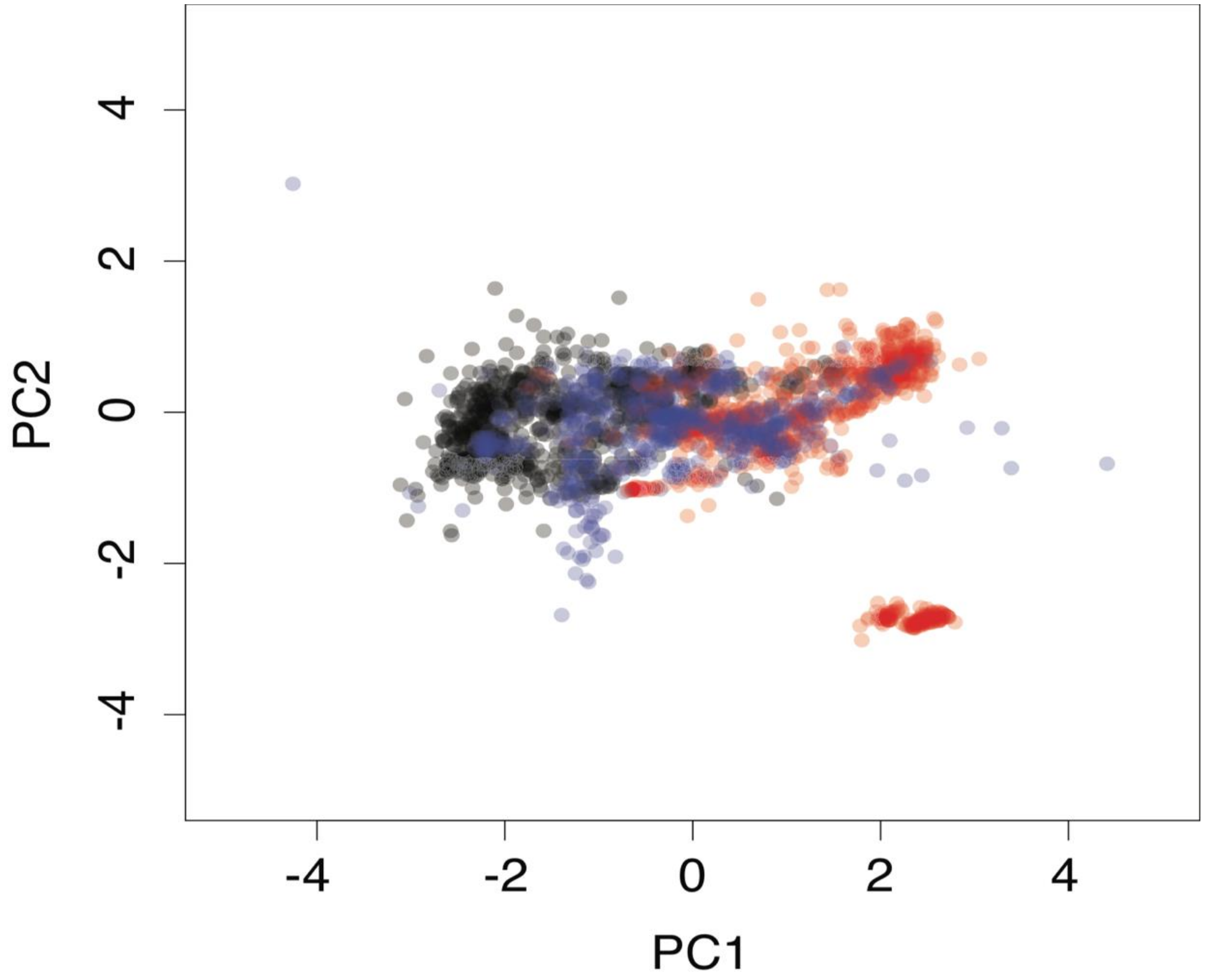


Figure 7

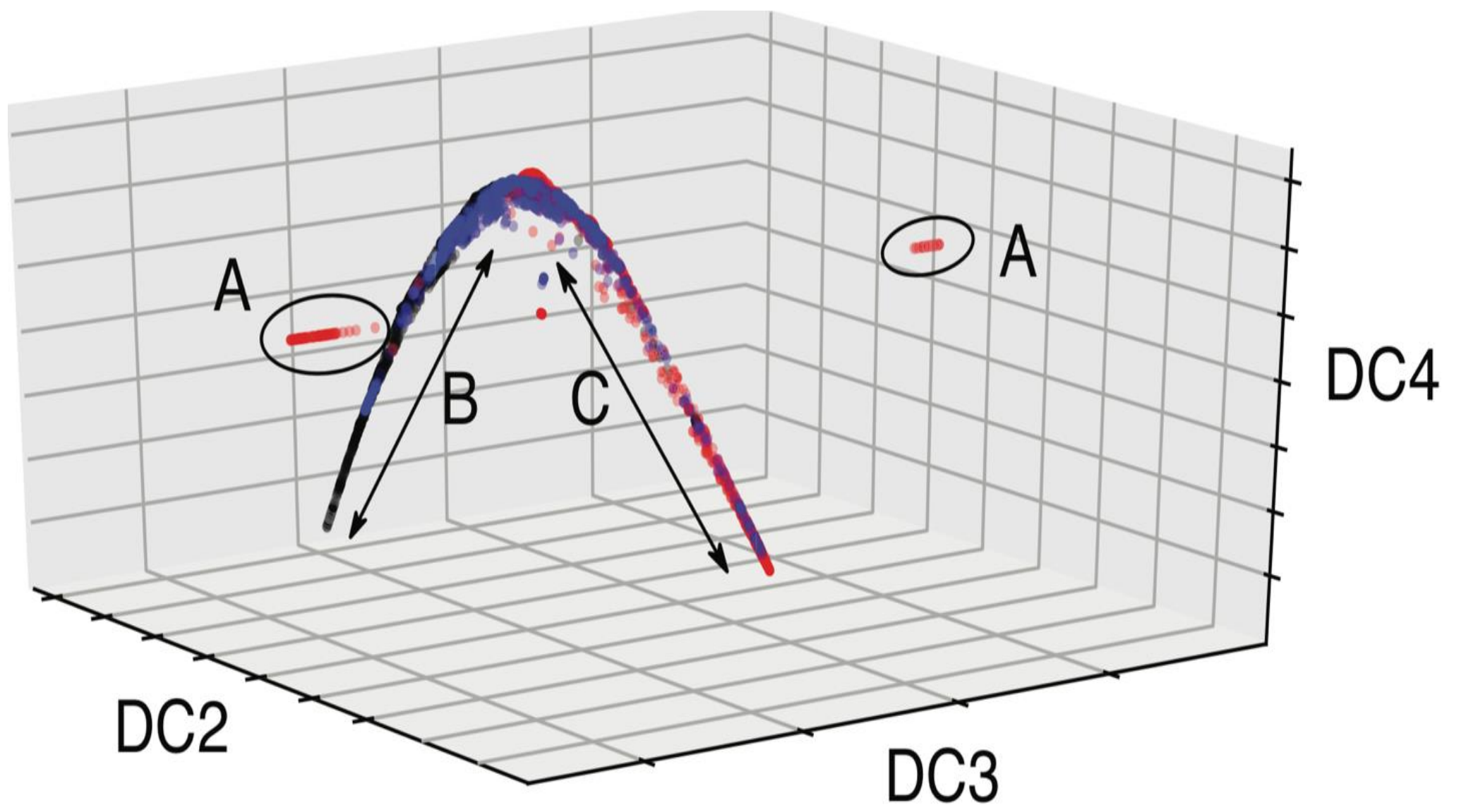


Figure 8

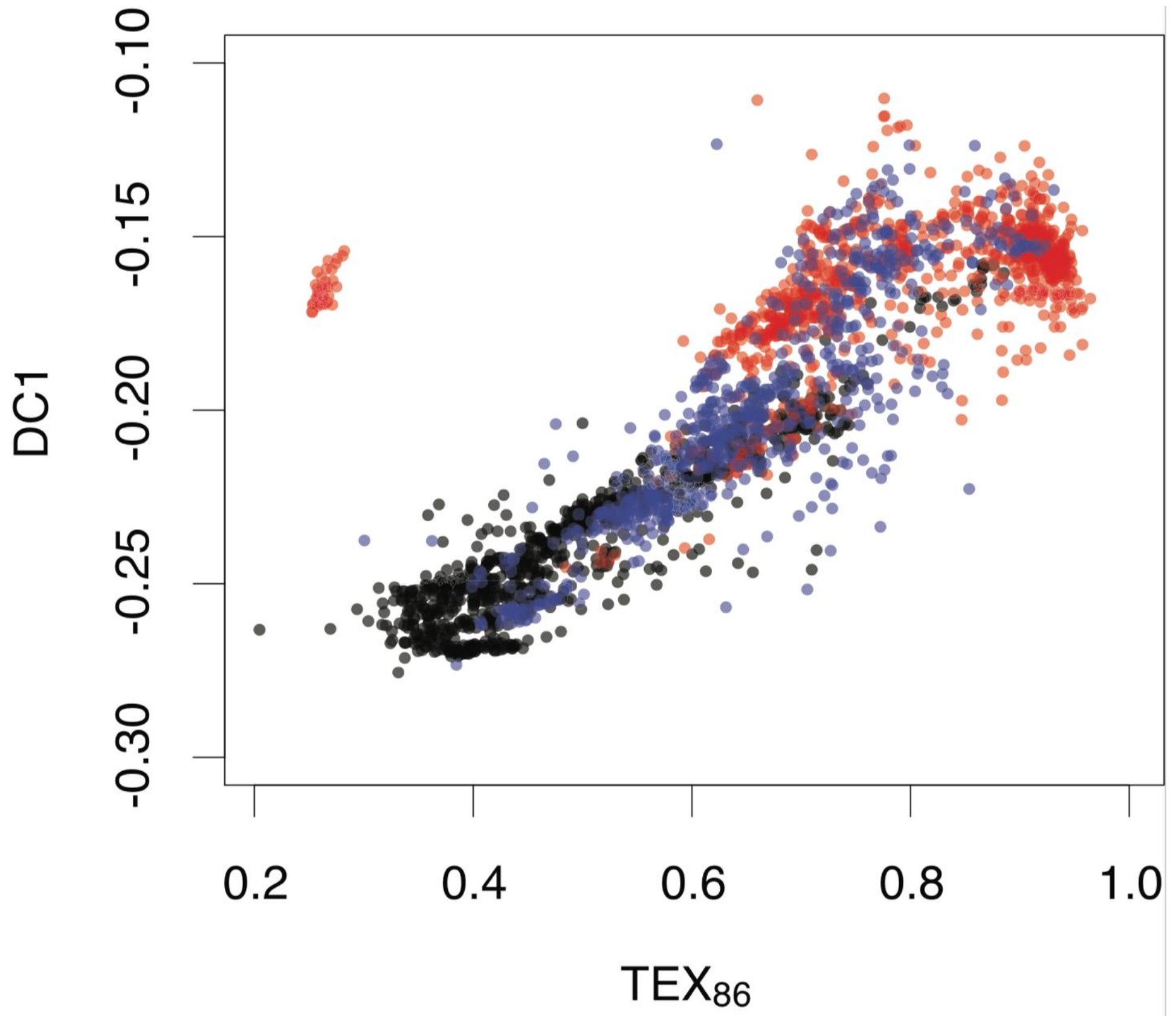


Figure 9

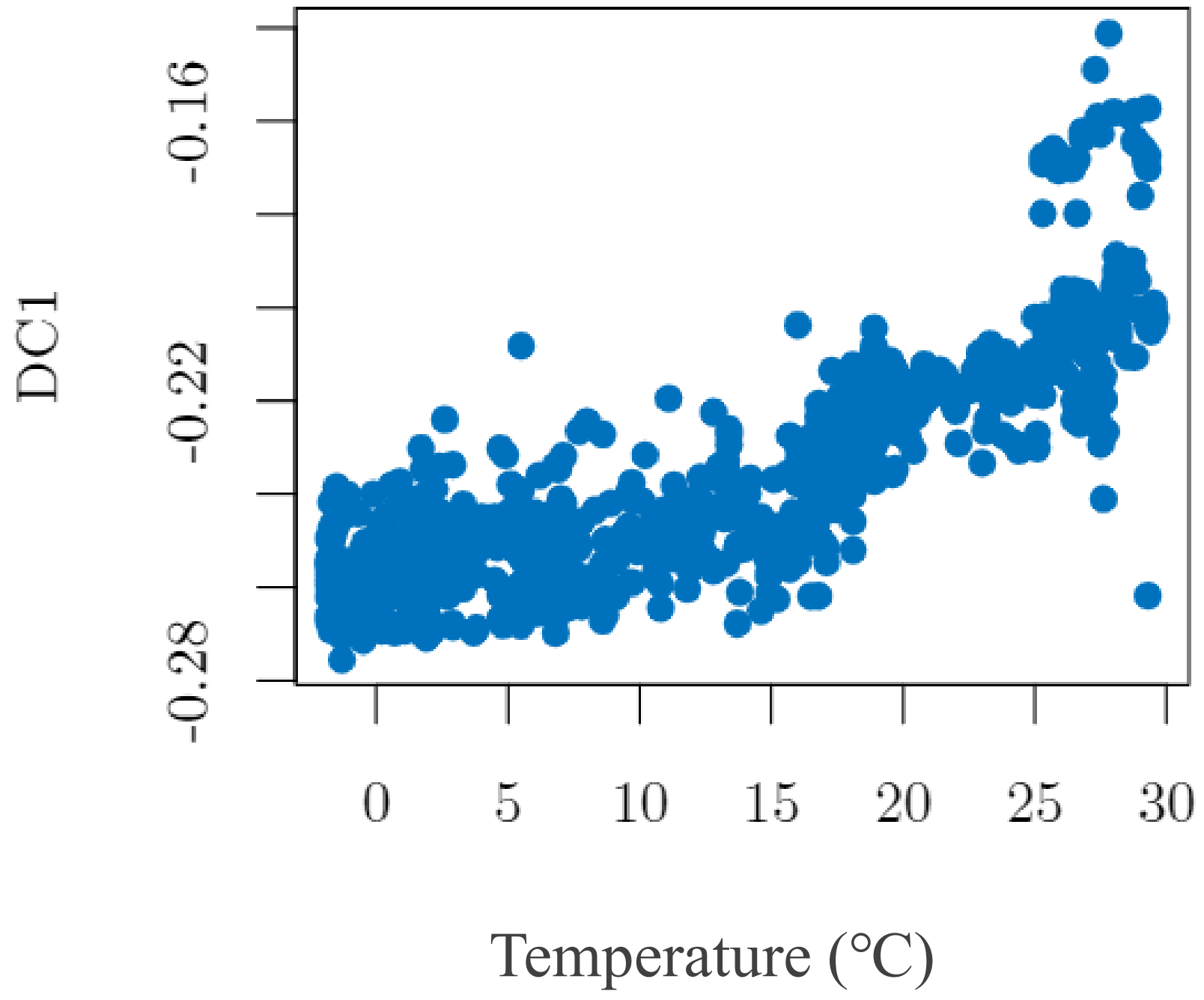


Figure 10

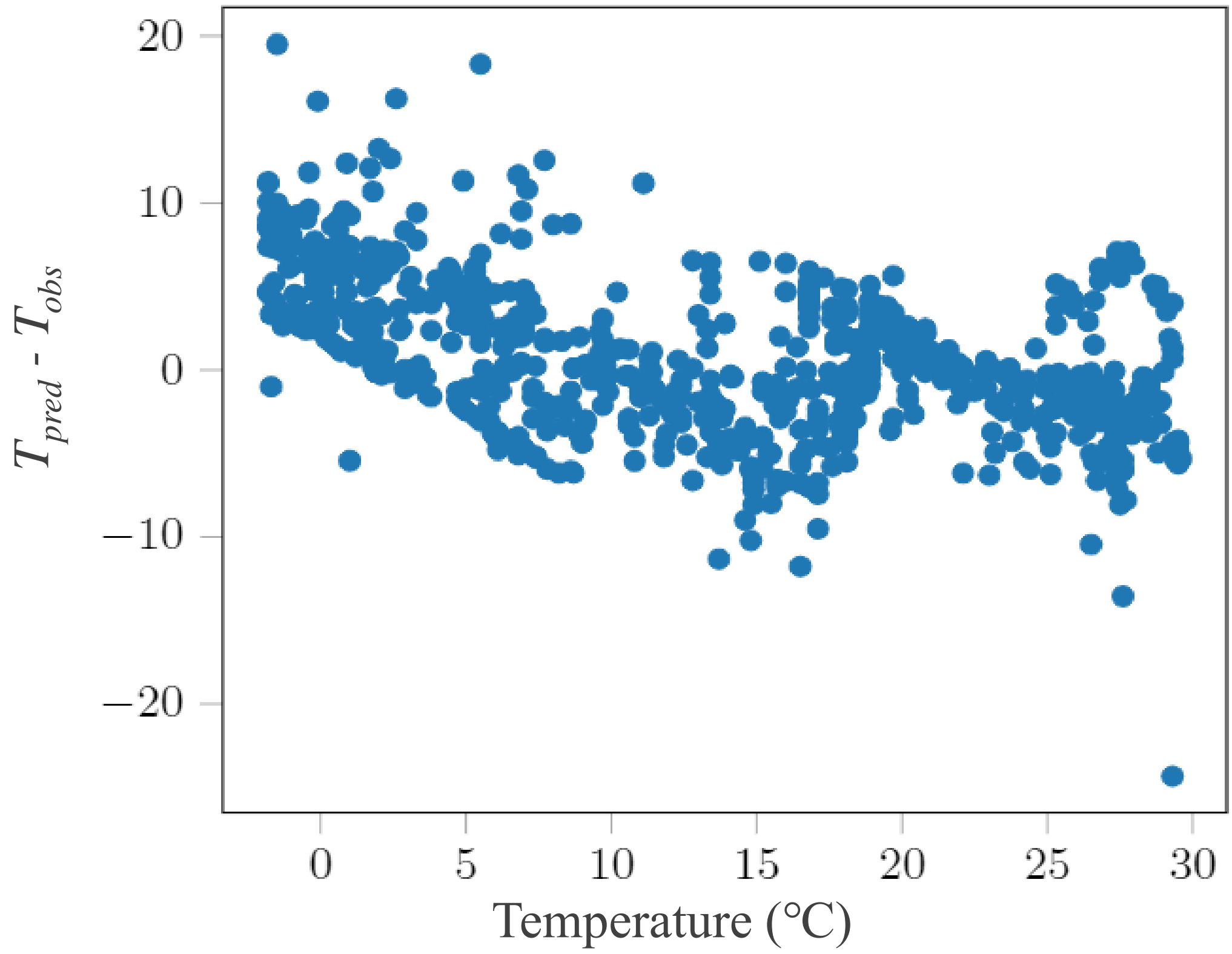


Figure 11

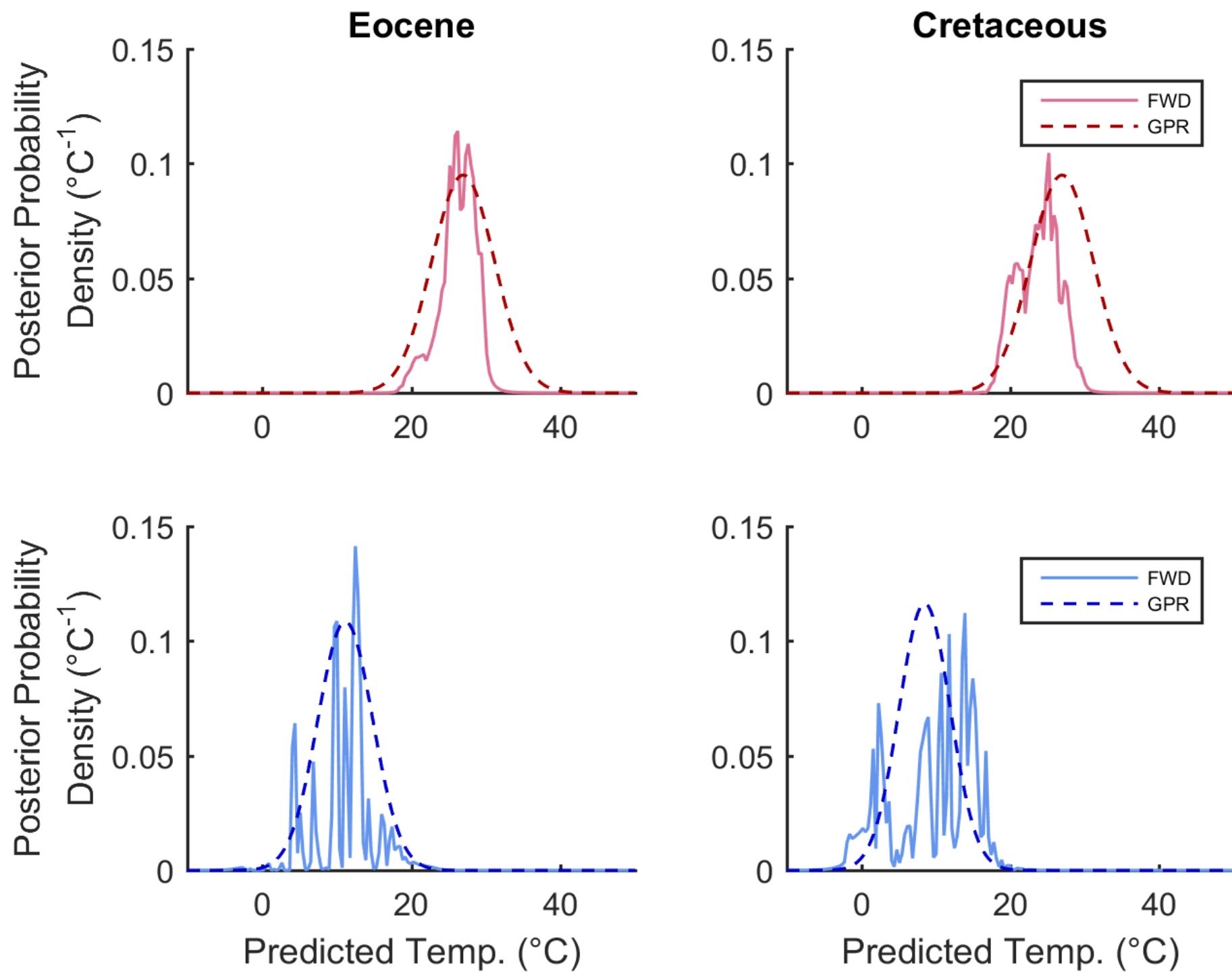


Figure 12

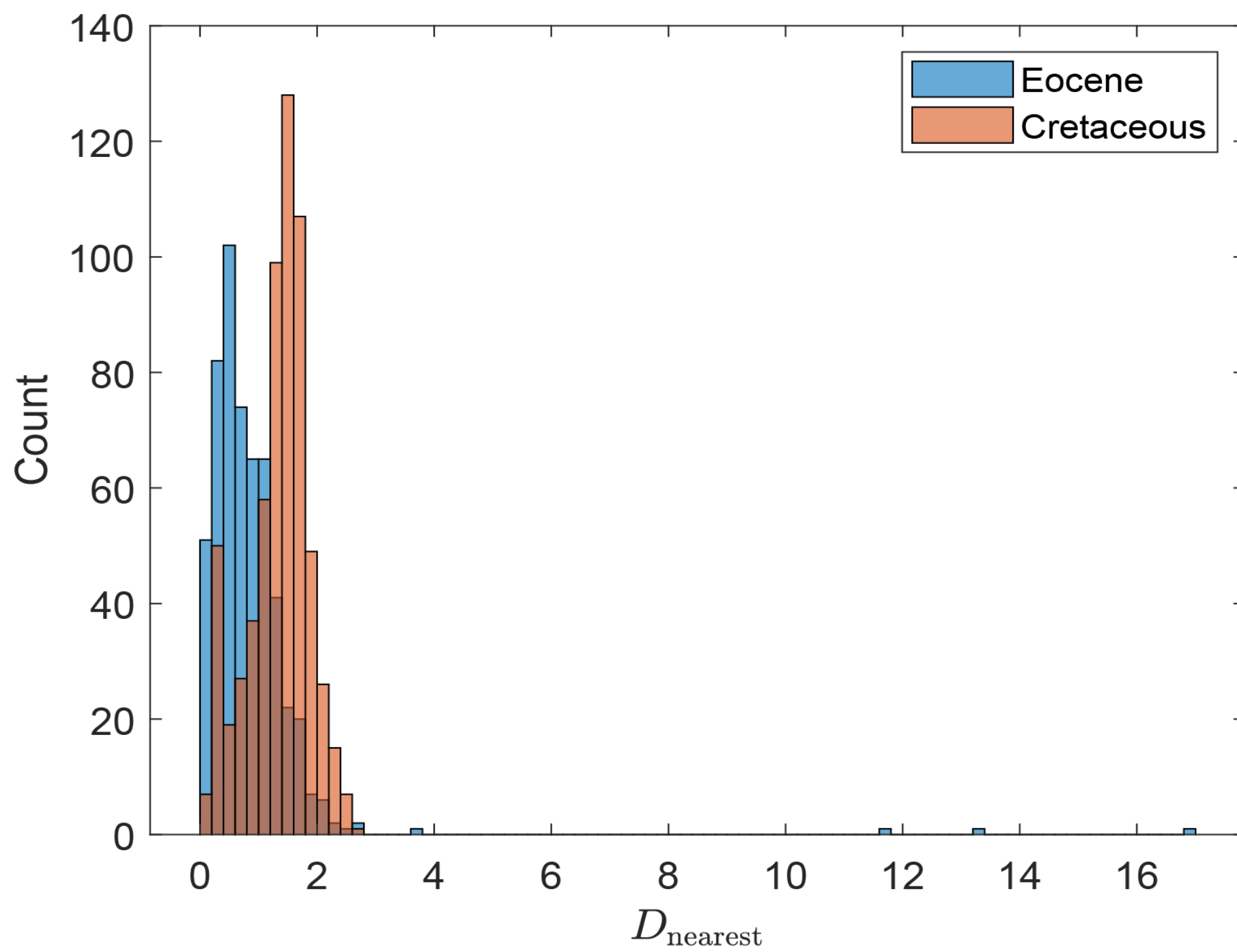


Figure 13

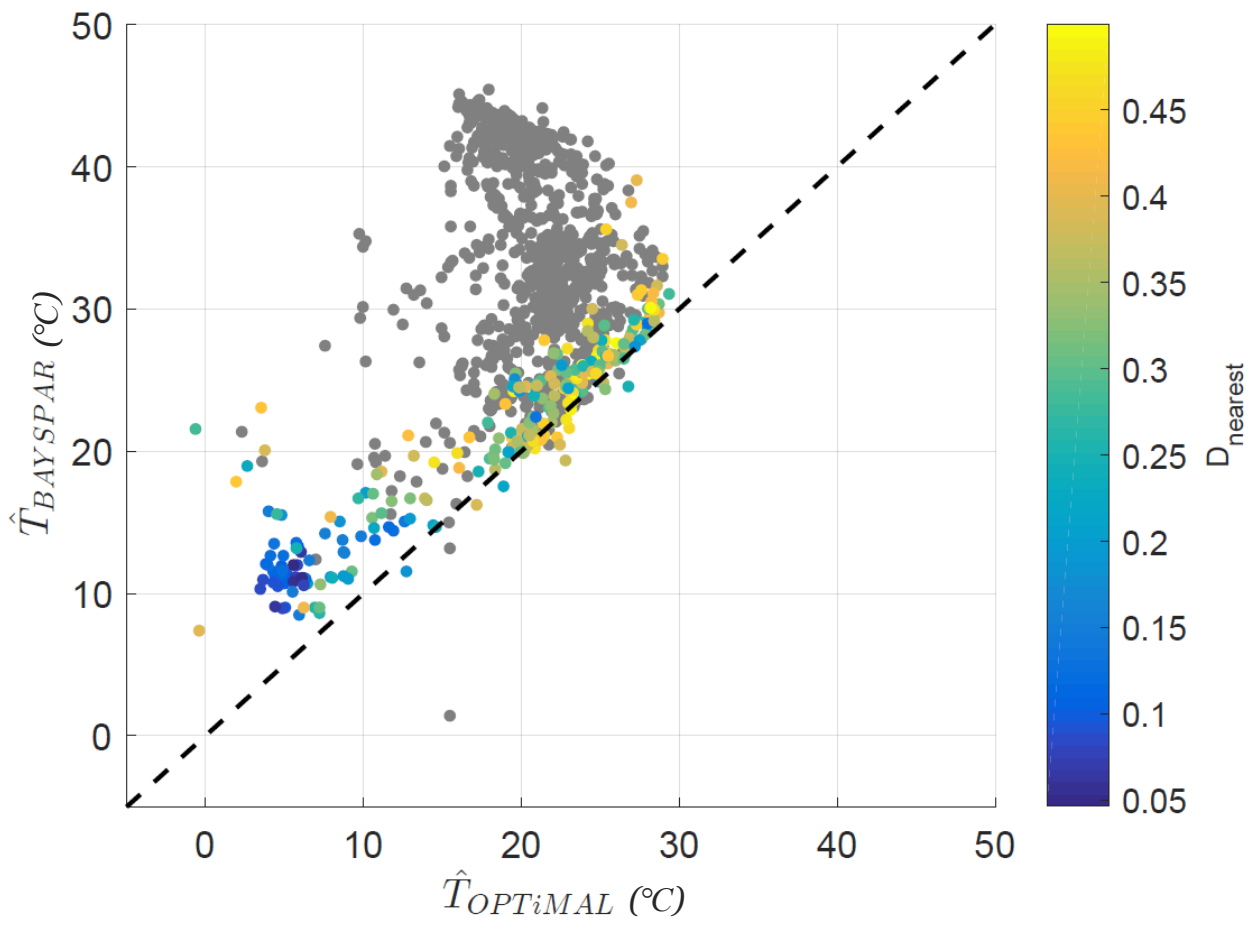


Figure 14

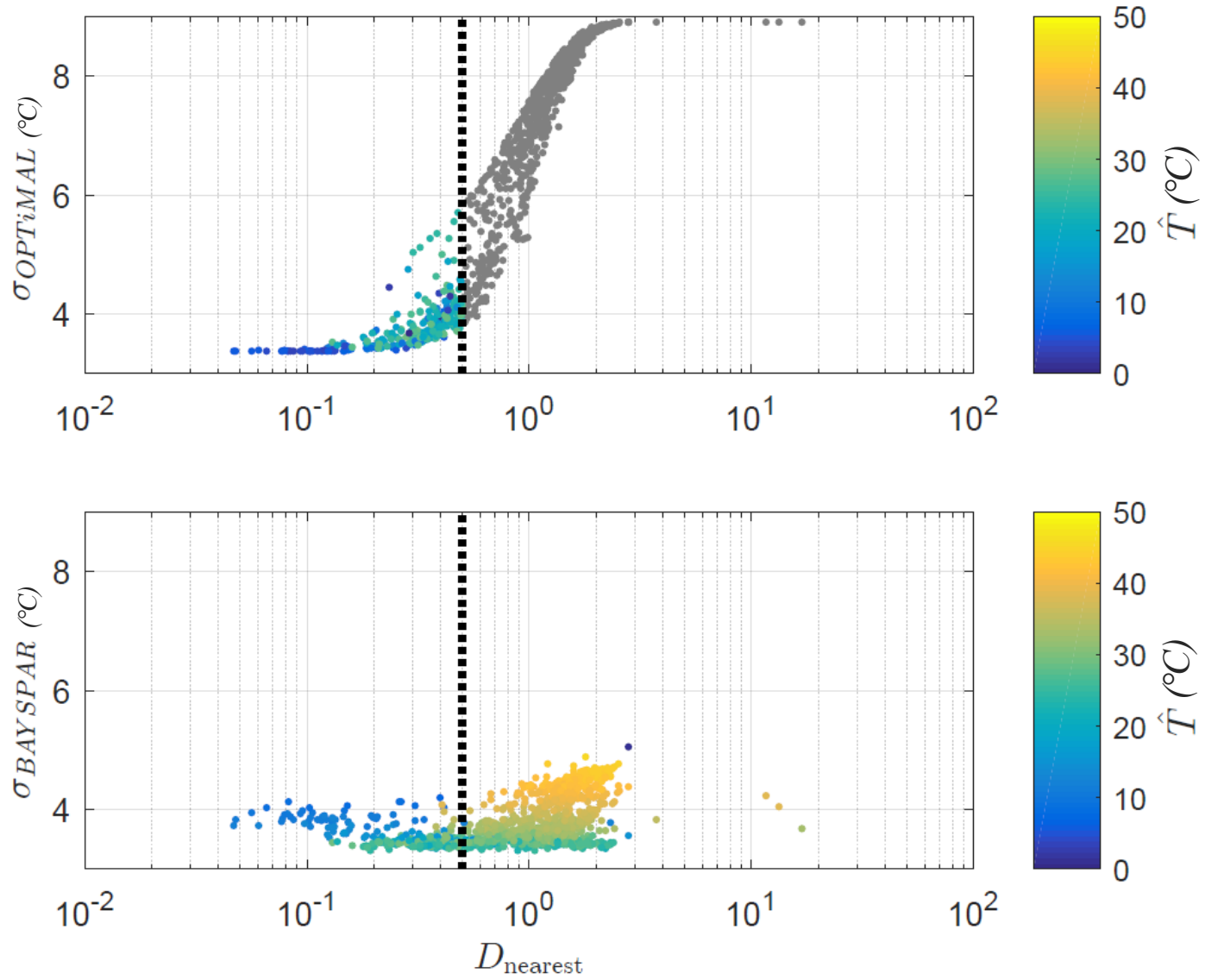


Figure 15

