

Response to Interactive comment on “OPTiMAL: A new machine learning approach for GDGT-based palaeothermometry” by Yvette L. Eley et al.

Dear Dr Reyes,

Thank you for your letter and for the reviews of our resubmission. We are pleased to note the “accept as is” judgement of Anonymous Reviewer 1, and their assessment of the manuscript as both excellent in its scientific significance and quality.

Here we address the points raised by new reviewer, Dr Tierney, which follow up on our response to her open discussion comments on our manuscript. Line numbers in our responses refer to the new “tracked changes” document.

Dr Tierney’s comments below are highlighted in bold with our responses in plain text.

Kind regards,

Dr Tom Dunkley Jones, on behalf of all the co-authors.

General comments:

This is a re-review of Eley et al., who present the OPTIMAL model for converting GDGT distributions to SST. In the first round of review I brought up a number of issues which came down to 1) tone of the paper, 2) concerns about over-fitting and/or mathematical clarity and 3) lack of applications. In reading the revised version I think that the tone is now much more appropriate and respectful (with a few small exceptions noted in the specific comments below). I am still concerned about overfitting. The authors argue that their in-sample validation technique (leave out a portion of the coretop data for validation iteratively) is proof enough that the method is not overfitted. This is not the case, because they are still validating on the training dataset (the coretop data). Out-of-sample validation is the gold standard. This is one reason why I suggested last time that they add a number of applications of OPTIMAL to downcore datasets and compare the results to previous calibrations. They now provide an example in Figure 14 (application to ODP 806 and ODP 850), but this Figure needs comparison with previous calibrations (i.e. BAYSPAR) so that the readers can see the differences. It is not an ideal choice b/c the alkenones at ODP 806 are saturated for much of the record. The authors need to add more examples beyond this one. I suggest 1 high-resolution late Quaternary record, and 1 deeper time record (Eocene perhaps). There should be data out there that have all of the fractional GDGTs, or else they can email the authors and get this. Note that Figures 13 and 14 are not really helpful - it is very difficult for readers to take much away from these. Paleoceanographers need to see applications. This is standard for new calibration papers and not, as the authors claim, out of the scope of the study. Only through this out-of-sample testing can we assess whether OPTIMAL is performing well.

We have included three downcore records as requested – two from the late Pleistocene, with direct comparisons against independent alkenone-based U^k_{37} SSTs, and one from the Eocene. These are shown in new Figures 16 and 17 and discussed in Section 4 of the revised manuscript.

Regarding mathematical clarity, the paper is still lacking on this front. I noticed there are more details in the SOM - this should all be in the main text, including Fig. S1 which is incredibly helpful. I suggest reorganizing their Methods section so that everything is much more clear and straightforward: describe the math, define the terms.

The revised manuscript is already extensive, running to over 13,000 words, and it is our judgment that the methodological details of the Forward Model – which is included as an exploratory way of trying to provide improved calibrations, but is not a part of the final OPTiMAL GPR model – is better placed within the Supplementary Materials. These methods are clearly laid out in the Supplement for those who want to follow this detail and explore this model further, whereas inclusion in the text would likely reduce the accessibility of the core message and purpose of this manuscript from the general reader of *Climates of the Past*.

It is also now evident that inversion of the GPR model is Bayesian. This should be made explicit in the main text along with a description of the priors used for SST.

Dr Tierney is confused – the inversion is only relevant to the forward modelling not the GPR model. Only the GPR model is involved in generating the OPTiMAL SST estimates. This supports our contention to keep the forward model detailed methodology in the SOM to prevent confusion with the primary GPR model (OPTiMAL).

It also seems that the model can in fact be extrapolated (?) but that the uncertainties will be high unless an informative prior is made.

The forward model can be extrapolated, and uncertainties will be large, unless “an informative prior is made”. That is exactly our point – that SST estimation out of range of the modern calibration, is largely controlled by the choice of “an informative prior”. We don’t think there are, yet, informative priors to provide confidence outside of the modern calibration range, but these might become available through either further culture or mesocosm studies, or paleo proxy-proxy calibrations.

Could you use constraints from other proxies here? That might provide some hope for deep-time users.

Yes, that could be attempted, but is beyond the scope of this paper. It would require a detailed assessment of all the uncertainty within other proxy systems for these constraints to be robustly integrated into a calibration model for GDGT paleothermometry. It would also involve the loss of independence between proxy systems.

All in all, it will be much easier to understand the steps that the authors took if the math is laid out in detail, along with (as I asked for last time) a description of what platforms and codes were used (Matlab, Python packages).

The underlying maths is laid out in detail within the paper; the full details of the coding are provided in GitHub: <https://github.com/carbonatefan/OPTiMAL>.

Also, my questions about collinearity (the fact that the GDGT predictors are not independent but in fact correlated with each other to a high degree) and possible regression dilution in the GPR model (or some other problem - maybe with the prior) were not explicitly addressed.

The machine learning methods we propose do not require inputs to be independent. Techniques such as Gaussian process regression learn from the data. If some inputs are uninformative (either because they are not correlated with the property we are trying to predict, or because they can be reconstructed from other inputs and therefore do not add value), they are recognised as such. They may not help, but they do not harm, either.

While we are, of course, sensitive to measurement noise in the independent variables, we do not perform linear regression and do not suffer from regression dilution in the usual sense of this term. Variability in both input and output variables is included in the uncertainty of the prediction.

GDGT abundance are definitely correlated, and that's expected — but not so correlated that you can effectively reduce the information to a single dimension, as TEX attempts to do. Here is the matrix of Spearman rank correlation coefficients (rows and columns are, in order, GDGT-0, GDGT-1, GDGT-2, GDGT-3, Cren, Cren'):

1.0000	-0.4165	-0.7898	-0.7986	-0.8937	-0.7266
-0.4165	1.0000	0.7807	0.4340	0.0364	0.6335
-0.7898	0.7807	1.0000	0.7443	0.4835	0.8844
-0.7986	0.4340	0.7443	1.0000	0.6690	0.6460
-0.8937	0.0364	0.4835	0.6690	1.0000	0.4399
-0.7266	0.6335	0.8844	0.6460	0.4399	1.0000

A lack of independence in the GDGT inputs is simply not an issue for Gaussian process regression (see C. E. Rasmussen & C. K. I. Williams, *Gaussian Processes for Machine Learning*, the MIT Press, 2006).

I welcome the inclusion of a comparison of Dnearest to BIT and MI, but deltaRI should be here too. In practice this is one of the most valuable metrics for identifying aberrant GDGT distributions.

Done.

Finally, it seems like the authors misunderstood my comment last time about providing some first-order constraints on their model. I was advocating for enforcing monotonicity, i.e. more rings = higher temperature. I was not advocating for a particular form for that (linear vs. non-linear). A monotonic constraint still seems appropriate to me. Couldn't this be built into OPTIMAL? Why would one not want to do so? The authors mention culture studies, most of which find more rings at higher T's. Even Bale et al. 2019 find more cren and cren' at higher T's, as expected. It seems like we have enough evidence in favor of monotonicity.

“Why would one not want to do so?” – because it would build-in unnecessary assumptions about the behaviour of the relationship into the calibration model. If this relationship is present in the data, our approach will find it; if it is not, then our approach won’t be biased by a prior user-assumption about the form of the temperature dependency.

If the authors can revise their paper to make the math clear, discuss any limitations as appropriate, and demonstrate that OPTIMAL can perform reasonably well on out-of-sample data (downcore time series) that would make me (and I suppose all readers) much more comfortable in terms of using this new approach.

With the new examples that we have provided, we hope we have provided some of this assurance, and look forward to Dr Tierney making use of OPTiMAL in near future.

Specific comments:

Line 131: "Like Qin et al. (2015), we note the non-linear nature of the individual experiments in Wuchter et al. (2004; see Fig. 5)." Need to clarify here that you mean Wuchter et al's Figure 5 and not your own. However this statement is deceiving. What Figure 5 in that paper shows is no response of TEX to SST b/t 5-15 degrees, and then a linear response thereafter in the series I incubation, plus no response in series II. It is not a non-linear (i.e. exponential) relationship. Combined with Schouten '07, the mesocosm data are linear b/t 10C and 40C (Schouten et al., 2007, Figure 4). I agree that this does not preclude a non-linear relationship in the real world - but it's important to not misstate what the data show.

We added a direct reference to the figure. We have not changed our text – if experiments show no or varied response of TEX₈₆ to temperature across the temperature range, as Tierney notes in this case, then we think it is valid to say that this is a “non-linear” response.

Line 150: "As such, these are collectively more representative of the community production contributing to samples in the global core-top TEX₈₆ calibrations of Kim et al., (2010) and BAYSPAR (Tierney & Tingley, 2014), which predominantly sample continental margin environments, rather than deep ocean / pelagic environments." I don't agree with this. How could a single-strain culture be more representative than environmental samples, which likely reflect multiple strains? There is no way that it could be. We don't know what strains contribute to the coretop dataset but it is certainly more diverse than just *N. maritimus*. Please remove.

This is a misreading of the sentence. The sentence meaning is that culture strains recovered from epi-continental shelf environments are more likely to represent the strains contributing to production of GDGTs going into core top calibrations, which are from the same epi-continental margin environments. However, this is a minor point, and to avoid confusion we have deleted these lines.

Line 164: "To use the responses of single, selected archaeal strains in culture to validate

a particular model of community-level responses to growth temperature is problematic even in the modern system (Elling et al., 2015)." I agree, and this directly contradicts the statement on Line 150.

See comment above – it was a misreading of this section, but it is now removed.

Line 195" "Powerful mathematical tools." I asked you to please delete this last time, as it is hyperbolic and non-specific. The tools here are no more powerful than other mathematical approaches. Replace with "an analysis of distance metrics", "machine learning" "GPR" or something similar.

Yes – and we disagree and retain. These machine learning tools presented are considerably more powerful at modelling complex relationships within multi-dimensional datasets than, say, linear regression. They are powerful tools.

Line 252: "For example, it may be that sea surface temperatures are very sensitive to one observable." I think you mean the reverse here - the observables are the GDGTs, and they might be more or less sensitive to SST. This paragraph would be easier to understand if you just use the term "each GDGT" vs. "observable".

Changed for clarity.

Eq. 7: The description of this distance metric is not totally clear. Can you clarify what x and y are here? Also as pointed out by Yi Ge, there are only five degrees of freedom, so doesn't this need to be adjusted accordingly?

Line 259: *"Thus, the normalised distance D between parameter data points x and y is:"*

Points x and y are samples within the GDGT space; $D_{x,y}$ is the normalised Euclidian distance between these two points. The sum in Eq. 7 should be taken over 6 parameters (GDGT-0, GDGT-1, GDGT-2, GDGT-3, cren, cren') and we apologize for a typo in the equation where the sum should be taken from 0 to 5, not 0 to 6. This has been corrected.

Lines 289-300: To what extent is the gain in information from the individual GDGTs due simply to the use of more than one parameters? As I pointed out in my first review, adding more parameters will improve RMSE but doesn't necessarily mean an improvement in skill. This should be addressed here. You can't say that the NN predictor "outperforms" TEX unless you rule out the effect of using more parameters, which incidentally are also not independent from each other.

There are no free parameters in the nearest neighbour approach: the predictor is just the temperature of the closest training set point in input space. And, indeed, as we say, the gain lies precisely in having more information available in 6 (5 independent) coordinates of the input point than in their one-dimensional combination, TEX.

Line 345: This R^2 of 0.83 is identical to what BAYSPAR can do with all of the data ($R^2 = 0.84$, Figure 5 in TT14). So it seems like both the random forest and BAYSPAR model perform similarly, even though BAYSPAR uses TEX86. Worth noting.

BAYSPAR computes R^2 by using the same data for calibration and validation, which allows for over-fitting and over-predicts R^2 .

Line 372: The authors haven't answered my question yet about the effects of collinearity (the correlation of the fractional GDGT abundances with each other). This seems like a good place to clarify this issue.

See points above, but have clarified in the main text (line 379-381): *“The accuracy of Gaussian process regression is not adversely affected by correlations between inputs (Rasmussen & Williams, 2006). Significantly correlated inputs that do not bring in new predictive power are appropriately down-weighted.”*

Line 425: The DCs look like they are showing evidence of the "horseshoe effect" common to standard PCA, in that branches B and C are part of the same horseshoe. This would signal a strong linear response in the multivariate data that isn't well-separated (?) When this occurs in standard PCA, the PCs are not interpretable. Can the authors comment here on this as it applies to the diffusion map technique?

PCA is fundamentally a linear transformation — a rotation of the coordinates. The so-called horseshoe effect is an indication of nonlinear correlations, which the PCA cannot account for. Diffusion maps are nonlinear dimensionality reduction tools, and so are generally capable of handling nonlinear correlations.

Line 477: "it does not account for the systematic uncertainty in the model when extrapolated beyond the calibration range" It does in that all possible regression lines are extrapolated, creating wide error bars and coalescing towards the prior if no other information is there. I would change to, "it still assumes a monotonic relationship between TEX and SST" which is a more accurate. However I don't think this is a bad assumption (see comment above).

Rephrased. Our point here is that BAYSPAR assumes that the linear model is intrinsically correct, and only uncertainties on the parameters within the model are accounted for, rather than any systematic uncertainty in the model itself.

Lines 481-486: So if I understand this correctly, the inversion of your model is Bayesian inference with priors placed on...what? This section needs some mathematical description to make this clear.

GPR need not be strictly Bayesian, though the calibration process can be viewed as Bayesian inference on the hyper-parameters of the model. For example, with the squared exponential kernel with automatic relevance determination that we use, these hyper-parameters are the kernel widths in each dimension. The priors referred to (e.g., in “Given that we have no

mechanistic model for the data generating process, we recommend the use of kernels which do not impose strong prior assumptions on the form of the GDGT-temperature relationship”) are not the classical Bayesian inference priors on model parameters of interest to the user that the Dr Tierney may be thinking of.

Line 488: "The clear pattern in the residuals does not necessarily indicate model misspecification, since no explicit noise model is specified for temperatures". I noted this last time - this looks like regression dilution. But if the model was specified as $TEX = f(SST) + \text{error}$ this shouldn't happen, unless...the prior is too tight? Please clarify what your priors are.

TEX is not used here at all. Regression dilution is simply not relevant as described above.

Line 521: This is the first mention of Dnearest, but it is not defined. Is Dnearest the same as D(x,y) described above? Please clarify.

No – as explained lines 523 - 525

Line 569: Bale et al. did observe an increase in cren and cren' in their culture experiments though, which suggests that there is some response of GDGTs to temperature, albeit not well-expressed in TEX₈₆.

Fine. We state there is no correlation with TEX₈₆.

Line 592: Figures 13 and 14 don't really communicate well how the Optimal model performs. It would be more useful to show a couple of time series and compare Optimal vs. BAYSPAR.

Have added timeseries as requested.

Line 610: It's not necessary to italicize "parametric" here. I would rephrase this to more accurately describe the difference b/t BAYSPAR and Optimal: "In contrast, while uncertainty bounds do increase when BAYSPAR is used to extrapolate beyond the modern calibration, they are not as large as Optimal because BAYSPAR still makes an assumption of a linear increase in SST at higher TEX values." As I said last time, BAYSPAR does account for model uncertainty, the issue is that the model we use is a linear one.

Parametric no longer italicized; rephrased in line with suggestion.

Line 615: deltaRI should be included here. It is arguably the most useful metric for ID'ing strange GDGT distributions.

Done.

Figure 16: I don't think this is an ideal application b/c UK37 is at its limit here at 806. Plus optimal should be compared to the previous calibrations (BAYSPAR, TEXH if you want, although the regression dilution in TEXH makes things hard to interpret).

Done

SOM: This should be in the main text, esp. Figure 1. Also the SOM appears to contain comments and incomplete references (refs).

See comments above about the forward model.

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 * Responses to Tierney's comments and additions are shown in blue text.

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 (Optimised Palaeothermometry from Tetraethers via MAchine Learning) 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 ~± 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).

80

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

82

83

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

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. (see Fig. 5 in Wuchter et al. 2004).
133 Moreover, the relatively lower Cren’ in these studies yield a very different intercept and slope compared to
134 core-top calibrations (e.g. Kim et al. 2010) making direct comparisons problematic.

135

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

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

170
171 Finally, the original uses of the TEX₈₆ proxy had a relatively poor representation of the true uncertainty
172 associated with palaeotemperature estimates, as they included no assessment of non-analogue behavior
173 relative to the modern core-top data. Instead, uncertainty was typically based on the residuals on the modern
174 calibration, with no reference to the relationship between GDGT distributions of an ancient sample and the
175 modern calibration data. An improved Bayesian uncertainty model “BAYSPAR” is now in widespread use
176 for SST estimation, which models TEX₈₆ to SSTs regression parameters, and associated uncertainty, as
177 spatially varying functions (Tierney and Tingley, 2015). The Bayesian approach, as with all approaches
178 based on the TEX₈₆ index, however, still does not include an uncertainty that reflects how well modelled
179 ancient GDGT assemblages are by the modern calibration – i.e. the degree to which they are non-analogue
180 - as it still functions on one-dimensional TEX₈₆ index values.

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

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

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

205

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

207

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

216

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

223

$$224 \quad \delta T = \sqrt{\frac{1}{N_v - 1} \sum_{k=1}^{N_v} (\hat{T}(x_k) - T(x_k))^2}$$

225 (Eq. 5)

226

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

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

236

237 *2.1 Nearest neighbours*

238

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

249

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

259

$$260 \quad D_{x,y}^2 \equiv \sum_{i=0}^5 \frac{(GDGT_i(x) - GDGT_i(y))^2}{var(GDGT_i)} \quad (Eq. 7)$$

261

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

264

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

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

276

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

285

286 *2.2 TEX₈₆ and Bayesian applications*

287

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

300

301 Indeed, this is what we find when we consider predictors of the form $\hat{T}_{1/TEX} = a + b/TEX_{86}$ and $\hat{T}_{TEXH} = c$
302 $+ d \log_{TEX86}$ (Liu et al., 2009; Kim et al., 2010), i.e., the established relationships between GDGT
303 distributions and SST. We fit the free parameters a , b , c , and d by minimising the sum of squares of the
304 residuals over the calibration data sets (least squares regression). We find that $\delta T_{1/TEX} = 6.1$ °C (note that
305 this is slightly better than using the fixed values of a and b from (Kim et al., 2010), which yield $\delta T_{1/TEX} =$
306 6.2 °C). We note that the corresponding R^2 value associated with these TEX_{86} based predictors is 0.64,
307 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
308 dataset based on Tierney and Tingley (2015), including data from the Red Sea (Kim et al. 2010).

309

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

321

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

330

331 *2.3 Machine learning Approaches – Random Forests*

332

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

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

346

347 *2.4 Gaussian Process Regression*

348

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

360

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

369 GDGT variant, then points with relatively nearby absolute input values in that coordinate are still distant.
370 We find that most GDGT parameters other than Cren are comparably useful in predicting temperature, with
371 GDGT-0 and GDGT-3 marginally the most informative. We considered whether interdependency of
372 percentage GDGT data could influence our calculations. Our analysis suggests that there are only five free
373 parameters. Machine learning tools should be able to pick up this correlation and effectively ignore one of
374 the parameters (or one parameter combination). For example, we do find that the GP emulator has a very
375 broad kernel in at least one dimension, signaling this. In principle, we could have considered only five of
376 six parameters. The smaller scale of some of the parameters is automatically accounted for by the trained
377 kernel size in GP regression, or by normalising to the appropriate dynamical range in our initial
378 investigation. [In short, the accuracy of Gaussian process regression is not adversely affected by correlations
379 between inputs \(Rasmussen & Williams, 2006\). Significantly correlated inputs that do not bring in new
380 predictive power are appropriately down-weighted.](#)

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

393 394 *2.5 Data Structure*

395
396 The random forest (Section 2.3) and GPR approaches (Section 2.4) are agnostic about any underlying bio-
397 physical model that might impart the observed temperature-dependence on GDGT relative abundances
398 produced by archaea. They are essentially optimized interpolation tools for mapping correlations between
399 temperature and GDGT abundances within the range of the modern calibration data set; they can make no
400 sensible inference about the behavior of this relationship outside of the range of this training data. To move
401 from interpolation within, to extrapolation beyond, the modern calibration requires an understanding of,
402 and model for, the temperature-dependence of GDGT production. To explore these relationships and the

403 extent to which the ancient and modern data reside in a coherent relationship within GDGT space, we
404 employed two forms of dimensionality reduction to enable visualisation of the data in two or three
405 dimensions. The fundamental point is that if temperature is the dominant control, all of the data should lie
406 approximately on a one-dimensional curve in GDGT space, and the arclength along this curve should
407 correspond to temperature; we will revisit this point below.

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

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

432
433 The relationship between the first diffusion component and TEX_{86} for all data is shown in Figure 8. There
434 is a clear correlation, despite the presence of some outlying Cretaceous points, some of which are not shown
435 because they lie so far outside the majority data range within this projection. This suggests that TEX_{86} is,
436 in one sense, a natural one-dimensional representation of the data. We also plot the first diffusion

437 component for the modern data as a function of temperature (Figure 9). We see a similar pattern emerging
438 to that displayed by TEX₈₆ - there is little sensitivity to temperature below 15 °C, and between ~20 and 25
439 °C. An interesting avenue for future research might be to explore the temperature-GDGT system from a
440 dynamical systems perspective, i.e. use simple mechanistic mathematical models to explore the
441 temperature-dependence of steady-state GDGT distributions. It may be that such models suggest that only
442 a few steady-states exist, and that temperature is a bifurcation parameter, i.e. it controls the switch between
443 the steady states. Note also the downward slope in the residual pattern in Figure 4 between 0 and 15-17
444 degrees celsius, and again at higher temperatures. This pattern is consistent with predictions that are biased
445 towards the centre of each 'cluster', i.e. a system which is not very sensitive to temperature, but can
446 distinguish between high and low temperatures reasonably well. This observation also links to recent culture
447 studies (Elling et al., 2015) and Pliocene-Pleistocene sapropel data (Polik et al., 2018), which support the
448 existence of discrete populations with unique GDGT-temperature relationships and that temporal changes
449 in population over time can drive changes in TEX₈₆.

450

451 *2.6 Forward Modelling*

452

453 Based on the analysis of the combined modern and ancient data structure outlined above, there appears to
454 be some consistency to underlying trends in the overall variance of GDGT relative abundances. These
455 trends provide some hope that models of this variance, and its relationship to sea surface temperature, within
456 the modern dataset could be developed to predict ancient SSTs. TEX₈₆ and BAYSPAR are such models,
457 but they are limited by, first, the reduction of six-dimensional GDGT space to a one-dimensional index;
458 and second, by an *ad hoc* model choice – linear, exponential – that does not account for uncertainty in
459 model fit to the modern calibration data, and the resultant uncertainty in the estimation of ancient SSTs
460 relating to model choice. To overcome these issues, we develop a forward model based on a multi-output
461 Gaussian Process (Alvarez et al., 2012), which models GDGT compositions as functions of temperature,
462 accounting for correlations between GDGT measurements. This model is then inverted to obtain
463 temperatures which are compatible with a measured GDGT composition. In simple terms, we posit that a
464 measured GDGT composition is generated by some unknown function of temperature and corrupted by
465 noise, which may be due to measurement error or some unmodelled particularity of the environment in
466 which the sample was generated. We proceed by defining a large (in this case infinite) set of functions of
467 temperature to explore and compare them to the available data, throwing away those functions which do
468 not adequately fit the data. This means, of course, that the behaviour of the functions we accept is allowed
469 to vary more widely outside the range of the modern data than within it. With no mechanistic underpinning,

470 choosing only one function (such as the inverse of TEX_{86}) based on how well it fits the modern data grossly
471 underestimates our uncertainty about temperature where no modern analogue is available.

472

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

490

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

498

499 An advantage of the forward modelling approach is that the inversion can incorporate substantive prior
500 information about temperatures for individual data points. In particular, other proxy systems can be used to
501 elicit prior distributions over temperatures to constrain GDGT-based predictions, particularly when
502 attempting to reconstruct ancient climates with no modern analogue in GDGT-space. We emphasise that
503 outside the range of the modern data, the utility of the models is almost solely due to the prior information

504 included in the reconstruction. At present, the only priors being used in the forward model prescribe a
505 reasonable upper limit and lower limit on temperatures (see Supplementary Information). The only way to
506 improve these reconstructions will be for future iterations to incorporate prior information from other
507 proxies. It is worth noting that the predictive uncertainty, while reasonably well-described by the standard
508 deviation in cases where ancient data lie quite close to the modern data in GDGT space, can be highly
509 multimodal (Fig. 11). This is the case when estimates are significantly outside of the modern calibration
510 dataset, such as low latitude data in the Cretaceous, or where there is considerable scatter in the modern
511 calibration data, for example in the low temperature range (<5 °C).

512

513 3. Non-analogue behavior and Extrapolation

514

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

538

539 Where ancient GDGT distributions lie far from the modern calibration data set ($D_{\text{nearest}} > 0.5$), we argue that
540 there is no suitable set of modern analogue GDGT distributions from which to infer growth temperatures
541 for this ancient GDGT distribution. Both the GPR and Fwd models revert to imposed priors once the
542 distance from the modern calibration dataset increases. We propose that this is more rigorous and justified
543 model behavior than extrapolation of TEX₈₆ or BAYSPAR predictors to non-analogue samples far from
544 the modern calibration data. As a result, the predictive models can only be applied to a subset of the Eocene
545 and Cretaceous data. We also note that there are two broad, non-mutually-exclusive categories of samples
546 that lie far from the modern calibration dataset ($D_{\text{nearest}} > 0.5$), the first are samples that seem to lie ‘beyond’
547 the temperature-GDGT calibration relationship, likely with (unconstrained) GDGT formation temperatures
548 higher than the modern core-top calibrations; the second are samples with anomalous GDGT distributions
549 lying on the margins of, or far away from the main GDGT clustering in 6-dimensional space (see outliers
550 in Fig. 8).

551

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

572 Nitrosopumilales (that contains many mesophilic marine strains). They found no correlation between
573 TEX₈₆ calibrations (either the Kim et al., core-top or Wuchter et al. 2004 and Schouten et al., 2008
574 mesocosm calibrations) with membrane lipid composition at different growth temperatures (37°C, 46°C,
575 and 50°C) and found that phylogeny generally seems to have a stronger influence on GDGT distribution
576 than temperature. In view of these existing data, we see no robust justification at present for the
577 extrapolation of modern core-top calibration data sets into the unknown above 30 °C, although the
578 coherent patterns apparent across GDGT space, between modern, Eocene and Cretaceous data (Figure 7),
579 do provide some grounds for hope that the extension of GDGT palaeothermometry beyond 30°C might be
580 possible in future.

581

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

583

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

592

593 Following Tierney and Tingley (2014) we use a reduced calibration data set, with the exclusion of Arctic
594 data with observed SSTs less than 3°C (“NoNorth / TT13” of Tierney and Tingley (2014)) but with the
595 inclusion of additional core top data from Seki et al. (2014). Full details of this calibration dataset are
596 provided in the Supplementary Information; to distinguish from the original OPTiMAL calibration data,
597 which included the Arctic data <3°C, we refer to the original data as “Op1” and the new calibration dataset
598 as “Op3”. An “Op2” is also available, which is the same as Op1 except that it excludes the Seki et al. (2014)
599 data. In sensitivity tests to a range of applications across Quaternary and deep-time datasets, calibration
600 Op1 and Op2 performed in almost identical fashion. The performance of Op1 and Op3 were very similar
601 in most applications, except in applications to the paleo-Arctic (see below), where the inclusion of modern
602 Arctic calibration data (Op1) provided closer calibration constraints to the paleo-data. Although
603 superficially this may be regarded as beneficial, in these instances the paleo-data have previously been
604 rejected because of a potential bias by non-marine inputs indicated by high BIT indices (Sluijs et al. 2020).
605 In this case, either the modern Arctic calibration data is impacted by similar non-thermal processes,

606 generating unusual GDGT abundance patterns, which are not appropriate to use for SST calibration, or,
607 there could be some consistency between the modern and ancient GDGT production by marine archaea in
608 the Arctic which may help in the understanding of GDGT-based paleothermometry in this unusual
609 environment (Sluijs et al. 2020). The D_{nearest} methodology may prove useful in quantifying analogue and
610 non-analogue behavior through time in such conditions. For the purposes of this study, however, we take
611 the conservative approach, and one that maintains a more consistent calibration basis with BAYSPAR, by
612 using OPTiMAL calibration Op3 in the remainder of this discussion, and recommend its use in future
613 applications of OPTiMAL.

614

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

635

636 We also provide an initial assessment of the inter-relationship between standard screening indices and
637 D_{nearest} , for the Eocene and Cretaceous compilations where the data are available to calculate these measures
638 (Figure 15). For ease of comparison between Eocene and Cretaceous datasets and visualization of the
639 majority of the data, extreme outliers ($D_{\text{nearest}} > 4.0$) are not shown. The metrics include the BIT index

640 (Hopmans et al., 2004; Weijers et al., 2006), the Methane Index (MI; Zhang et al., 2011), the deviation
641 between TEX_{86} and the Ring Index (ΔRI ; Zhang et al., 2016) and the %GDGT-0 (Blaga et al., 2009;
642 Sinninghe Damsté et al., 2012). The standard screening levels for each of these metrics, as used in previous
643 paleo-compilations (O'Brien et al. 2017), are shown in the blue shaded areas on Figure 15 ($\text{BIT} > 0.5$; MI
644 > 0.5 ; $\Delta\text{RI} > 0.3$; %GDGT-0 $> 67\%$) – data points within these areas fail the standard screening. Also shown
645 on Figure 15 is the region where data pass our D_{nearest} screening requirement (grey shaded vertical region).
646 In nearly all cases GDGT assemblages that fail these traditional screening tests also have D_{nearest} values that
647 exceed 0.5 – i.e. “abnormal” GDGT assemblages are well screened D_{nearest} . The main exception to this is
648 the BIT index in the Eocene data set, where 15 samples have high BIT values (>0.5) but have GDGT
649 assemblages that are close to modern analogues in the calibration dataset ($D_{\text{nearest}} < 0.5$). Of these samples,
650 9 are from the Arctic Ocean between the PETM and ETM2, an interval noted for its relatively high BIT
651 index values (Sluijs et al. 2020), 3 are from the Eocene-Oligocene transition of ODP Site 1218 (eastern
652 Equatorial Pacific) (Liu et al. 2009), 2 are from the middle Eocene of Seymour Island (Douglas et al. 2014),
653 and 1 is from the late Eocene of DSDP Site 511, which has been already noted as an individual sample with
654 anomalous high BIT in this dataset (Liu et al. 2009; Inglis et al. 2015). Although high BIT at ODP Site
655 1218 has been inferred to represent “relatively high terrestrial input” (Inglis et al. 2015) this seems unusual
656 for a fully pelagic site situated on oceanic crust >3000 km away from the nearest continental landmass.
657 Interpreting high BIT values as exclusively caused by terrestrial organic components appears problematic
658 in this instance, especially as $D_{\text{nearest}} < 0.5$ give some assurance that these GDGT assemblages from ODP
659 Site 1218 are well-modelled by the modern calibration dataset. GDGT assemblages from Seymour Island
660 associated with high BIT values (>0.4) appear to have an impact on the $\text{TEX}_{86}^{\text{H}}$ SST proxy (Inglis et al.
661 2015), but the 2 samples that fail BIT (>0.5) but pass D_{nearest} (<0.5) give OPTiMAL SSTs consistent (5-
662 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
663 relationship between D_{nearest} and BIT suggests that BIT is not always closely coupled to GDGT assemblages
664 that are strongly divergent from the modern calibration dataset.

665

666 With respect to the other screening indices there are clear indications that increased distance from the
667 modern calibration (increased D_{nearest}) is associated with a trend towards the “thresholds of failure” in the
668 screening indices. This pattern is most clear with the ΔRI in both the Cretaceous and the Eocene data, as
669 increasing numbers of samples fail ΔRI as D_{nearest} increases. This supports ΔRI as a robust methodology for
670 identifying samples that strongly diverge from the expected temperature-dependence of GDGT
671 assemblages as modelled by TEX_{86} in the modern calibration dataset. There are, however, samples that pass
672 $D_{\text{nearest}} < 0.5$ but fail ΔRI in both the Eocene and Cretaceous datasets – these must have “near neighbours”
673 in the modern calibration data, but yet have a temperature-sensitivity that is less well-modelled by TEX_{86}

674 (divergence between RI and TEX_{86}). Conversely there are many Eocene and Cretaceous data points with
675 $\Delta\text{RI} < 0.3$, but which fail $D_{\text{nearest}} (>0.5)$. These data most likely represent GDGT assemblages formed at
676 high temperatures, beyond the range of the modern calibration data.

677
678
679 ~~These plots show little relationship between the BIT and MI screening indices and D_{nearest} values. Whilst in~~
680 ~~the Eocene, samples with the highest D_{nearest} values (>3) also show very elevated BIT values (>0.8), in the~~
681 ~~Cretaceous the exceptionally anomalous assemblages (D_{nearest} values >100) are not anomalous in either BIT~~
682 ~~or MI. Conversely, in the Eocene there are many samples with relatively high BIT (>0.3) that are below the~~
683 ~~D_{nearest} threshold of 0.5. The behaviour of these systems needs to be examined in detail in future studies, but~~
684 ~~a conservative approach would be to apply all three screening indices (BIT, MI and D_{nearest}) to have the~~
685 ~~most confidence in resulting temperature estimates.~~

686
687 To investigate these behaviours requires the publication of the full range GDGT abundance data. Whilst
688 key compilations of Eocene and Cretaceous GDGT data have strongly encouraged the release of such
689 datasets (Lunt et al. 2012; Dunkley Jones et al. 2013; Inglis et al. 2015; O'Brien et al. 2017), most Neogene
690 studies only publish TEX_{86} values. Without full GDGT assemblage data neither OPTiMAL nor other
691 detailed assessments of GDGT behaviour and type can be made, and we would strongly encourage authors,
692 reviewers and editors to ensure the publication of full GDGT assemblages in future.

693
694 Finally, to test the behavior of OPTiMAL within established SST time series, we provide three examples
695 two from the late Pleistocene to Holocene (Figure 16) and one from the Eocene (Figures 17 and 18). ~~where~~
696 ~~full GDGT assemblage data were made available, and there are comparison alkenone based $U^{k'}_{37}$ data from~~
697 ~~the same sampling location — ODP 806 and ODP 850, respectively in the West and Eastern Equatorial~~
698 ~~Pacific (Figure 16; Zhang et al. 2014). For the Pleistocene to Holocene examples OPTiMAL SSTs are~~
699 ~~shown against estimates from BAYSPAR and the alkenone-based $U^{k'}_{37}$ temperature proxy. The first of~~
700 ~~these timeseries is from GeoB 7702-3 in the Eastern Mediterranean and spans the last 26 kyr, including~~
701 ~~data spanning Termination I (Castañeda et al., 2010). The second is from ODP Site 1146 in the South China~~
702 ~~Sea and spans the last 350 kyr (Thomas et al. 2014). In both records the long-term dynamics are consistent~~
703 ~~between the independent $U^{k'}_{37}$ SST proxy and both BAYSPAR and OPTiMAL. In the Eastern~~
704 ~~Mediterranean OPTiMAL SSTs are slightly cooler in the glacial and warmer in the Holocene than the other~~
705 ~~proxies. In the South China Sea, OPTiMAL is again cooler than BAYSPAR during glacial intervals, but at~~
706 ~~this location is in closer agreement than BAYSPAR with the $U^{k'}_{37}$ SST proxy through most of the record.~~

707 In both these examples, we show the 5th and 95th percentiles for OPTiMAL and those reported by the
708 BAYSPAR methodology.

709

710 The final example is from the latest Paleocene to early Eocene of IODP Expedition 302 Hole 4A on
711 Lomonosov Ridge (Sluijs et al. 2006; Sluijs et al. 2009; Sluijs et al. 2020). This site is useful as it has been
712 the focus of detailed reassessment and reanalysis, using most of the available screening methodologies to
713 detect aberrant GDGT assemblages (Sluijs et al. 2020). Here we use this recently published data to compare
714 the new D_{nearest} screening metric against multiple other screening protocols (Figure 17). We also show both
715 D_{nearest} values and OPTiMAL SST estimates for two models – one with modern Arctic data with SST < 3°C
716 included in the calibration (OPTiMAL_{Arctic}; equivalent to calibration dataset Op1 first present by Eley et al.
717 2019) and one with this data excluded (OPTiMAL_{noArctic}; equivalent to the new calibration dataset Op3). It
718 is clear from the pattern of D_{nearest} for these two options, that the inclusion of modern Arctic data provides
719 more calibration data that are closer to the Eocene paleo-Arctic, to the extent that substantially more
720 samples pass the $D_{\text{nearest}} < 0.5$ constraint, especially in pre-ETM2 interval from ~372 to 376 mcd. This
721 interval contains, however, samples with the highest BIT values of the succession (> 0.4), and elevated ΔRI
722 (> 0.3). With these other “warning signs” concerning the reliability of GDGT assemblages for SST
723 estimation in this interval, the relatively low D_{nearest} values are most likely to represent some similarity in
724 the non-thermal controls on GDGT assemblages between the modern and paleo-Arctic. More work needs
725 to be done to constrain the reliability of temperature-dependence and archaeal GDGT production in these
726 modern high latitude systems before we can have confidence in their inclusion in calibration datasets for
727 paleo-SST estimation. It is on the basis that we recommend users of OPTiMAL use the the “noArctic”
728 (Op3) calibration for the time being. The OPTiMAL methodology does, however, offer a simple means to
729 integrate new robust calibration data, and a method to explore the distance relationships between modern
730 and ancient GDGT production.

731

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

741
742 ~~are not at the limit of alkenone~~ saturation in ODP 806, OPTiMAL and $U^{k^2}_{37}$ agree well in the Plio-
743 Pleistocene. In ODP 850, there is a strong agreement in the reproduction of a long term late Miocene to
744 Recent cooling trend in both $U^{k^2}_{37}$ and OPTiMAL of $\sim 5^\circ\text{C}$. There is, however a consistent $\sim 2^\circ$ offset between
745 cooler OPTiMAL and warmer $U^{k^2}_{37}$ temperatures at this location. This offset is very similar in magnitude
746 and direction as that between TEX_{86}^H and $U^{k^2}_{37}$ used in the original study, and is more likely due to an
747 inherent feature (seasonality or depth) of archaeal versus eukaryotic production at this site (Zhang et al.
748 2014).

749 750 **5. Conclusions**

751
752 Although the fundamental issue of non-analogue behaviour is a key problem for GDGT-temperature
753 estimation, it has an undue impact on the community's general confidence in this method. In part, this is
754 because these issues have not been clearly stated and circumscribed - rather they have been allowed to erode
755 confidence in the GDGT-based methodology through the use of GDGT-based palaeothermometry far
756 outside the modern constraints on the behavior of this system. The use of GDGT abundances to estimate
757 temperatures in clearly non-analogue conditions is, at present, problematic on the basis of the available
758 calibration constraints or a good understanding of underlying biophysical models. We hope that this study
759 prompts further investigations that will improve these constraints for the use of GDGTs in deep-time
760 paleoclimate studies, where they clearly have substantial potential as temperature proxies. Temperature
761 estimates based on fossil GDGT assemblages that are within range of, or similar to, modern GDGT
762 calibration data, do, however, rest on a strong, underlying temperature-dependence observed in the
763 empirical data. With no effective means of separating the "good from the bad" can lead to either false
764 confidence and inappropriate inferences in non-analogue conditions, or a false pessimism when ancient
765 samples are actually well constrained by modern core-top assemblages.

766
767 In this study, we apply modern machine-learning tools, including Gaussian Process Emulators and forward
768 modelling, to improve temperature estimation and the representation of uncertainty in GDGT-based SST
769 reconstructions. Using our new nearest neighbour test, we demonstrate that $>60\%$ of Eocene, and $>90\%$ of
770 Cretaceous, fossil GDGT distribution patterns differ so significantly from modern as to call into question
771 SSTs derived from these assemblages. For data that does show sufficient similarity to modern, we present
772 OPTiMAL, a new multi-dimensional Gaussian Process Regression tool which uses all six GDGTs (GDGT-
773 0, -1, -2, -3, Cren and Cren') to generate an SST estimate with associated uncertainty. The key advantages
774 of the OPTiMAL approach are: 1) that these uncertainty estimates are intrinsically linked to the strength of

775 the relationship between the fossil GDGT distributions and the modern calibration data set, and 2) by
776 considering all GDGT compounds in a multi-dimensional regression model it avoids the dimensionality
777 reduction and loss of information that takes place when calibrating single parameters (TEX_{86}) to
778 temperature. The methods presented above make very few assumptions about the data. We argue that such
779 methods are appropriate with the current absence of any reasonable mechanistic model for the data
780 generating process, in that they reflect model uncertainty in a natural way. Finally, we note the potential
781 for multi-proxy machine learning approaches, synthesising data from other palaeothermometers with
782 independent uncertainties and biases, to improve calibration of ancient GDGT-derived SST reconstructions.

783
784

785 **Acknowledgements:**

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

796
797

798 **Figure Captions:**

799

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

802

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

805

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

809

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

811

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

814

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

817

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

822

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

826

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

828

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

830

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

834

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

838

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

844

845 **Figure 14.** Inter-comparison of temperature estimates and standard errors (y-axis) for compiled Eocene
846 and Cretaceous data calculated using OPTiMAL (top) and BAYSPAR (bottom). Greyed out data fails the
847 $D_{nearest}$ test (>0.5), and the colour scaling reflects $D_{nearest}$ values for those datapoints that pass. The black
848 dashed line shows the $D_{nearest}$ threshold (>0.5).

849
850 **Figure 15.** Comparison of $D_{nearest}$ against standard screening indices, BIT and MI index, ΔRI and
851 %GDGT-O for the Eocene (Inglis et al., 2015) and Cretaceous (O'Brien et al., 2017) datasets. Blue
852 shaded regions show the standard cut-off points for these indices (see text); grey shaded region highlights
853 data that are below the $D_{nearest}$ threshold of 0.5. The outlined black box is the region of data that fails
854 traditional screening indices but passes $D_{nearest}$ (<0.5).

855
856 **Figure 16.** Late Pleistocene to Holocene GDGT-derived OPTiMAL palaeotemperatures compared to
857 BAYSPAR and $U^{k'}_{37}$ SSTs. Shaded regions represent reported 5th and 95th percentile confidence intervals.
858 Top panel - Eastern Mediterranean data from core GeoB 7702-3 (Castaneda et al. 2010); bottom panel –
859 South China Sea data from ODP Site 1146 (Thomas et al. 2014).

860
861 **Figure 17.** Comparison of GDGT screening indices, TEX_{86}^H , BAYSPAR and OPTiMAL SSTs from the
862 Eocene Arctic Site IODP Expedition 302 Hole 4A. Data and figures modified from the most recent
863 reassessment by Sluijs et al. (2020).

864

865 **References:**

866 Aitchison, J.: The Statistical Analysis of Compositional Data. J. R. Stat. Soc. Series B Stat. Methodol. 44,
867 139–160, 1982.

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

869 Aitchison, J., Greenacre, M.: Biplots of compositional data. J. R. Stat. Soc. Ser. C Appl. Stat. 51, 375–
870 392, 2002.

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

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

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

879 Bijl, P. K., Bendle, J.A.P., Bohaty, S.M., Pross, J., Schouten, S., Tauxe, L., Stickley, C., McKay, R.M.,
880 Röhl, U., Olney, M., Sluijs, A., Escutia, C., Brinkhuis, H. and Expedition 318 Scientists.: Eocene

881 cooling linked to early flow across the Tasmanian Gateway. *Proc. Natl. Acad. Sci. U.S.A.*, 110,
882 9645–9650, 2013.

883 Brassell, S. C.: Climatic influences on the Paleogene evolution of alkenones, *Paleoceanography*, 29, 255-
884 272, doi:10.1002/2013PA002576, 2014.

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

890 Cadillo-Quiroz, H., Didelot, X., Held, N. L., Herrera, A., Darling, A., Reno, M. L., Krause, D. J.,
891 Whitaker, R. J.: Patterns of gene flow define species of thermophilic Archaea, *PLOS*, Cadillo-
892 Quiroz, H. et al. (2012) *PLOS Biology*, <https://doi.org/10.1371/journal.pbio.1001265>, 2012.
893

894 Castañeda, I. S., E. Schefuß, J. Pätzold, J. S. Sinninghe Damsté, S. Weldeab, and S. Schouten.:
895 Millennial-scale sea surface temperature changes in the eastern Mediterranean (Nile River Delta
896 region) over the last 27,000 years, *Paleoceanography*, 25, PA1208, doi:10.1029/2009PA001740,
897 2010.
898

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

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

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

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

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

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

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

918 Eley, Y. L., Thompson, W., Greene, S. E., Mandel, I., Edgar, K., Bendle, J. A., and Dunkley Jones, T.:
919 OPTiMAL: A new machine learning approach for GDGT-based palaeothermometry, *Clim. Past*
920 *Discuss.*, <https://doi.org/10.5194/cp-2019-60>, in review, 2019.

- 921 Elling, F. J., Könneke, M., Lipp, J. S., Becker, K. W., Gagen, E. J., and Hinrichs, K.-U.: Effects of growth
922 phase on the membrane lipid composition of the thaumarchaeon *Nitrosopumilus maritimus* and
923 their implications for archaeal lipid 20 distributions in the marine environment, *Geochimica et*
924 *Cosmochimica Acta*, 141, 579-597, 2014.
- 925 Elling, F. J., Könneke, M., Mußmann, M., Greve, A., and Hinrichs, K.-U.: Influence of temperature, pH,
926 and salinity on membrane lipid composition and TEX₈₆ of marine planktonic thaumarchaeal
927 isolates, *Geochimica et Cosmochimica Acta*, 171, 238-255, 2015.
- 928 Elling, F.J., Konnecke, M., Nicol, G. W., Stieglmeier, M., Bayer, B., Spieck, E., de la Torre, J. R.,
929 Becker, K. W., Thomm, M. Prosser, J. I., Herndl, G., Schleper, C., Hinrichs, K-U.:
930 Chemotaxonomic characterisation of the thaumarchaeal lipidome, *Environmental Microbiology* 19,
931 2681–2700 , 2017.
- 932
- 933 Filzmoser, P., Hron, K., Reimann, C.: Principal component analysis for compositional data with outliers.
934 *Environmetrics* 20, 621–632, 2009a.
- 935 Filzmoser, P., Hron, K., Reimann, C., Garrett, R.: Robust factor analysis for compositional data. *Comput.*
936 *Geosci.* 35, 1854–1861, 2009b.
- 937 Filzmoser, P., Hron, K., Reimann, C.: Interpretation of multivariate outliers for compositional data.
938 *Comput. Geosci.* 39, 77–85, 2012.
- 939 Gliozzi, A., Paoli, G., De Rosa, M., Gambacorta, A.: Effect of isoprenoid cyclization on the transition
940 temperature of lipids in thermophilic archaeobacteria, *Biochimica et Biophysica Acta (BBA)*
941 *Biomembranes*, 735, 234 – 242, 1983.
- 942 Haghverdi, L., Buettner, F., Theis, F.J: Diffusion maps for high-dimensional single-cell analysis of
943 differentiation data. *Bioinformatics* 31, 2989–2998, 2015.
- 944 Haghverdi, L., Büttner, M., Wolf, F.A., Buettner, F., Theis, F.J.: Diffusion pseudotime robustly
945 reconstructs lineage branching. *Nat. Methods* 13, 845–848, 2016.
- 946 Haslett, J., Whitley, M., Bhattacharya, S., Salter-Townshend, M., Wilson, S.P., Allen, J.R.M., Huntley, B.,
947 Mitchell, F.J.G.: Bayesian palaeoclimate reconstruction. *J Royal Statistical Soc A* 169, 395–438;,
948 2006.
- 949 Herfort, L., Schouten, S., Boon, J. P., and Sinninghe Damsté, J. S.: Application of the TEX₈₆ temperature
950 proxy to the southern North Sea, *Organic Geochemistry*, 37, 1715-1726, 2006.
- 951 Hertzberg, J. E., Schmidt, M. W., Bianchi, T. S., Smith, R. K., Shields, M. R., & Marcantonio, F.:
952 Comparison of eastern tropical Pacific TEX₈₆ and Globigerinoides ruber Mg/Ca derived sea surface
953 temperatures: Insights from the Holocene and Last Glacial Maximum. *Earth and Planetary Science*
954 *Letters*, 434, 320–332, 2016.
- 955 Hollis, C. J., Taylor, K. W. R., Handley, L., Pancost, R. D., Huber, M., Creech, J. B., Hines, B. R.,
956 Crouch, E. M., Morgans, 25 H. E. G., Crampton, J. S., Gibbs, S., Pearson, P. N., and Zachos, J. C.:
957 Early Paleogene temperature history of the Southwest Pacific Ocean: Reconciling proxies and
958 models, *Earth and Planetary Science Letters*, 349–350, 53-66, 2012.
- 959 Hollis, C. J., Dunkley Jones, T., Anagnostou, E., Bijl, P. K., Cramwinckel, M. J., Cui, Y., Dickens, G. R.,
960 Edgar, K. M., Eley, Y., Evans, D., Foster, G. L., Frieling, J., Inglis, G. N., Kennedy, E. M.,

961 Kozdon, R., Lauretano, V., Lear, C. H., Littler, K., Meckler, N., Naafs, B. D. A., Pälike, H.,
962 Pancost, R. D., Pearson, P., Royer, D. L., Salzmann, U., Schubert, B., Seebeck, H., Sluijs, A.,
963 Speijer, R., Stassen, P., Tierney, J., Tripathi, A., Wade, B., Westerhold, T., Witkowski, C., Zachos,
964 J. C., Zhang, Y. G., Huber, M., and Lunt, D. J.: The DeepMIP contribution to PMIP4:
965 methodologies for selection, compilation and analysis of latest Paleocene and early Eocene climate
966 proxy data, incorporating version 0.1 of the DeepMIP database, *Geosci. Model Dev. Discuss.*,
967 <https://doi.org/10.5194/gmd-2018-309>, in review, 2019.

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

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

974 Huguet C, Kim J-H, Sinninghe Damsté J.S., Schouten S: Reconstruction of sea surface temperature
975 variations in the Arabian Sea over the last 23 kyr using organic proxies (TEX₈₆ and UK 0 37).
976 *Paleoceanography* 21(3): PA3003, 2006.

977 Hurley, S. J., Elling, F. J., Könneke, M., Buchwald, C., Wankel, S. D., Santoro, A. E., Lipp, J.S.,
978 Hinrichs, K., Pearson, A.: Influence of ammonia oxidation rate on thaumarchaeal lipid composition
979 and the TEX₈₆ temperature proxy. *Proceedings of the National Academy of Sciences*, 113, 7762–
980 7767, 2016.

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

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

987 Kim, J.-H., Schouten, S., Hopmans, E. C., Donner, B., and Sinninghe Damsté, J. S.: Global sediment
988 core-top calibration of the TEX₈₆ paleothermometer in the ocean, *Geochimica et Cosmochimica*
989 *Acta*, 72, 1154-1173, 2008.

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

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

997 Liu, X-L., Zhu, C., Wakeham, S.G., Hinrichs, K-U.: In situ production of branched glycerol dialkyl
998 glycerol tetraethers in anoxic marine water columns, *Marine Chemistry*, 166, 1 – 8, 2014.

999 Lunt, D. J., Dunkley Jones, T., Heinemann, M., Huber, M., LeGrande, A., Winguth, A., Loptson, C.,
1000 Marotzke, J., Tindall, J., 15 Valdes, P., Winguth, C.: A model-data comparison for a multi-model

- 1001 ensemble of early Eocene atmosphere-ocean simulations: EoMIP, *Clim. Past Discuss.*, 8, 1229-
1002 1273, 2012.
- 1003 Mentch, L., Hooker, G.: Quantifying Uncertainty in Random Forests via Confidence Intervals and
1004 Hypothesis Tests. *Journal of Machine Learning Research*, 17, 1-41, 2016.
- 1005 O'Brien, C. L., Robinson, S. A., Pancost, R. D., Sinninghe Damsté, J. S., Schouten, S., Lunt, D. J.,
1006 Alsenz, H., Bornemann, J. A., Bottini, C., Brassell, S. C., Farnsworth, A., Forster, A., Huber, B.
1007 T., Inglis, G. N., Jenkyns, H. C., Linnert, C., Littler, K., Markwick, P., McAnena, A., Mutterlose,
1008 J., Naafs, B. D. A., Püttmann, W., Sluijs, A., van Helmond, N. A. G. M., Vellekoop, J., Wagner, T.,
1009 Wrobel, N. E.: Cretaceous sea-surface temperature evolution: Constraints from TEX₈₆ and
1010 planktonic foraminiferal oxygen isotopes, *Earth-Science Reviews*, 172, 224-247, 2017.
- 1011 Park, E., Hefter, J., Fischer, G., Mollenhauer, G.: TEX₈₆ in sinking particles in three eastern Atlantic
1012 upwelling regimes. *Organic Geochemistry*, 124, 151–163, 2018.
- 1013 Pearson, P. N., van Dongen, B. E., Nicholas, C. J., Pancost, R. D., Schouten, S., Singano, J. M., Wade, B.
1014 S.: Stable warm tropical climate through the Eocene Epoch. *Geology*, 35, 211-214, 2007.
- 1015 Polik, C. A., Elling, F. J., Pearson, A.: Impacts of Paleoecology on the TEX₈₆ Sea Surface Temperature
1016 Proxy in the Pliocene-Pleistocene Mediterranean Sea. *Paleoceanography and Paleoclimatology*, 33,
1017 1472–1489, 2018.
- 1018 Qin, W., Amin, S. A., Martens-Habbena, W., Walker, C. B., Urakawa, H., Devol, A. H., Ingalls, A.E.,
1019 Moffett, J.W., Ambrust, E.V., Stahl, D. A.: Marine ammonia-oxidizing archaeal isolates display
1020 obligate mixotrophy and wide ecotypic variation. *Proceedings of the National Academy of
1021 Sciences of the United States of America*, 111, 12504–12509, 2014.
- 1022 Qin, W., Carlson, L. T., Ambrust, E. V., Stahl, D. A., Devol, A. H., Moffett, J. W., Ingalls, A. E.:
1023 Confounding effects of oxygen and temperature on the TEX 86 signature of marine
1024 Thaumarchaeota. *Proceedings of the National Academy of Sciences*, 112, 10979–10984, 2015.
- 1025 Rasmussen, C.E., Nickisch, H.: Gaussian Processes for Machine Learning (GPML) Toolbox. *J. Mach.
1026 Learn. Res.* 11, 3011–3015, 2010.
- 1027 Rasmussen, C. E. & Williams, C. K. I.: *Gaussian Processes for Machine Learning*, the MIT Press, 2006,
1028 ISBN 026218253X.
- 1029 Sangiorgi, F., van Soelen, E., Spofforth, D., J. A., Pälike, H., Stickley, C., St. John, K.,
1030 Koç, N., Schouten, S., Sinninghe Damsté, J. A., Brinkhuis, H.: Cyclicity in the middle Eocene
1031 central Arctic Ocean sediment record: Orbital forcing and environmental response,
1032 *Paleoceanography*, 23, 10.1029/2007PA001487, 2008.
- 1033 Schouten, E., Hopmans, E.C., Forster, A., Van Breugel, Y., Kuypers, M.M.M., Sinninghe Damsté, J.S.:
1034 Extremely high seasurface temperatures at low latitudes during the middle Cretaceous as revealed
1035 by archaeal membrane lipids. *Geology*, 31, 1069–1072, 2003.
- 1036 Schouten, S., Forster, A., Panoto, F. E., and Sinninghe Damsté, J. S.: Towards calibration of the TEX₈₆
1037 palaeothermometer for 20 tropical sea surface temperatures in ancient greenhouse worlds, *Organic
1038 Geochemistry*, 38, 1537-1546, 2007.

- 1039 Schouten, S., Hopmans, E. C., Sinninghe Damsté, J. S.: The organic geochemistry of glycerol dialkyl
1040 glycerol tetraether lipids: A review, *Organic Geochemistry*, 54, 19-61, 2013.
- 1041 Schouten, S., Hopmans, E. C., Schefuß, E., Sinninghe Damsté, J. S.: Distributional variations in marine
1042 crenarchaeotal membrane lipids: a new tool for reconstructing ancient sea water temperatures?
1043 *Earth and Planetary Science Letters*, 204, 15 265-274, 2002.
- 1044 Seki, O., Bendle, J. A., Harada, N., Kobayashi, M., Sawada, K., Moossen, H., Sakamoto, T.: Assessment
1045 and calibration of TEX₈₆ paleothermometry in the Sea of Okhotsk and sub-polar North Pacific
1046 region: Implications for paleoceanography. *Progress in Oceanography*, 126, 254–266, 2014.
- 1047 Sluijs A, Schouten S, Pagani M, Woltering, M., Brinkhuis, H., Sinninghe Damsté, J.S., Dickens, G.R.,
1048 Huber, M., Reichart, G., Stein, R., Matthiessen, J., Lourens, L.J., Pedentchouk, N., Backman, J.,
1049 Moran, K. and the Expedition 320 Scientists: Subtropical arctic ocean temperatures during the
1050 Palaeocene/Eocene thermal maximum. *Nature* 441, 610–613, 2006.
- 1051 Sluijs, A., Schouten, S., Donders, T. H., Schoon, P. L., Rohl, U., Reichart, G.-J., Sangiorgi, F., Kim, J.-
1052 H., Sinninghe Damsté, J. S., Brinkhuis, H.: Warm and wet conditions in the Arctic region during
1053 Eocene Thermal Maximum 2, *Nature Geosci*, 2, 777-780, 2009.
- 1054 Taylor, K. W. R., Willumsen, P. S., Hollis, C. J., Pancost, R. D.: South Pacific evidence for the long-term
1055 climate impact of the Cretaceous/Paleogene boundary event, *Earth-Science Reviews*, 179, 287-302,
1056 2018.
- 1057 Taylor, K. W., Huber, M., Hollis, C. J., Hernandez-Sanchez, M. T., Pancost, R. D.: Re-evaluating modern
1058 and Palaeogene GDGT distributions: Implications for SST reconstructions, *Global and Planetary
1059 Change*, 108, 158-174, 2013.
- 1060 Thomas, E. K., S. C. Clemens, W. L. Prell, T. D. Herbert, Y. Huang, Z. Liu, J. S. S. Damsté, Y. Sun, and
1061 X. Wen.: Temperature and leaf wax $\delta^2\text{H}$ records demonstrate seasonal and regional controls on
1062 Asian monsoon proxies, *Geology*, 42(12), 1075–1078, doi:10.1130/G36289.1, 2014.
- 1063 Tierney, J. E.: GDGT Thermometry: Lipid Tools for Reconstructing Paleotemperatures. Retrieved from
1064 https://www.geo.arizona.edu/~jesst/resources/TierneyPSP_GDGTs.pdf, 2012
- 1065 Tierney, J. E., and Tingley, M. P.: A Bayesian, spatially-varying calibration model for the TEX₈₆ proxy.
1066 *Geochimica et Cosmochimica Acta*, 127, 83-106, 2014.
- 1067 Tierney, J. E., and Tingley, M. P.: A TEX₈₆ surface sediment database and extended Bayesian calibration,
1068 *Scientific data*, 2, 150029, 2015.
- 1069 Williams, C.K.I., and Rasmussen, C.E.: Gaussian processes for machine learning. MIT Press Cambridge,
1070 MA, 2006.
- 1071 Wuchter, C., Schouten, S., Coolen, M. J. L., and Sinninghe Damsté, J. S.: Temperature-dependent
1072 variation in the distribution 30 of tetraether membrane lipids of marine Crenarchaeota: Implications
1073 for TEX₈₆ paleothermometry, *Paleoceanography and Paleoclimatology*.
1074 doi:10.1029/2004PA001041, 2004.
- 1075 Zhang, Y. G., and Liu, X.: Export Depth of the TEX₈₆ Signal. *Paleoceanography and Paleoclimatology*.
1076 doi.org/10.1029/2018PA003337, 2018.

1077 Zhang, Y. G., Pagani, M., Wang, Z.: Ring Index: A new strategy to evaluate the integrity of TEX₈₆
1078 paleothermometry, *Paleoceanography*, 31, 220-232, 2016.

1079 Zhang, Y. G., Zhang, C. L., Liu, X.-L., Li, L., Hinrichs, K.-U., Noakes, J. E.: Methane Index: a tetraether
1080 archaeal lipid 15 biomarker indicator for detecting the instability of marine gas hydrates, *Earth and*
1081 *Planetary Science Letters*, 307, 525- 534, 2011.

1082 Zhang, Y.G., Pagani, M., Liu, Z.: A 12-million-year temperature history of the tropical Pacific
1083 Ocean: *Science*, 343, 84-86, 2014.

1084 Zhu, J., Poulsen, C. J., Tierney, J.: Simulation of Eocene extreme warmth and high climate sensitivity
1085 through cloud feedbacks, *Science Advances*, 5(9), eaax1874, 2019.

1086

1087

1088

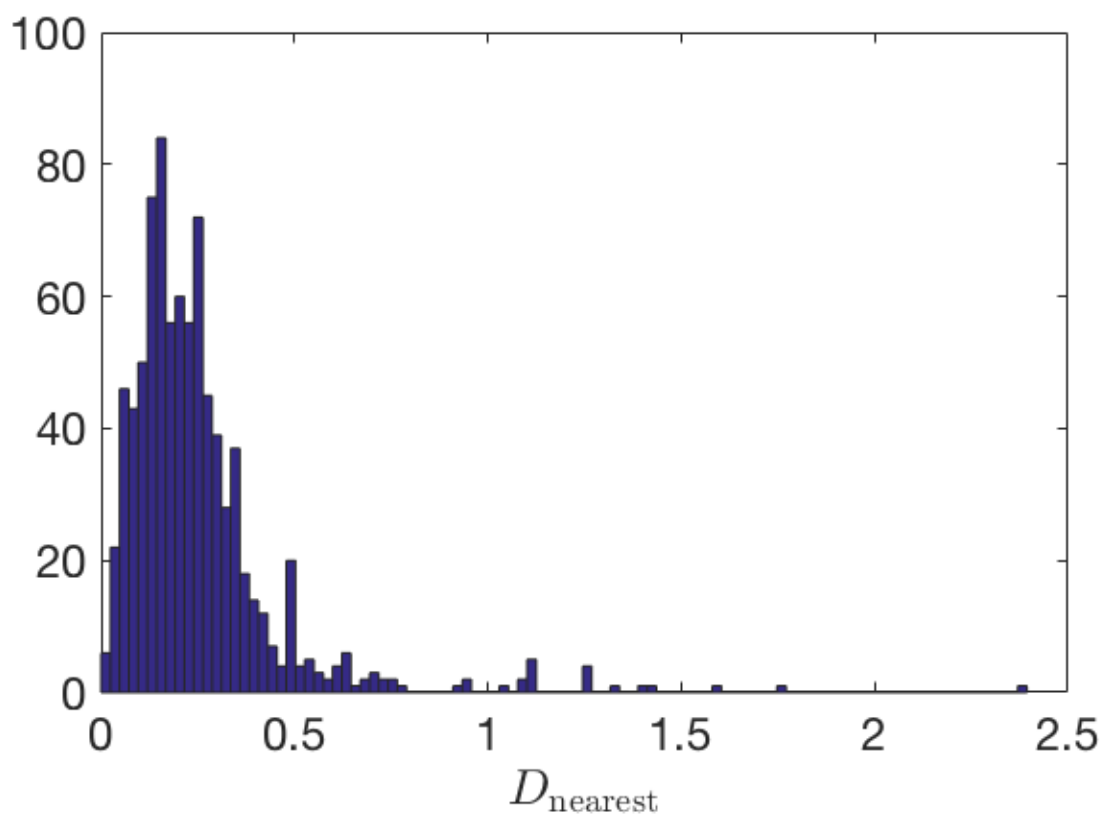
1089

1090

1091

1092 **Figure 1**

1093

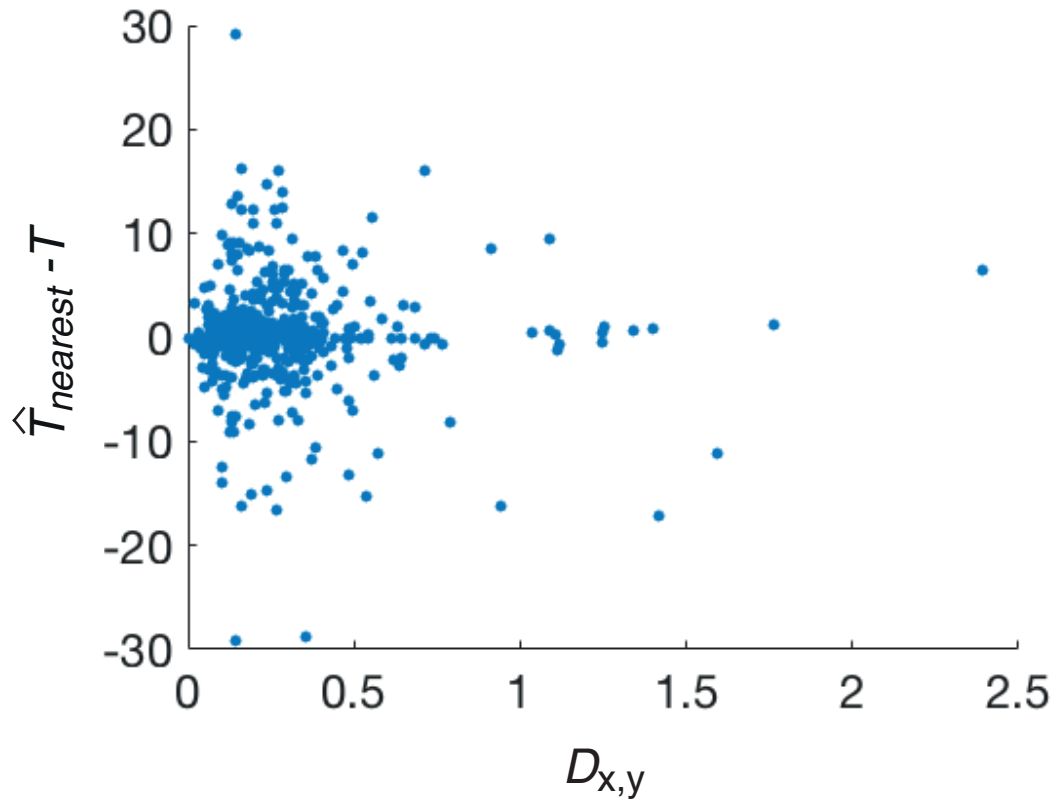


1094

1095

1096

Figure 2



1097

1098

1099

Figure 3

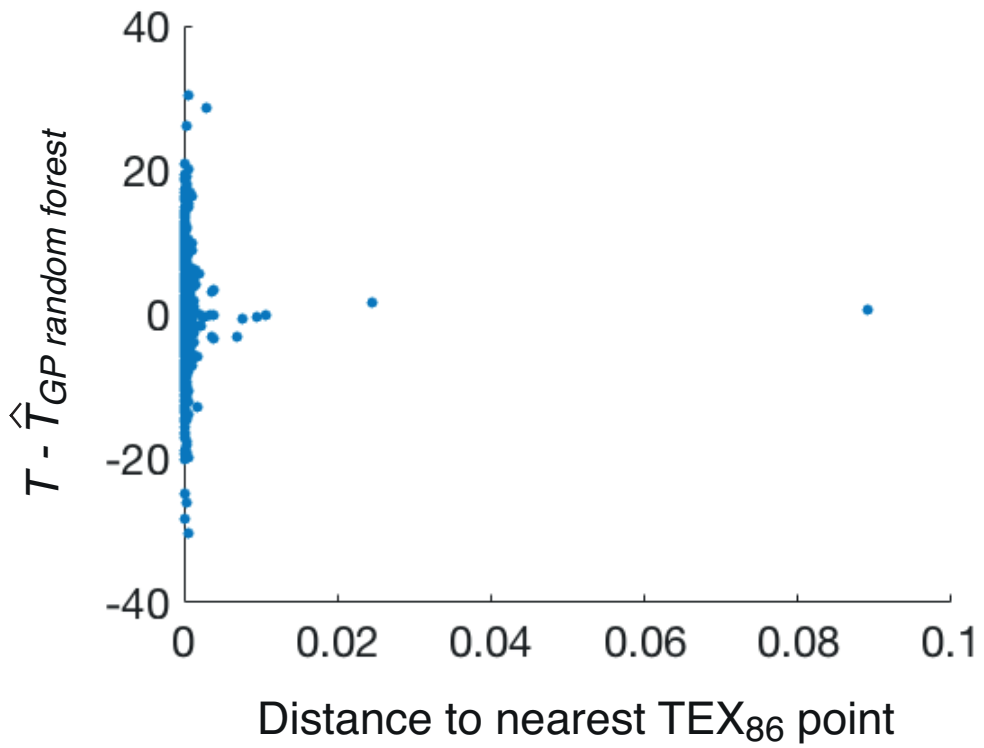
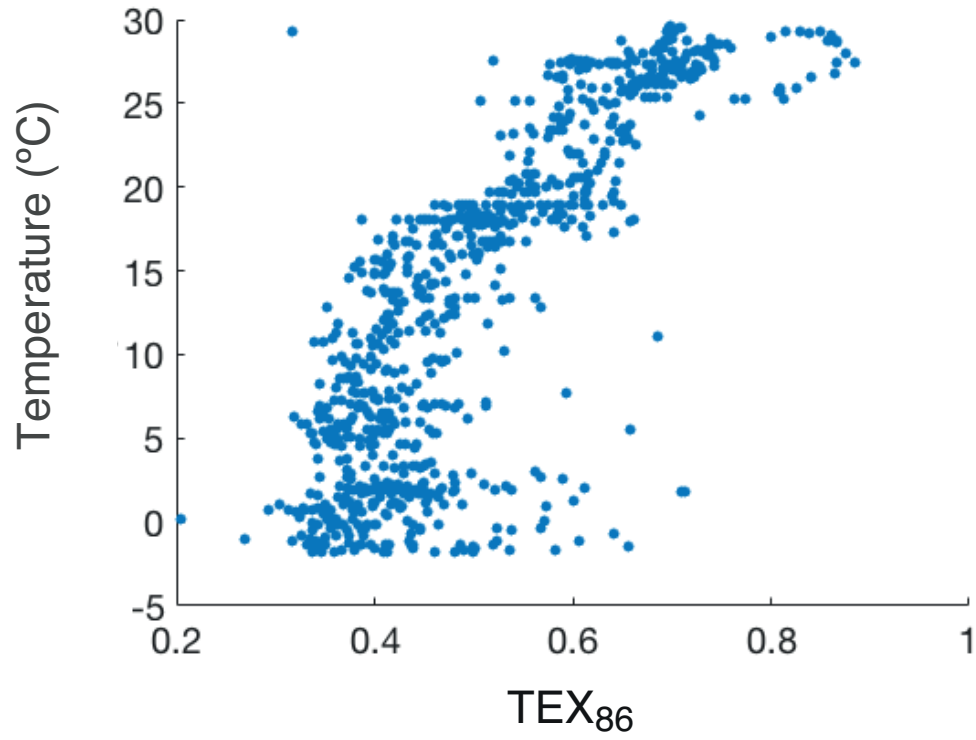
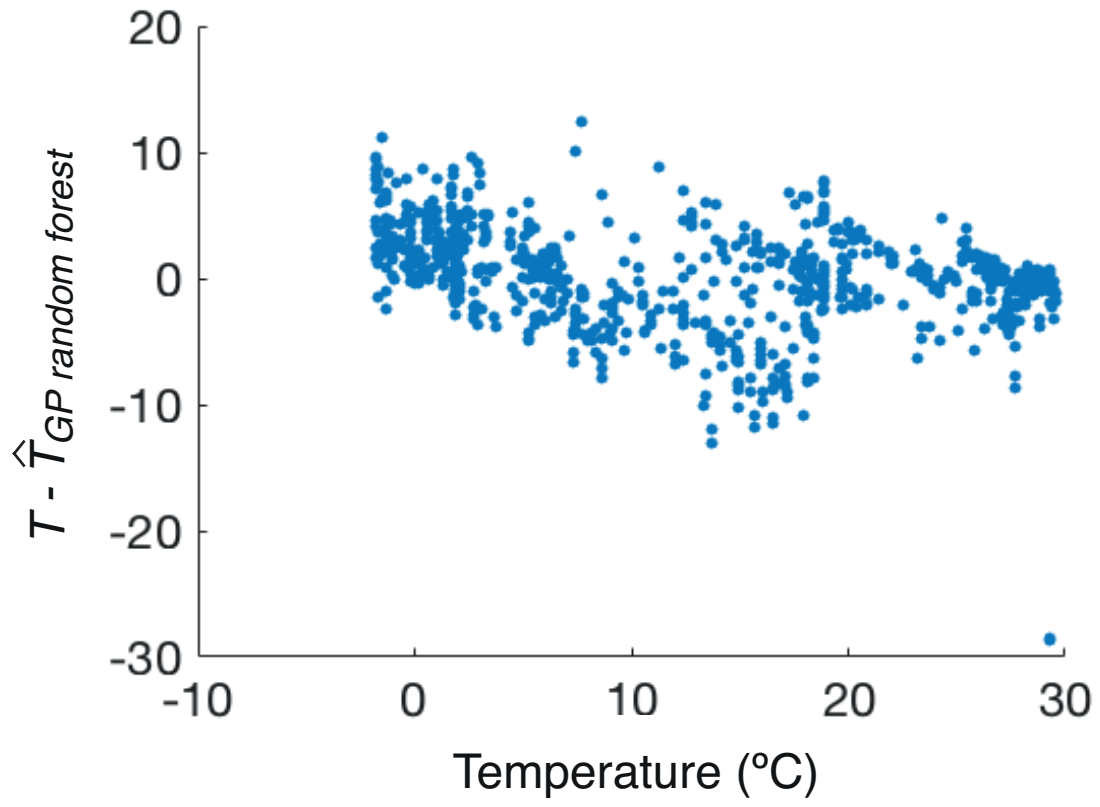


Figure 4



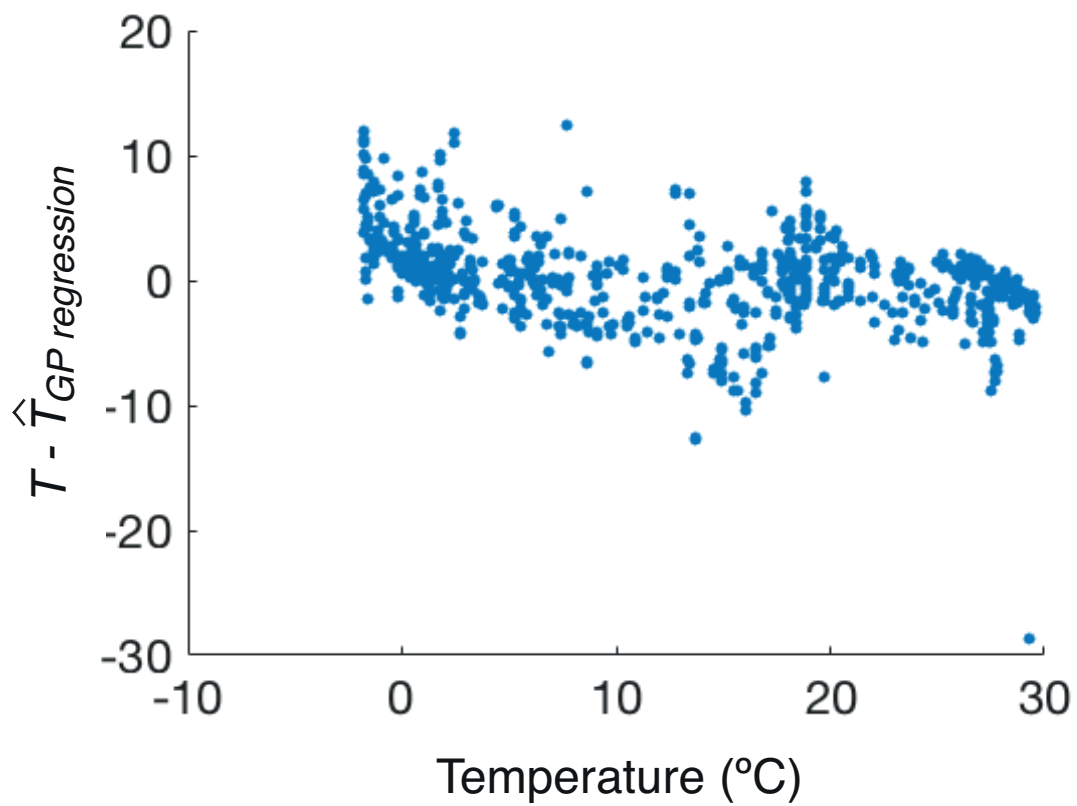
1101

1102

1103

1104

Figure 5

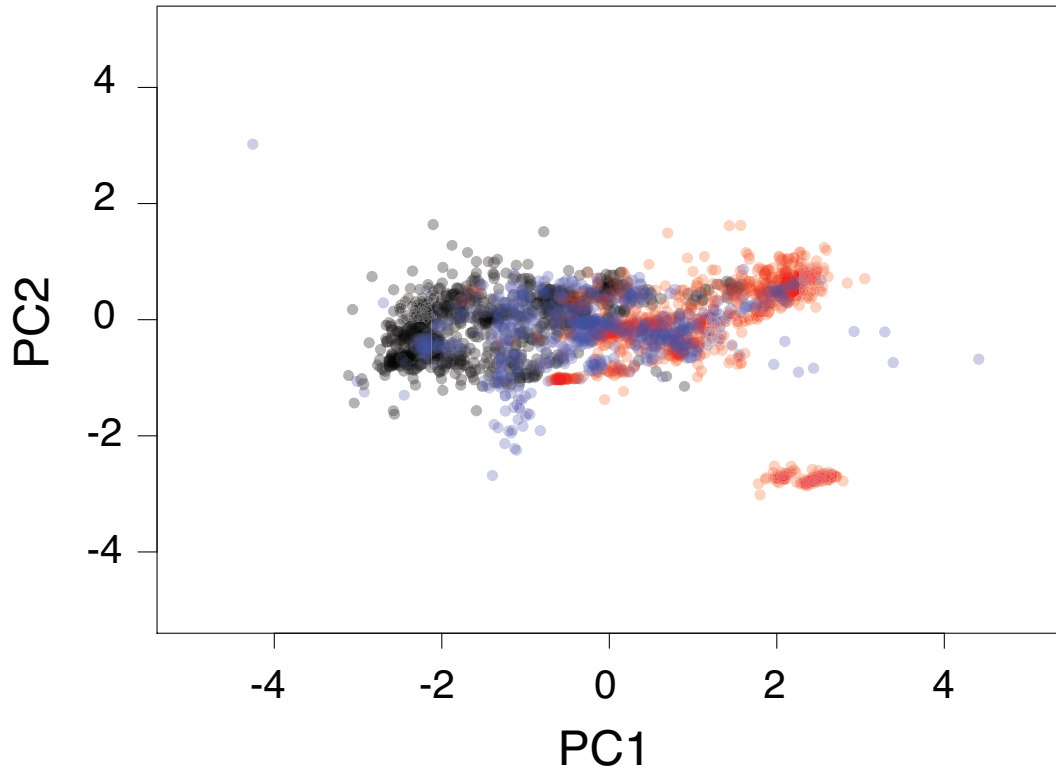


1105

1106

1107

Figure 6

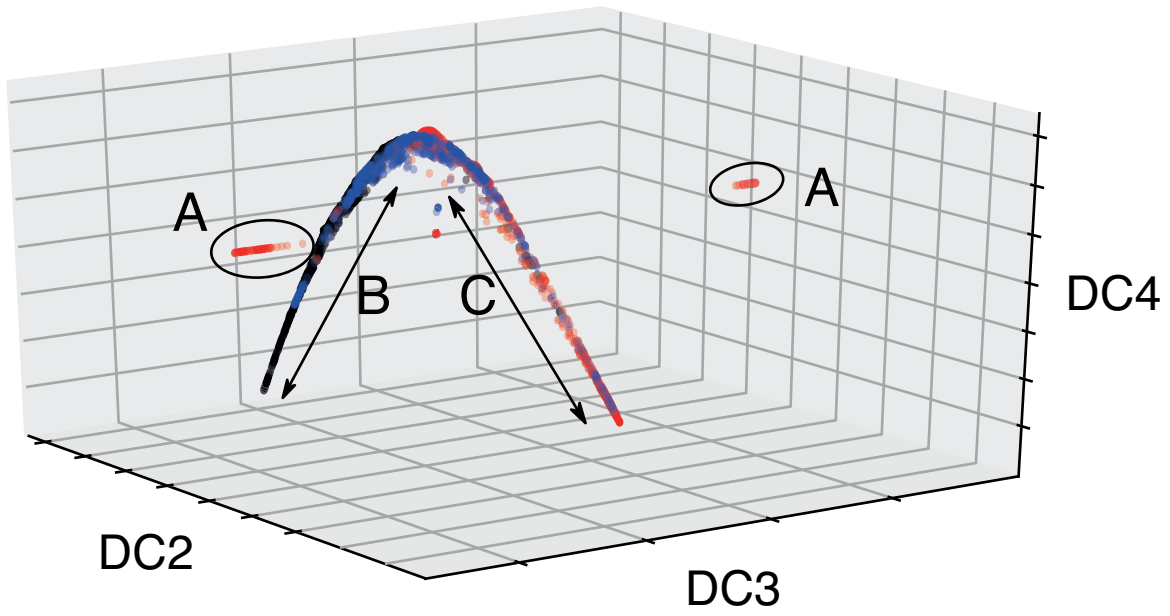


1108

1109

1110

Figure 7



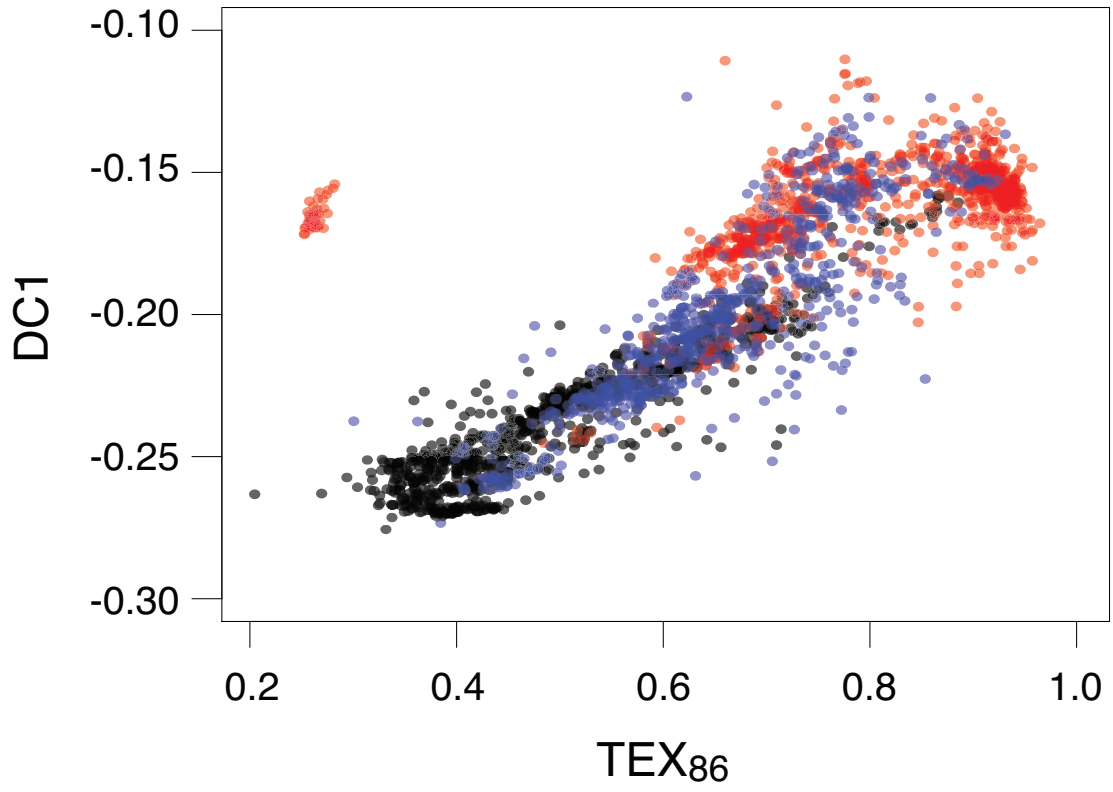
1111

1112

1113

1114

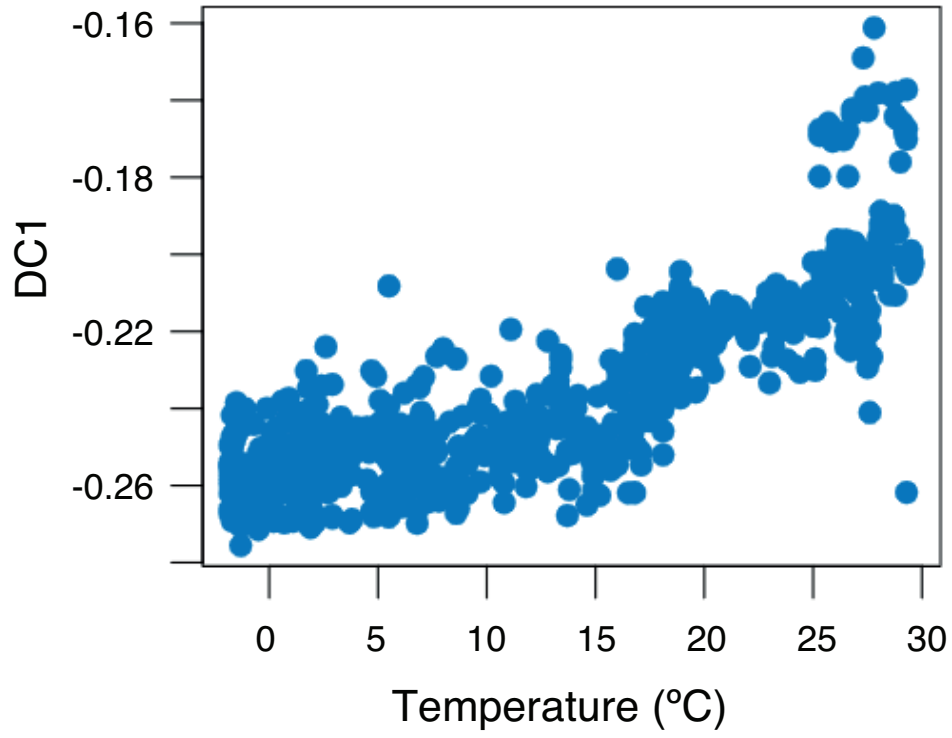
Figure 8



1115

1116

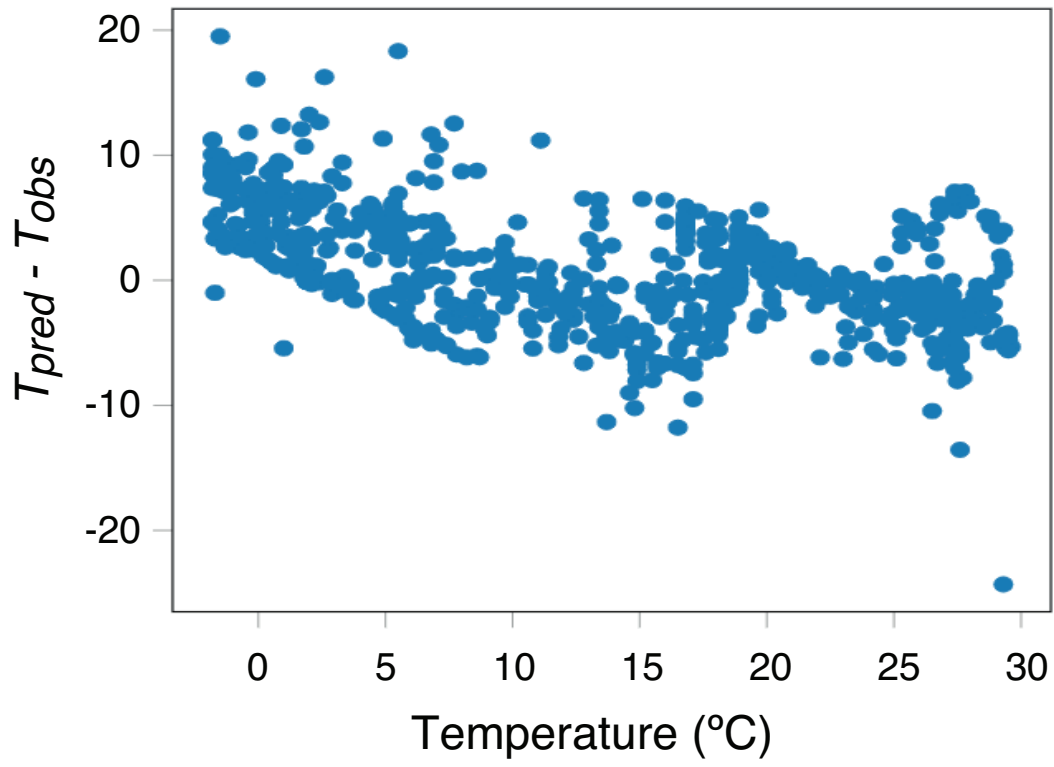
Figure 9



1117

1118

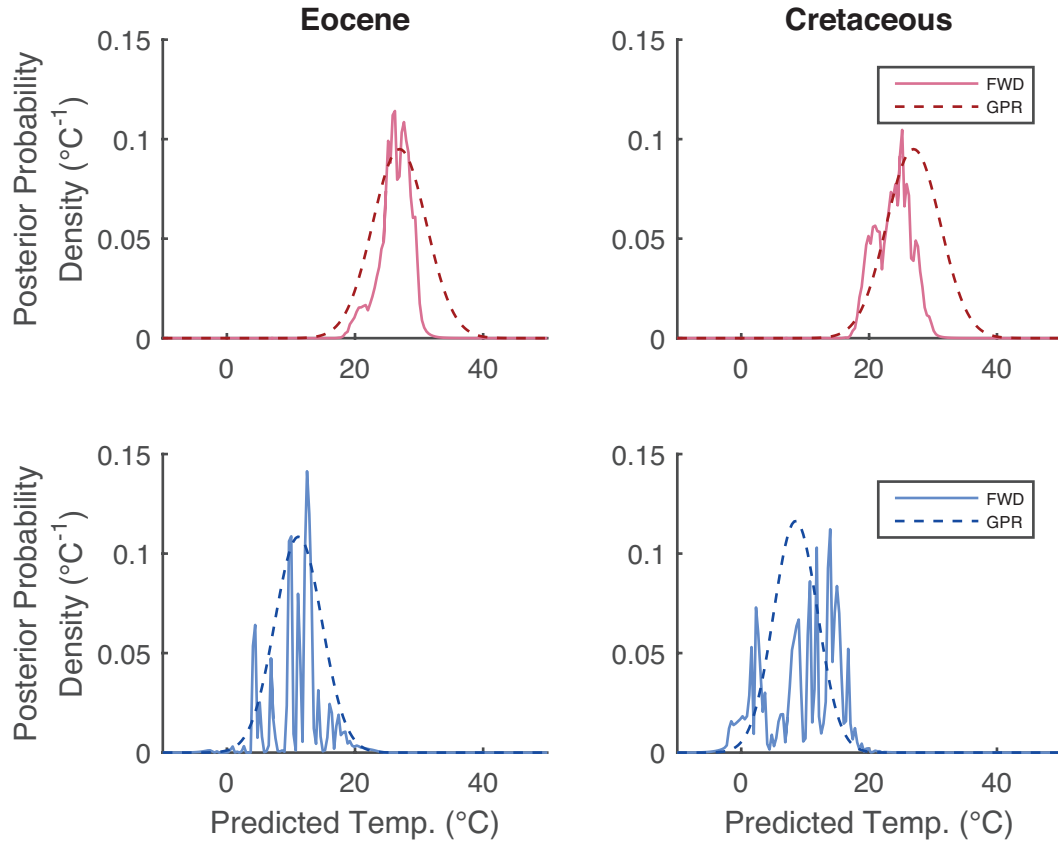
Figure 10



1119

1120

Figure 11

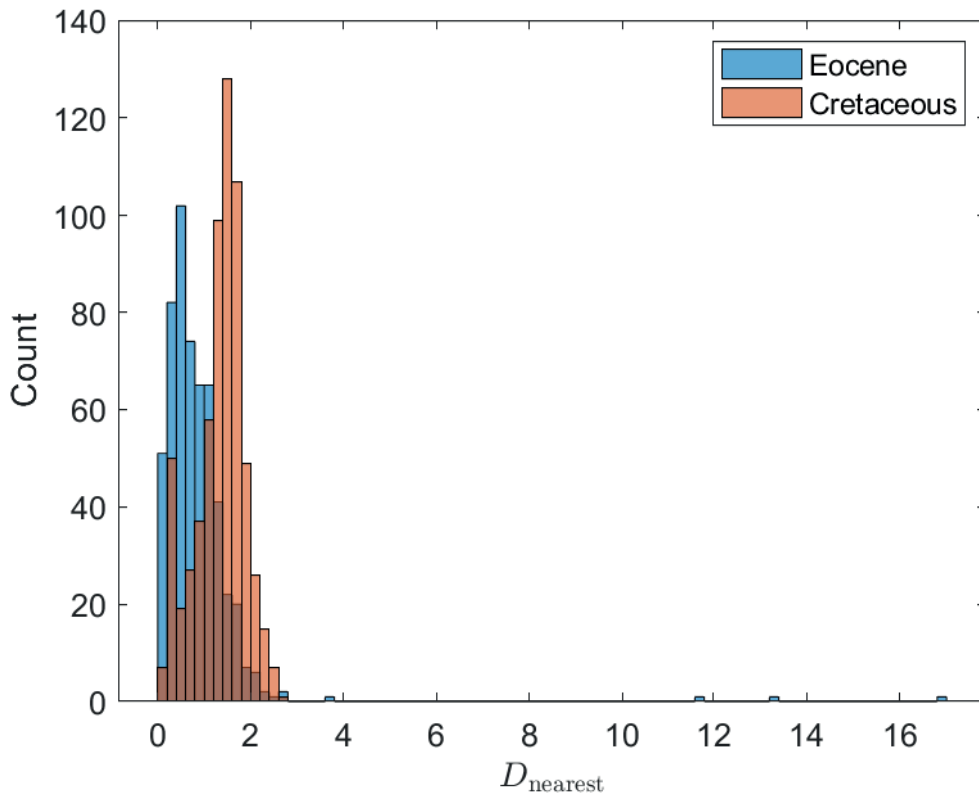


1121

1122

1123

Figure 11



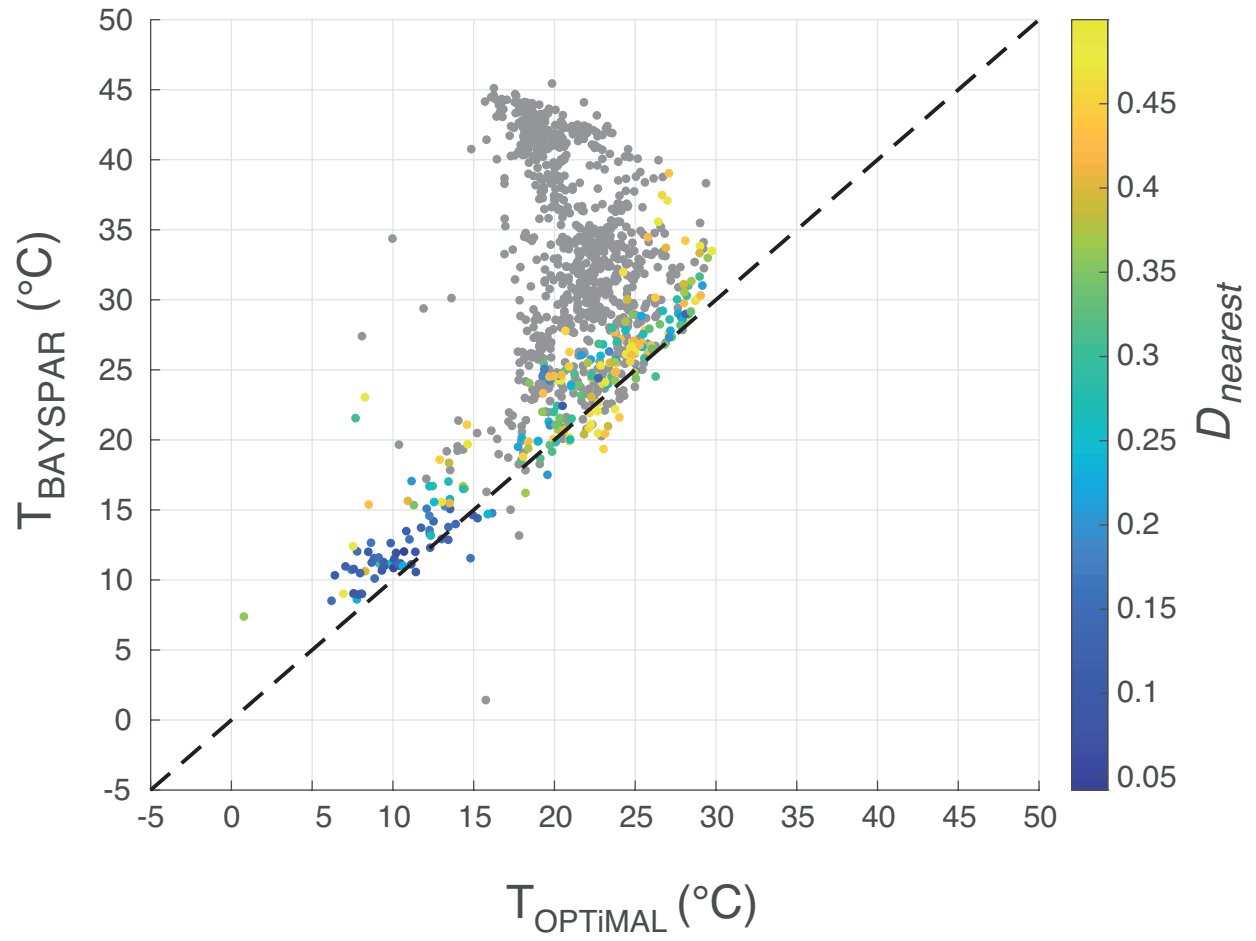
1124

1125

1126

1127

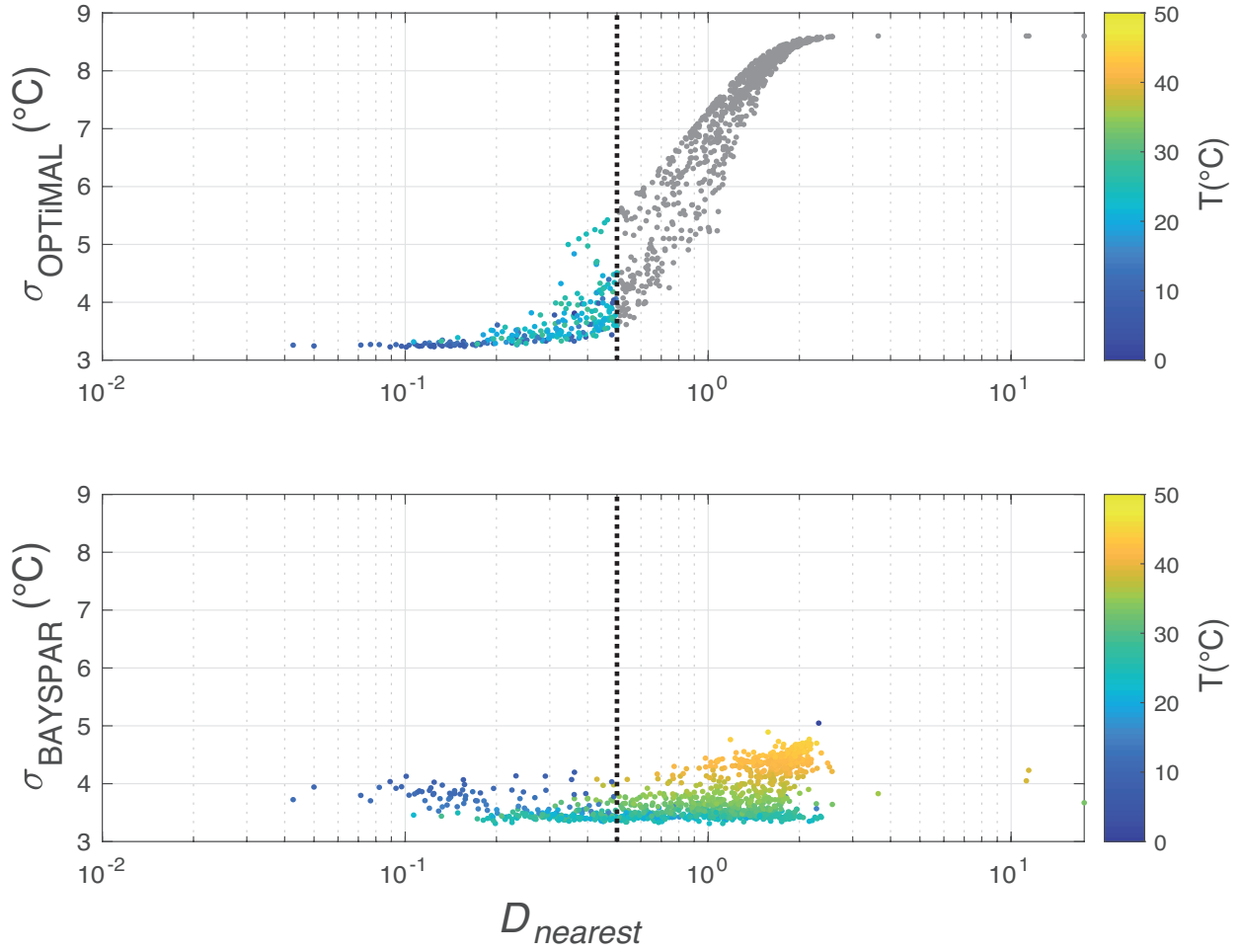
Figure 13



1128

1129

Figure 14



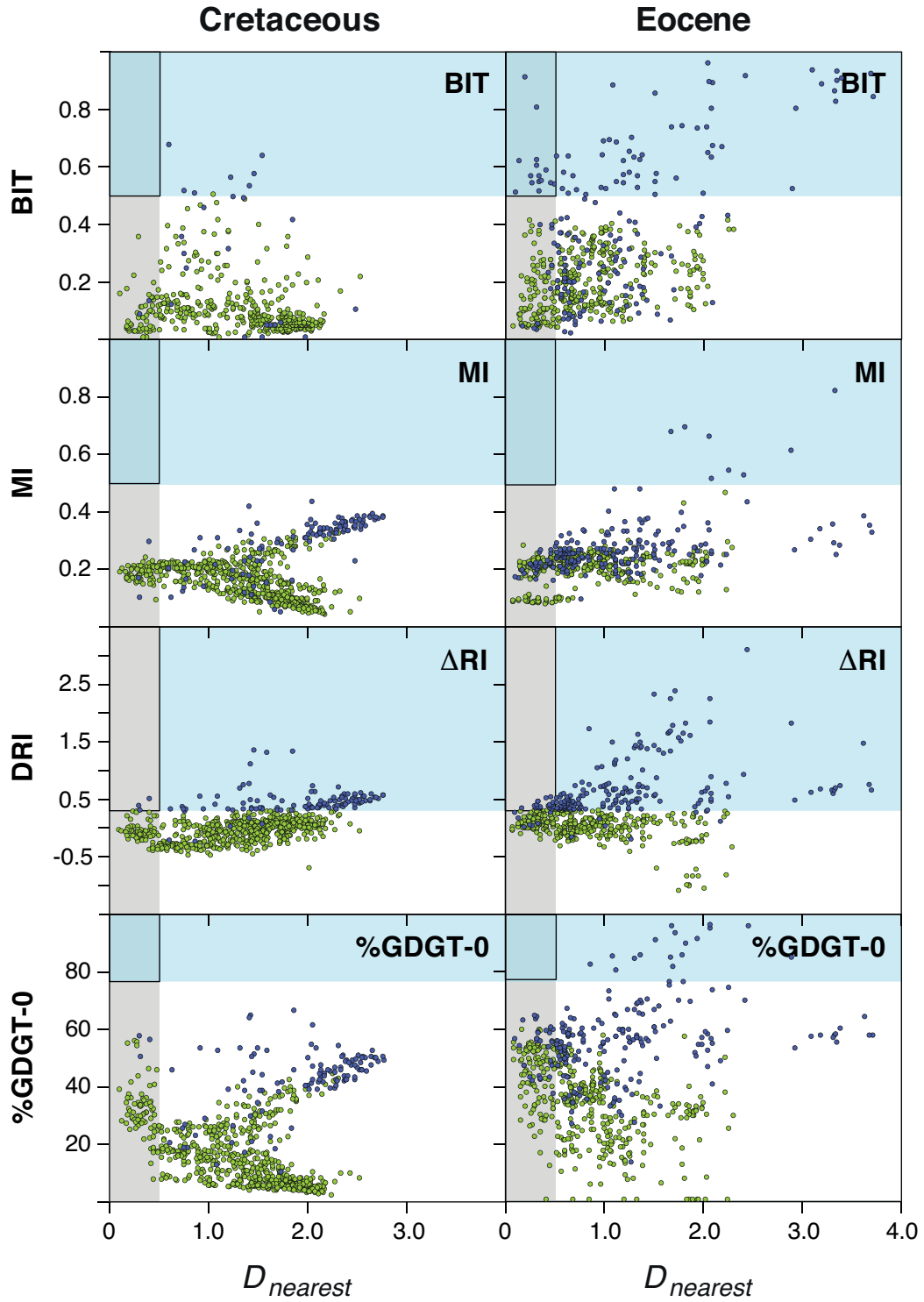
1130

1131

1132

1133

1134 Figure 15

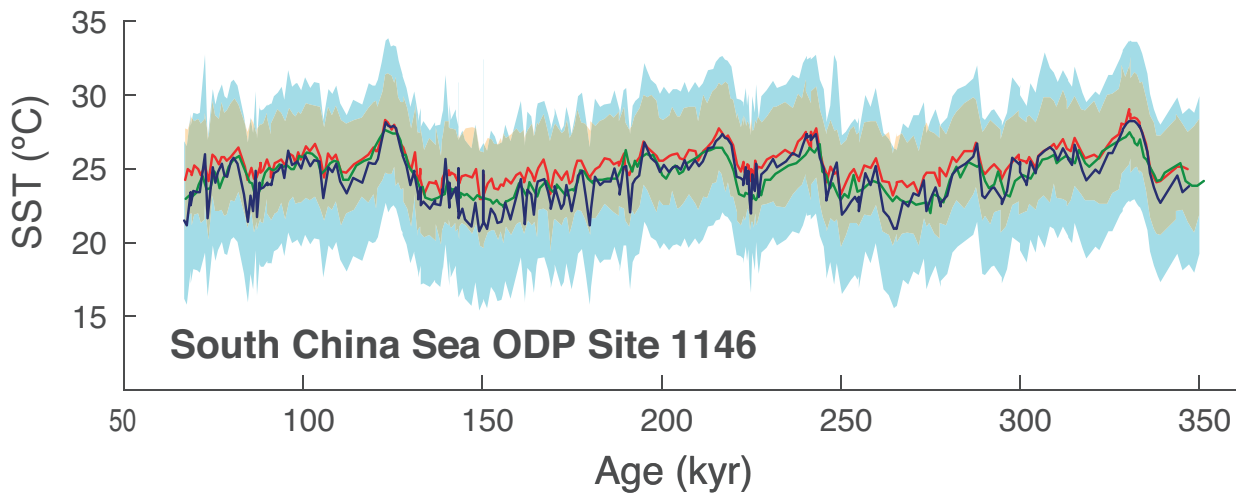
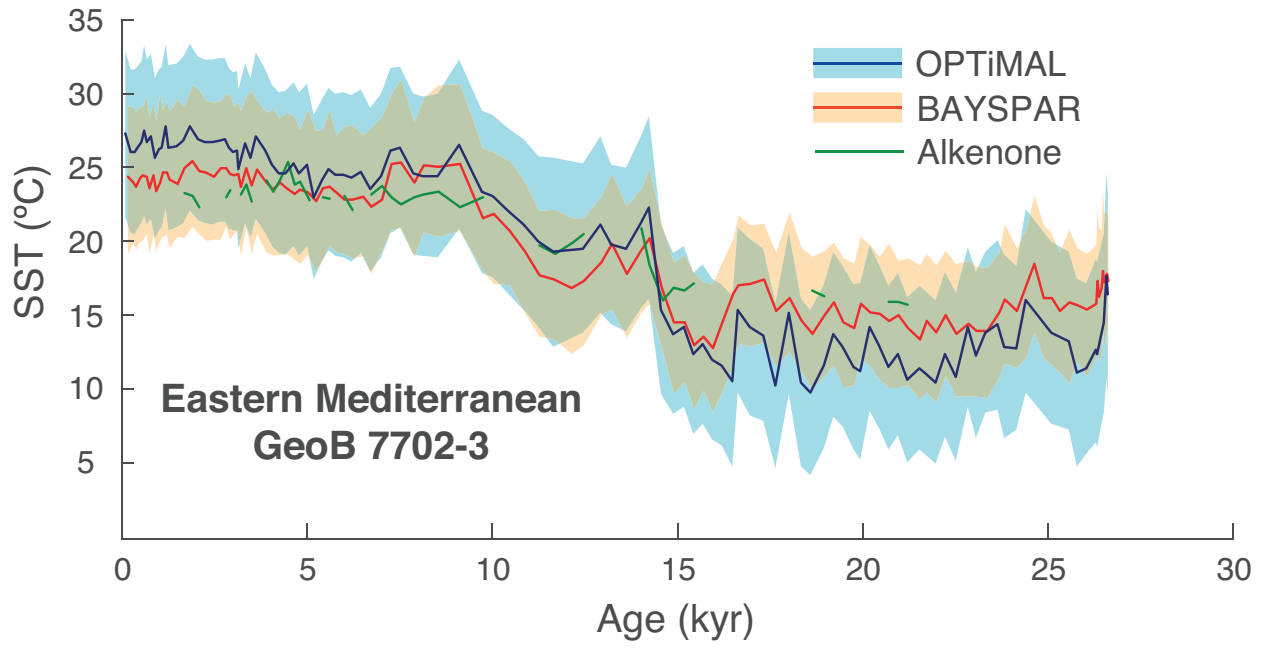


1135

1136

1137 **Figure 16**

1138

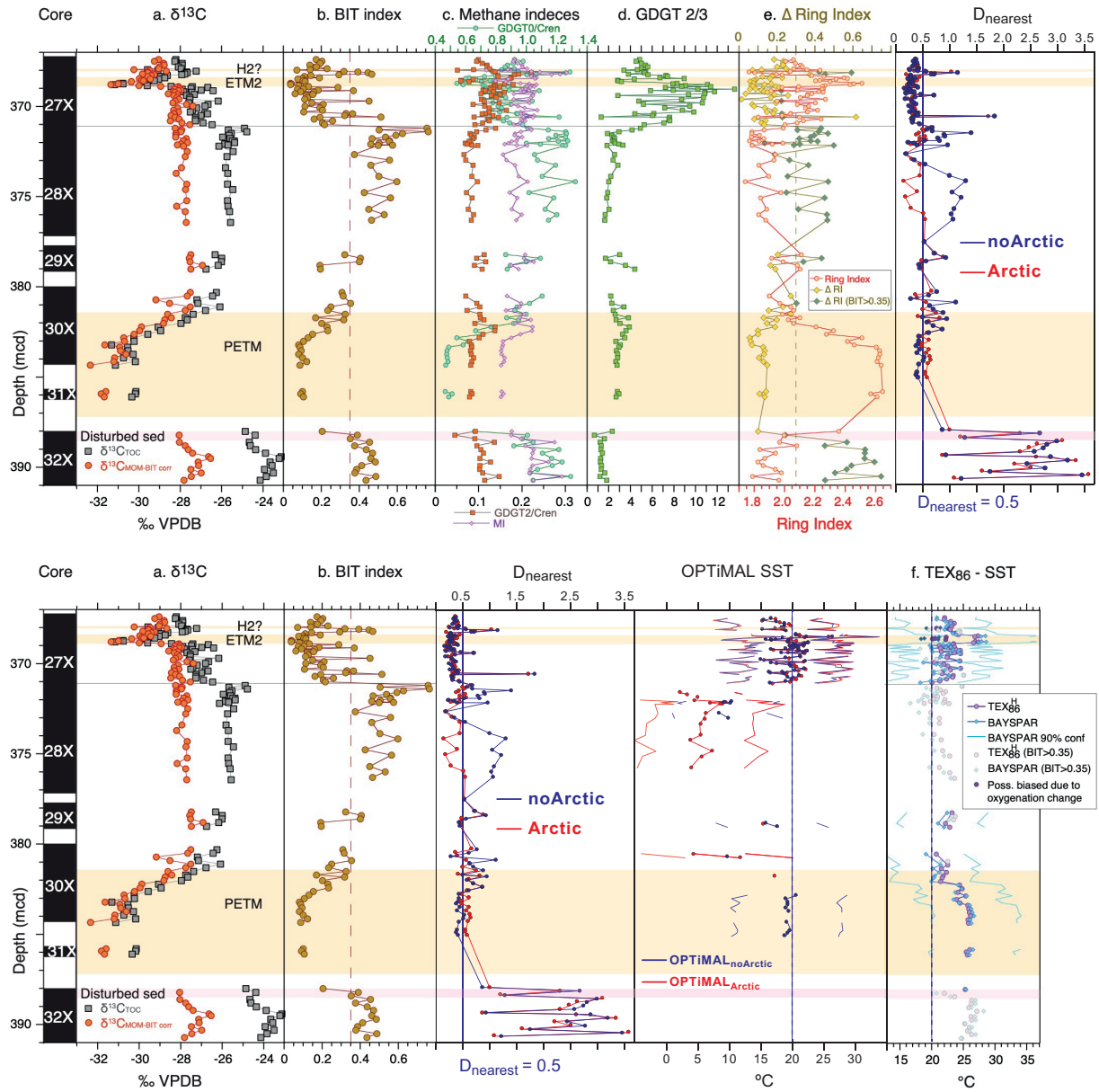


1139

1140

1141 **Figure 17**
 1142

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



1143