1 A Predictive Algorithm For Wetlands In Deep Time Paleoclimate Models

- 2 David J. Wilton¹, Marcus Badger^{2,3,4}, Euripides P. Kantzas¹, Richard D. Pancost³, Paul J.
- 3 Valdes⁴, David J. Beerling¹
- 4
- ¹Dept. Animal and Plant Sciences, The University of Sheffield, Sheffield, S10 2TN, UK
- ⁶ ²School of Environment, Earth and Ecosystem Sciences, The Open University, Milton
- 7 Keynes, MK7 6AA
- ³Organic Geochemistry Unit, The Cabot Institute, School of Chemistry, School of Earth
- 9 Sciences, The University of Bristol, Bristol, BS8 1TH, UK
- ⁴Bristol Research Initiative for the Dynamic Global Environment (BRIDGE), The Cabot
- 11 Institute, School of Geographical Sciences, The University of Bristol, BS8 1TH, UK
- 12 *Correspondence to*: David J. Wilton (d.j.wilton@shef.ac.uk)
- 13

Abstract. Methane is a powerful greenhouse gas produced in wetland environments via 14 15 microbial action in anaerobic conditions. If the location and extent of wetlands are unknown, such as for the Earth many millions of years in the past, a model of wetland fraction is 16 required in order to calculate methane emissions and thus help reduce uncertainty in the 17 18 understanding of past warm greenhouse climates. Here we present an algorithm for predicting 19 inundated wetland fraction for use in calculating wetland methane emission fluxes in deep 20 time paleoclimate simulations. The algorithm determines, for each grid cell in a given 21 paleoclimate simulation, the wetland fraction- predicted by a nearest neighbours search of 22 modern day data in a space described by a set of environmental, climate and vegetation variables. To explore this approach, we first test it for a modern day climate with variables 23 obtained from observations and then for an Eocene climate with variables derived from a 24 25 fully coupled global climate model (HadCM3BL-M2.2, Valdes et al., 2017). Two independent dynamic vegetation models were used to provide two sets of equivalent 26 vegetation variables which yielded two different wetland predictions. As a first test the 27 method, using both vegetation models, satisfactorily reproduces modern data wetland fraction 28 at a course grid resolution, similar to those used in paleoclimate simulations. We then applied 29 the method to an early Eocene climate, testing its outputs against the locations of Eocene coal 30 deposits. We predict global mean monthly wetland fraction area for the early Eocene of 8 to 31 10×10^6 km² with corresponding total annual methane flux of 656 to 909 Tg CH₄ year⁻¹, 32 depending on which of two different dynamic global vegetation models are used to model 33 wetland fraction -and methane emission rates. Both values are significantly higher than 34 estimates for the modern-day of 4×10^6 km² and around 190 Tg CH₄ year⁻¹ (Poulter et- al., 35 2017, Melton et. al., 2013). 36

37

38 **1 Introduction**

- 39 Methane (CH₄) is a powerful greenhouse gas. As well as absorbing infrared radiation from
- 40 the Earth's surface it also contributes to additional indirect warming through its
- 41 photochemistry and oxidation to CO₂ in the atmosphere (IPCC 2013). <u>Along with other trace</u>
- 42 gases, methane is therefore an important component of the Earth's climate system, but for
- 43 studies of the past, such as warm greenhouse paleoclimates, we lack suitable geochemical or
- 44 <u>biological proxies for methane concentration.</u> Therefore, Earth system models used to
- 45 reconstruct ancient climate or develop future climate scenarios must either assume
- 46 atmospheric methane concentrations as a boundary condition and/or incorporate dynamic
- 47 methane fluxes from natural sources <u>and sinks</u> (Beerling et al. 2011). The main natural source
- 48 of methane is wetland environments via microbial action in anaerobic conditions (Whiticar,
- 49 1999), but methane fluxes from wetlands are also modulated by climatic factors such as
- 50 temperature (Westermann, 1992). Therefore, in order to model fluxes of methane to the
- atmosphere both the extent and locations of wetlands need to be known. For modern day,
- 52 recent past and near future scenarios, maps of observed wetland extent (Prigent et al. 2007,
- 53Papa et al. 2010, Schroeder et al., 2015, Poulter et al, 2017) can be used or wetland extent can
- be calculated at a sub-grid level from fine resolution topographical data (as in the
- 55 TOPMODEL approach of Beven and Kirkby (1979), Lu and Zhuang (2012), Stocker et al.
- 56 (2014), Lu et al. (2016)), as wetlands only form where the ground is relatively flat.

57 For the study of deep time paleoclimates (many millions of years in the past) there are no

- 58 direct observations of wetland extent, although we may use a proxy such as coal deposit
- 59 <u>locations as we discus in section 3.2.1</u>, and the topography is only known on relatively coarse
- 60 resolutions of around 0.5 $^{\circ}$ at best. Therefore, any model calculation of wetland extent must
- 61 either rely on using approximate knowledge of the topography or not rely on the topography
- at all. Previous studies (Beerling et al., 2011, Valdes et al., 2005), the only current model-
- 63 <u>based approach for deep-time paleoclimates</u>, classified grid cells as either producing or not
- 64 producing methane, based on either: i) a month being within a defined melt season, for grid
- 65 cells where mean monthly temperature drops below 0 $^{\circ}$ C at some point infor at least one
- 66 <u>month of</u> the year; or ii) precipitation being greater than evapotranspiration. They then scaled
- 67 emissions by empirically derived functions of the variance or standard deviation of
- orography, at the best resolution available. The scaling effectively reduces methane emission
- 69 rates in grid cells where elevation varies significantly and are therefore unlikely to have
- substantial wetlands within them, but relies on what may be quite coarse resolution
- topography not able to resolve sub-grid scale variations. <u>The goal of this paper is to explore</u>
- 72 <u>other methodologies for calculating wetland extent in the context of a deep time</u>
- 73 <u>paleoclimates.</u>
- 74 In this work we develop a nearest neighbour-based algorithm to predict the fraction of a
- specified area that is wetland (FW). We base this on <u>a</u> modern day reference data set of FW
- and corresponding environmental variables, empirically associating the FW observations with
- corresponding observed climate data and vegetation data calculated using one of two
- 78 dynamic global vegetation models (DGVMs)<u>, the Sheffield Dynamic Global Vegetation</u>
- 79 Model (Woodard et al., 2009; Beerling and Woodward, 2001) and the Lund-Postdam-Jenna
- 80 model (Wania et al., 2009).- Wetland is defined in the same manner as for our reference data
- 81 (Poulter et al. 2017), discussed in the following section. It includes both permanently and
- 82 seasonally flooded soils but excludes lakes, reservoirs, rivers, areas of rice cultivation saline
- 83 <u>estuaries and salt marshes.</u> We demonstrate its application by predicting FW and CH₄ fluxes

84 for an early Eocene (52 Ma) model climate, an interval of greenhouse warming (Zachos et al.,

- 85 2008) when sedimentary records indicate the existence of large areas of wetlands (Sloan et
- al., 1992, Beerling et al., 2009). For the Eocene, the same climate variables are obtained from
- 87 a fully coupled global climate model and vegetation variables are derived from the same
- DGVMs. We then predict FW for the Eocene by analysis and comparison to the modern-day
 reference data. We note that different reference sets, vegetation models or climate models
- will likely yield different results and these should be explored in future work, but our aim
- 91 here is to demonstrate this approach and its potential rather than to produce a model-model
- 92 intercomparson intercomparison.
- 93 In the Data and Methods section we fFirstly, we describe modern day wetland data -at 0.5°
- spatial resolution and a monthly time step for a mean modern day year, along with climate
- and vegetation data which we later use as a reference data set. We then describe two test data
- sets at lower spatial resolution, equivalent to that used in paleoclimate models, again for a
 single year. The first of these is for the modern day and derived by interpolation of the
- 97 single year. The first of these is for the modern day and derived by interpolation of the 98 reference data and the second is derived from a paleoclimate model of the early Eocene. We
- 98 briefly describe unsuccessful attempts to model FW through analysis of the reference data
- set. The main conclusion of these unsuccessful attempts being to indicate that any
- 101 relationship between FW and various environmental variables must be quite complex. We
- 102 <u>then introduce before moving on to the nNearest nNeighbours method we later</u> found to be
- 103 successful <u>and</u>. We <u>finally in this section also</u> describe the model used to calculate wetland
- 104 methane emissions.
- 105 <u>In the Results and Discussion section</u>We then we first discuss-the model results for the
- 106 modern day test data set where we expect and then Early Eocene climate. For the modern day
- 107 test data set the nearest neighbour method should perform wellyield strong agreement, since
- 108 <u>the test data it is simply a downscaled version of the reference data interpolated to lower</u>
- 109 <u>spatial resolution</u>; these results, therefore, serve to demonstrate whether or not <u>some a</u>
- 110 generalised form of the <u>nearest neighbour</u> method c<u>ouldan</u> be successfully applied to
- 111 prediction of FW for a climate very different to the modern day. We then apply this method
- to prediction of FW for the Eocene, and show that we can tune it by using the locations of
- 113 coal deposits as wetland proxies.
- 114

115 **2 Data and Methods**

116 2.1 Modern day reference data

117 We use a modern-day reference data set of observed FW, the term observed being used to

- 118 <u>distinguish this from our later model results</u>, with corresponding environmental data to
- develop an algorithm for the prediction of FW in the past, i.e. we assume that there exists a
- relationship between FW and the environmental variables compiled in the reference data and
- then apply that relationship to predicting FW in the past. We use the recently developed
- 122 SWAMPS-GLWD (Poulter et al., 2017), which improves on the Surface Water Microwave
- 123 Product Series (SWAMPS) (Schroeder et al., 2015) by adding using the static inventory of
- 124 wetland area from the Global Lakes and Wetlands Database (GLWD) (Lehner and Doll
- 125 2004)-data, correcting the SWAMPS dataset in regions where this satellite derived dataset
- fails to detect water beneath closed canopies. We calculated the average monthly FW at each

- 127 $0.5^{\circ} \times 0.5^{\circ}$ grid cell for the years 2000 to 2012 on a monthly time step to give a modern-day 128 FW (FW_{obs}; annual max shown in Figure 1). Corresponding climate data on the same spatial
- and temporal resolution were obtained from CRU-NCEP v4.0 (Wei et al. 2014) and averaged
- to give monthly values for a mean modern-day year over the same time interval. The climate
- data for this mean year were then used to drive two DGVMs: the Sheffield Dynamic Global
- 132 Vegetation Model (SDGVM) (Woodward et al., 1995; Beerling and Woodward, 2001) and
- the Lund-Postdam-Jenna model (LPJ) (Wania et al., 2009) to produce corresponding
- 134 vegetation data. The combination of these yielded a reference data set of FW, climate
- 135 (temperature and precipitation) and vegetation (leaf area index, net primary productivity,
- transpiration, evapotranspiration, soil water content and surface runoff) variables (either
- 137 SDGVM or LPJ) for a set of $0.5^{\circ} \times 0.5^{\circ}$ spatial and monthly temporal resolution sites for a 138 single modern-day average year. Some variables, such as transpiration and
- 139 evapotranspiration, are available from both climate and vegetation models. In such cases we
- 140 use those from the vegetation model as these will be calculated from a more advanced
- 141 <u>vegetation scheme</u>. To ensure that wetlands in areas dominated by agriculture or where one
- 142 of our vegetation models, SDGVM, predicts bare land, did not bias our FW predictions, such
- 143 grid cells were removed from the reference data. For the latter, this was done simply by
- removing those grid cells that SDGVM predicted to be bare land. For the former, we
- removed those that were 50 % or more, by cover, classed as cultivated and managed or
- 146 mosaic cropland (Global Land Cover 2000 database, 2003).
- 147 Many of the methods that can be used to analyse the reference data and predict FW require 148 that the data are scaled, so that each variable covers a similar range of values. Therefore, we 149 scaled the values of each environmental variable, *X*, using their <u>global</u> mean, μ_x , and <u>global</u> 150 standard deviation, σ_x , i.e. for a given grid cell, *J*, each variable was scaled as:

151
$$X'(J) = \frac{X(J) - \mu_X}{\sigma_X}$$
(1)

152 This scales all variables such that they have <u>global</u> mean of 0 and standard deviation 1.

153 2.2 Test data sets

154 A modern-day test data set was made by interpolating the reference climate data to $2.5^{\circ} \times$ 155 3.75°, the spatial resolution often used for paleoclimate models. -The DGVMs simulations 156 were conducted ondriven by this interpolated data to yield the vegetation outputs. All climate 157 and vegetation variables were scaled in the same way as the reference data, using the global means and standard deviations of the reference data. The palaeoclimatic assessment of our 158 model was performed using an early Eocene test data set made using a single year of output, 159 on a monthly time step, from a three dimensional fully dynamic coupled ocean-atmosphere 160 global climate model HadCM3BL-M2.2 (Valdes et al., 2017), on a 2.5° latitude by 3.75° 161 162 longitude grid and at a monthly time step for a single year. To simulate the early Eocene a Ypresian paleogeography and high CO₂ (4x modern; 1120 ppm; Agnostous et al., 2016) was 163 used. SDGVM and LPJ were both run with these model-simulated climate data to produce 164 the vegetation variables required, as was done for the reference data set, whereas temperature 165 and precipitation were derived directly from the climate model. All variables were again 166 scaled using the means and standard deviations of the reference data. Therefore, for each 167 climate, modern day and early Eocene, we have two test data sets for a mean year on a 168

monthly time step, at 2.5° x 3.75° spatial resolution, both with the same climate data, one
with SDGVM vegetation data and one with LPJ vegetation data. Predictions for each test data
set were made with the corresponding vegetation model's reference data set. <u>The reference</u>
and test data sets are summarized in Table 1.

173

174 2.3 Initial unsuccessful models of wetland fraction

Before discussing the model we employed to predict paleoclimate FW, it is useful to describe 175 briefly other strategies that we attempted but that did not yield robust predictions when 176 evaluated against modern-day data. The first of these was to examine FW vs individual 177 178 environmental variables graphically from the reference data, to ascertain if we could define ranges for those variables that corresponded to predominantly low or high FW; this is similar 179 to the approach of Shindell et al. (2004), who proposed threshold values of standard deviation 180 181 of topography, ground temperature, ground wetness and downward shortwave flux for wetland development. However, this proved unsuccessful, revealing only the rather obvious 182 183 relationship that wetlands do not usually occur when mean monthly temperature is below 0 184 °C. Although we expected to identify relationships for FW with other environmental 185 variables (i.e. ground wetness), none were found. This is due to the combined effects of wetland occurrence being the function of multiple factors and the fact that most grid cells 186 have FW ≈ 0 for all months of the year and the number of grid cells with significantly non-187 zero FW is quite small. Therefore, environmental variables associated with high values of 188 FW also tend to be associated with FW ≈ 0 . Poor correlation of FW with environmental 189 variables is also due to the important control exerted by the topography; regardless of 190 climate, wetlands cannot form in landscapes where excess water flows away rather than 191 remaining in situ. Collectively, these factors caused significant overlap in the range of 192

193 environmental variables associated with both low and high FW.

194 Another approach was a multiple linear regression using the reference data in order to derive an equation for FW in terms of linear functions of multiple environmental variables. 195 However, this yielded equations that predicted a widespread occurrence of very low FW, 196 including those areas where FW_{obs} is very high either seasonally or throughout the year. 197 Similarly, poor predictive models were obtained whether derived for all sites or just those 198 restricted to specific plant functional types. These outcomes likely occur because linear 199 regression optimises a function by minimising the error between predicted and observed 200 values. As most grid cells have FW ≈ 0 (Figure 1) the 'best' regression equation is one that 201 predicts FW very low almost everywhere, since in the majority of cases this is quite accurate. 202 Efforts were made to use other optimisation criteria with customised functions that attempted 203 to put more weight on predicting high FW correctly at the expense of larger errors where FW 204 is low. However, these simply over predicted FW. Therefore, we were unable to find any 205 satisfactory solution based on linear regression. That we do not find a satisfactory regression 206 equation for FW on the reference data suggests that any relationship between FW and the 207 environmental variables must be complex and therefore another approach is required if we 208 are to be able to predict FW. 209

- 210
- 211

212 2.4 FW predicted by a nearest neighbour search

- 213 Given that we were unable to find simple mathematical formula with which to predict FW we
- 214 <u>must consider another approach. Nearest neighbour searches can be used to predict a property</u>
- 215 for a query by comparing data for that query to similar such data from a reference data set.
- 216 We find the entry in the reference data set that is most similar to, i.e. the nearest neighbour of,
- 217 the query and predict the query has the same value in the property of interest as its nearest
- 218 <u>neighbour.</u> The reference data set of FW and environmental variables sites on a 0.5° grid at a
- 219 monthly time step can be viewed as a set of data points yielding FW at many different
- 220 locations in a multi-dimensional space. The eight dimensions of that space are the two
- climate and six vegetation variables; temperature, precipitation, leaf area index, net primary
 productivity, transpiration, evapotranspiration, soil water content and surface runoff. It is
- 222 productivity, transpiration, evaportalispiration, son water content and surface runoin. It is
 223 logical to assume that points close to each other in such a space probably have similar FW.
- Therefore, iIf we have the same environmental variables for a site of unknown FW, we can
- search the reference data set for its nearest neighbour and , i.e. the point in the dataset nearest
- to it. We then predict it would have the same FW as that for the nearest neighbour in the
- 227 reference set, as illustrated schematically below.
- The set of N environmental variables, suitably scaled, X₁, X₂ ... X_N, defines an N dimensional space
- 230 2. The Euclidean distance between two points, I and J, in this space is given by D_{IJ}

231 •
$$D_{IJ} = \sqrt{\sum_{k=1,N} (X_k(I) - X_k(J))^2}$$
 (2)

- 232 3. We calculate D_{IJ} for site *I* of unknown FW and all sites, *J*, in the reference data set, 233 for each of which we know FW(*J*)
- 4. We find J_{min} , the nearest neighbour, that which gives the lowest D_{IJ}
- 235 5. We then predict FW (I) = FW (J_{min})
- 236 6. If site *I* is classed as bare land by the DGVM, thereby having all vegetation variables 237 = 0, we predict FW(*I*) = 0
- 238 This nearest neighbour (NN) method can, if necessary, be extended to a KNN method,
- whereby rather than predicting FW based solely on the single nearest neighbour we insteadconsider some function of the K nearest neighbours.
- 241

242 2.5 Calculating wetland methane emissions

The aim of this study was to derive an algorithm for predicting wetland fraction that can then be used to calculate methane emissions. For the latter, we use the empirical method described by Cao et al. (1996), where methane production, mp, and methane oxidation, mo, rates for a specific grid cell and month, both in g CH₄ m⁻² month⁻¹, are given by:

$$247 mtext{ } mp = R_h f_t (3)$$

248
$$mo = mp \left(0.6 + 0.3 \frac{GPP}{GPP_{max}}\right)$$
(4)

249 Where R_h is <u>absolute</u> soil respiration and <u>absolute</u> *GPP* is gross primary productivity, both <u>in</u> 250 <u>g C m⁻² month⁻¹ and</u> obtained from the respective vegetation model. *GPP_{max}* is the maximum value of GPP for that grid cell for any month of the year. f_t is a function that scales for <u>air</u> temperature, *TMP*, in °C.

253
$$f_t = \frac{\exp(0.04055 \, TMP)}{3.375}$$
 (5)

This is capped at a maximum value of 1. In principle there would also be a scaling function for water table depth, but this is defined as 1 for inundated wetlands and we are only modelling inundated wetland fraction, as that is how the SWAMPS-GLWD FW dataset is

- 257 defined.
- 258 Methane emission rate, *me*, is then the difference between methane produced and methane 259 oxidised, scaled by the wetland fraction for that grid cell and month

$$260 \quad me = (mp - mo) FW$$

(6)

261

262 **3 Results and Discussion**

263 3.1 Modern day test data set

264 The modern-day test set explained in Sect. 2.2 was used as a first, simple, test of the nearest neighbour algorithm for predicting FW described in Sect. 2.4. Since the modern-day test set 265 is simply the reference climate data downscaled interpolated from 0.5° to the courser 266 HadCM3BL-M2.2 model grid of 2.5° by 3.75° (with vegetation from the DGVMs), we 267 expect the NN algorithm to yield predicted FW reasonably consistent with a similar 268 downscaling of the SWAMPS-GLWD observed FW. If the NN predicted FW does not 269 achieve this, then that would indicate that the NN algorithm has failed to predict FW 270 sufficiently accurately. Therefore this test is primarily designed to indicate that a nearest 271 neighbour algorithm either does or does not have the potential to be applied to paleoclimates. 272

Fig. 2 shows maps of seasonal, June–July–August and December–January–February, average 273 FW from the observed SWAMPS-GLWD data interpolated to 2.5° x 3.75° along with the 274 predicted FW using either SDGVM or LPJ vegetation data test sets. For both vegetation 275 models, the predicted FW maps are similar to the observed-interpolated data. Sparse patches 276 of high FW occur in the tropics, especially the Amazon, throughout the year, and large areas 277 of seasonal summer wetlands occur in Alaska, Canada and Siberia. The monthly variation of 278 FW north and south of 30° N, i.e. essentially comparing boreal and tropical wetlands is 279 shown in Figure 3. We split the global values into these two zones because there are virtually 280 no southern hemisphere boreal wetlands, and any division based purely on latitude is 281 arbitrary. The nearest-neighbour algorithm generates the correct seasonal FW pattern in 282 boreal regions and, as expected, a relatively constant monthly FW in the tropics. However, 283 284 SDGVM consistently underestimates the amount of tropical wetland, whilst LPJ agrees reasonably well with observations; mean monthly values are 2.11, 1.47 and 1.90 x 10^6 km² 285 for the observed, SDGVM and LPJ respectively. This is due to the fact that SDGVM classes 286 some grid cells as bare land, assumed to have FW = 0 in our algorithm, even though some of 287 these have non-zero FW in the SWAMPS-GLWD database. LPJ does not classify these grid 288 cells as bare land but instead treats them as very low amounts of vegetation, therefore 289 290 vielding higher global FW that is more consistent with observations. If we exclude from the observed data those grid cells SDGVM predicts as bare land, then the SDGVM prediction 291

- matches better the observed data and LPJ predictions (Table <u>2</u>4). These results give
- confidence that a nearest neighbour algorithm is able to reproduce acceptable FW based on
- these specific climate and vegetation variables.

295 Figure 4 shows the monthly variation in wetland methane emissions for boreal and tropical 296 areas, calculated using: the observed or predicted FW, both vegetation models' outputs and 297 Eq. 3 to 6. The annual methane emissions totals are summarised in Table 32, along with other 298 recent estimates from model intercomparisons. The annual and monthly zonal methane emissions are broadly similar for a given vegetation model regardless of whether the 299 observed or predicted FW is used. SDGVM gives global emissions in line with the other 300 modelling studies, whereas those from LPJ are somewhat lower. This is mainly due to 301 differences in tropical emissions. SDGVM yields higher tropical emissions than LPJ but 302 slightly lower emissions north of 30°N. The main factors influencing the modelled methane 303 emissions (other than FW) are, according to equations (3) to (5), temperature (which is the 304 same for both vegetation models), soil respiration (R_h) and gross primary productivity (GPP), 305 the latter two differing between the two vegetation models. It appears that differences in R_h 306 lead to the different zonal methane totals. South of 30° N SDGVM and LPJ model annual 307 total R_h of 46,000 Tg C year⁻¹ and 35,000 Tg C year⁻¹ respectively and, using the same 308 observed FW, SDGVM and LPJ model annual methane emissions of 123 Tg CH₄ year⁻¹ and 309 69 Tg CH₄ year⁻¹ respectively. Therefore, in the tropics the differences in the predicted 310 methane emissions seem to be due to differences in calculated R_h . North of 30° N both 311 DGVMs have similar $R_{h,2}$, 20,000 Tg C year⁻¹ and 22,000 Tg C year⁻¹ respectively for 312 SDGVM and LPJ, and similar values of methane emissions, 64 Tg CH₄ year⁻¹ and 65 Tg CH₄ 313 vear⁻¹ respectively. 314

315 We stress that this was simple test for a nearest-neighbour approach, for reasons outlined at 316 the beginning of this section, and the satisfactory results obtained here merely indicate this is

- an approach that has potential to be useful in predicting FW for a paleoclimate.
- 318

319 **3.2 Early Eocene climate**

320 In the previous section we have shown that a NN method can reproduce FW for a modern

day climate, justifying its application to the early Eocene climate described in section 2.2.

However, as noted at the end of section 2.4 a NN method can be extended to KNN, whereby

323 we predict FW based on some function of the FW of K nearest neighbours (noting that in 3.1, NN is simply 1NN is KNN with K = 1). A 1NN showith we that works well to any dist modern

NN is simply 1NN, i.e. KNN with K=1). A 1NN algorithm that works well to predict modern day FW may not work as well for a paleo climate of many millions of years in the past. The

reference data set we use, section 2.1, is very similar to the modern day test set, the latter's

climate data is simply obtained by interpolating the former to a courser spatial grid.

- 328 Therefore, we expected and observed high correlation between modern day FW predicted
- from the nearest neighbour in the reference data and the actual FW. The early Eocene test
- data has significant differences to the reference data since the climate of the early Eocene is
- obviously not the same as the modern day. Therefore, it will be harder for a nearest neighbour
- based method, searching a space described by climate and vegetation data, to find a nearest
- neighbour in the modern day reference data with the correct early Eocene FW, whatever that
- may be. It may be that for a high FW early Eocene grid cell the nearest neighbour happens to
- have quite low FW and vice versa. Figure.1 shows that FW can change from very high to

- almost zero over relatively small distances, for example in the Amazon basin, and that
- therefore sites with similar climate and vegetation can have very different FW. The greater
- the degree of difference between the early Eocene and the modern day reference data sets, the
- more likely it is that the first nearest neighbour does not have the correct FW.
- FW calculated for the Early Eocene using the exact same 1NN method as used for the 340 modern day test set yields values of global monthly mean wetland area of 4.07 x 10⁶ km² 341 using SDGVM. This is around 33% higher than that for the modern day, $3.00 \times 10^6 \text{ km}^2$ from 342 Table 24. However, this includes a contribution of $1.53 \times 10^6 \text{ km}^2$ from areas south of 30° S , 343 which have an almost negligible contribution for the modern day, so the tropics and northern 344 Boreal regions actually have lower FW for the Early Eocene. Given that the Early Eocene 345 was significantly warmer and wetter than the modern day (Carmichael et- al. 2017), we 346 expect greater wetland area than the modern day. Beerling et al. (2011) reported global 347 wetland area for an Early Eocene climate using SDGVM; employing their method to our 348 Early Eocene climate, so as to eliminate differences arising from the specific HadCM3 model 349 climate and spatial resolution, yields global monthly mean FW area of $16.29 \times 10^6 \text{ km}^2$, four 350 times higher than the value we would calculate from a 1NN method. Therefore, based on 351 comparison with both the modern day and a previous Eocene study, it appears that a 1NN 352
- 353 method may be unsuitable for a paleoclimate that is very different to our modern day
- reference climate, and we consider KNN with higher values of K.
- 355

356 3.2.1 maxKNN Maximum of K nearest neighbours FW prediction

357 If indeed the 1NN results are too low then that implies that for some hypothetical high FW 358 sites from the Early Eocene, the first nearest neighbours in the reference data have very low 359 FW. Therefore, if we consider higher values of K we may improve our estimate by predicting FW to be the maximum FW of K nearest neighbours (maxKNN) in the reference data. 360 However, applying this approach will yield increasingly higher FW as K increases, 361 requiring a data-constrained optimisation of K. Clearly there are no observations of Eocene 362 wetland distributions with which to properly train any predictive algorithm, but we may 363 utilise a suitable proxy for wetlands to try and obtain such a constraint. Here we use the 364 distribution of coal deposits in the Eocene, (Boucot et al., 2013) shown in Figure 5 as such 365 constraints. There are some limitations to this approach. Coal is formed in wetlands, but can 366 also form in other settings such as lakes; and of course, these datasets do not document where 367 wetlands were present but the sedimentary record is missing or has not been published. In the 368 tropics, coal may not have formed in wetland environments due to a very high rate of carbon 369 cycling and in northern latitudes subsequent glaciations could have eroded coal deposits 370 away. Moreover, data will be sparse or non-existent for remote or inaccessible modern day 371 regions, such as under the Antarctic ice sheet. We also note that precise age and location, 372 especially when comparing to low resolution climate simulations, could cause disagreement 373 for grid-by-grid comparisons. A final and critical complication is that FW is a number 374 between 0 and 1, corresponding to the fraction of a site that is wetland, whereas the coal data 375 is a binary measure: either a grid cell has or does not have a coal deposit within it. For all of 376 these reasons, data-model comparisons must be done cautiously; nonetheless, these data are 377 useful for identifying the most effective K value for reconstructing likely wetlands. 378

- We defined two functions to assess how well a model FW matched the locations of Eocene
- coal deposits. Firstly, fl is defined as the mean distance, in km, of a coal deposit location to a
- 381 grid cell with model FW predicted to be > 0.2. The choice of 0.2 representing significant FW
- is arbitrary but the analysis was repeated with other values and the same conclusions were
- found. Secondly, f^2 is defined as the mean FW of the grid cell closest to each coal deposit
- 384 location, providing that site is within 2 grid points of that coal deposit location, to allow some 385 leeway with regard to different projected locations of land masses in the early Eocene. Again
- the choice of a 2-pixel limit is arbitrary but the analysis was repeated with other limits and
- 387 the same conclusions found.
- Figure 6 shows the values of f1 and f2 for maxKNN predictions of FW with increasing K for both the SDGVM and LPJ Early Eocene data sets, compared to a data set of coal deposit
- locations. As explained, since FW increases with K then by extension, so does the likelihood
- of a site with a coal deposit in or close to it coinciding with a site of significant FW.
- Therefore, we do not seek to find the value of K that will give the lowest value of f1 and
- highest value of f^2 as that would simply be K equal to the size of the entire reference data set.
- Instead, we try to find the lowest value of K that gives a "good" prediction for both f1 and f2.
- Although "good" is a subjective measure, we define it based on where increases in K result in marginal improvements in f1 and f2. For both vegetation models as K increases from 1 to 3 f1
- decreases significantly and f2 increases significantly. For K > 3 the decrease in f1 levels out
- and the increase in f^2 also declines. Therefore, we conclude that based on comparison of
- 399 predicted FW and locations of coal deposits, K=3 is a reasonable choice to make predictions 400 for our early Eocene climate via a maxKNN algorithm.
- 401

402 **3.2.2 FW predicted by max3NN**

Figure 7 shows annual maximum FW (i.e. for each pixel the highest of the 12 monthly 403 values) calculated by a max3NN model using SDGVM or LPJ vegetation data, as described 404 above, with the locations of early Eocene coal deposits also shown. The annual maximum 405 FW is shown here as FW might only need to be high at some point during the year to give 406 407 rise to coal deposits. -The areas of predicted high FW are much larger than for the modern day (Fig. 1); moreover, at this spatial resolution there are often abrupt changes from low-408 medium (yellow) to much higher (red) values leading to some isolated patches of high FW. 409 The approach makes it difficult to interrogate specific factors that drive the increase in 410 Eocene FW compared to today but given the wetter climate of the Early Eocene higher FW 411 than the modern day is to be expected. The patchiness is partly a consequence of using annual 412 maximum FW but also reflects the challenge of predicting a characteristic of a 413 paleoenvironment based on modern day reference data. Considering zonal total FW and 414

- seasonal average FW maps, i.e. averaging out some of the small scale spatial and temporal
- 416 variability, is likely a better approach for understanding ancient methane cycling and these
- 417 are discussed later.
- 418 The maps of predicted FW are quite different for the two vegetation models, but the greatest
- 419 differences are in areas with very little or no coal deposits, e.g. the tropics, north eastern
- 420 North America and Antarctica, making it difficult to critically evaluate them against the data.
- 421 However, the monthly variations given by the two vegetation models in total FW (Figure 8)
- 422 and methane emissions (Figure 9), for the three latitudinal zones are reasonably similar with

423 respect to seasonal variations, in that both have their highest values in the late spring and summer months for zones north of 30° N and south of 30° S and no clear seasonal variation 424 in the tropics. In the tropical zone, predictions of monthly FW area are similar in magnitude 425 for the two vegetation models, with SDGVM usually predicting higher FW than LPJ. 426 However, in the zone north of 30° N LPJ predicts much higher FW than SDGVM throughout 427 June to October with a peak in September, whereas SDGVM peaks in May. A similar but less 428 striking pattern occurs for the zone south of 30°S where again LPJ predicts higher summer 429 FW area than SDGVM. These differences between the two vegetation models are also 430 evident in maps of seasonal average predicted FW (Figure 10). In June to August, SDGVM 431 predicts very little wetland area in the northern hemisphere, whereas LPJ predicts moderate to 432 high FW areas over much of the land north of around 50° N. In December to February both 433 models predict almost zero FW north of around 50° N. In the tropics and the southern 434 hemisphere, the two models predict similar amounts of wetland area, but with SDGVM 435 predicting slightly higher FW overall between 30° S to 30° N and LPJ predicting slightly 436

437 higher FW south of 30° N.

This differs from the modern day distribution of wetlands (Figure 1) and likely arises from a 438 variety of method-dependent factors. First, the coarser resolution leads to more patchy 439 distribution, as is evident in the modern day data in Figures 1 and 2 (top row) at $0.5^{\circ} \times 0.5^{\circ}$ 440 and 2.5° x 3.75° spatial resolution. This is particularly true for the tropics where wetlands do 441 occur in small areas. Secondly, the nature of the nearest neighbour algorithm relies on the 442 principle that a grid cell in a paleoclimate with specific values of environmental variables will 443 have the same FW as a grid cell in a modern day reference data set with similar values for 444 those environmental variables; however, other factors influence wetland fraction, such as the 445 topography. Therefore, a nearest neighbour method predicting FW for a paleoclimate from a 446 modern day reference data may well have errors for a given grid cell and month. These errors 447 should reduce when averaged over latitudinal zones or seasonal averages. 448

449 The differences between methane emissions from the two vegetation models likely arise from 450 their respective impacts onf soil water balance, via the magnitude of evapotranspiration 451 (EVT) relative to precipitation (PRC). As the vegetation model, used to calculate EVT, and climate model, used to calculate PRC, s are not dynamically coupled, PRC will be the same in 452 all Eocene simulations, but EVT will vary; thus, vegetation models that yield elevated EVT 453 454 in a given grid cell are more likely to yield negative water balance (PRC-EVT) and low FW. Figure 11 shows the June to August mean PRC-EVT for SDGVM and LPJ, revealing that it 455 is negative in most places north of 30° N for SDGVM but is slightly positive or at least much 456 closer to zero for LPJ. Therefore, SDVGM will generally predict lower FW by identifying 457 modern day nearest neighbours where PRC < EVT and unlikely to be wetland. The lack of 458 extensive coal deposits in the high northern latitudes, especially where the LPJ-based 459 approach predicts wetlands, could indicate that the LPJ approach has over-predicted FW. 460 However, we caution that this could be a data limitation issue and future work is required to 461 interrogate the forecasts of these two methods. Regardless, both models yield broadly similar 462 463 results on global and zonal terms (Table 43) indicating that the KNN algorithm could be a 464 useful complementary approach for interrogating ancient wetland extent and methane emissions. Global monthly mean FW for the Eocene is $8.5 \times 10^6 \text{ km}^2$ and $10.3 \times 10^6 \text{ km}^2$ 465 predicted by SDGVM and LPJ respectively. Both of these values are larger than for the 466 modern day value of $3.0 \times 10^6 \text{ km}^2$, as we would have expected. 467

468 **4.** Conclusions

We have presented a nearest neighbour method by which FW can be calculated at sites on the 469 470 Earth's surface for an Eocene paleoclimate based on a set of environmental variables obtained from climate and vegetation models and comparison of these to a modern day 471 472 reference data set. This has been used as an offline tool using data obtained from climate and 473 vegetation models, rather than by embedding this within existing Earth systems models, as 474 the goal of this work was to explore and improve on methods of predicting FW for deep time paleoclimates. The precise formulation of the nearest neighbour approach was determined 475 through comparison to locations of Eocene coal deposits and indicated that a max3NN 476 477 method was best suited in this case. That should not be taken to imply that a max3NN would be the best in general; for another paleoclimate a similar analysis to that performed here 478 would be required to determine the optimum implementation of KNN. It would therefore be 479 480 of interest in future work to apply this methodology to other paleoclimates to see if similar results are obtained, perhaps using different environmental variables to those we have used to 481 find nearest neighbours and perhaps other proxies for paleo-FW, should they become 482 available. The predicted distributions of FW are much higher than those of today, as we 483 would expect. We have assessed this using two different global vegetation models, and whilst 484 these do yield some geographical differences in FW arising from different evapotranspiration 485 estimates, they are broadly similar when considering zonal means. For both vegetation 486 models, global monthly mean modelled FW area is less than, around half to two thirds, that 487

- 488 of Beerling et al., 2011, as are the values of the wetland methane emissions. However, our
- new method does not rely on the standard deviation of orography, a variable which is onlyknown to a relatively coarse resolution for deep paleoclimates.
- 491

492 Code and Data

493 This study presents a methodology using existing data and climate and vegetation models.

494 Information relating to these is already included in this article. Code implementing the

495 maxKNN prediction of FW is included as supplement.

496 Author Contribution

497 DJW and DJB planned the work with advice from all co-authors. DJW carried out most of
498 the experimental work with MB providing the HadCM3BL-M2.2 and EPK the LPJ model
499 data. DJW prepared the manuscript with contributions from all co-authors.

500 **Competing Interests**

501 The authors declare that they have no conflict of interest.

502 Acknowledgements

503 Funding was provided by the Natural Environmental Research Council (NERC) grant

504 NE/J00748X/1. The authors would like to thank Chris Scotese for access to and advice on

Eocene coal deposit data. We also thank two anonymous referees for their comments and
 advice on improving this manuscript.

508 **References**

- 509 Anagnostou, E., John, E. H., Edgar, K. M., Foster, G. L., Ridgwell, A., Inglis, G. N., Pancost,
- 510 R. D., Lunt, D. J. and Pearson, P. N.: Changing atmospheric CO2 concentration was the
- 511 primary driver of early Cenozoic climate, Nature, 533(7603), 380–384,
- 512 doi:10.1038/nature17423, 2016.
- 513 Beerling, D. J., and Woodward, F. I.: Vegetation and the Terrestrial Carbon Cycle:
- 514 Modelling the First 400 Million Years. Cambridge University Press, Cambridge, 2001
- 515
- Beerling, D., Berner, R. A., Mackenzie, F. T., Harfoot, M. B., and Pyle, J. A.: Methane and
 the CH₄-related greenhouse effect over the past 400 million years, Am. J. Sci., 309, 97-113,
- 518 DOI 10.2475/02.2009.01, 2009.
- 519 Beerling, D. J., Fox, A., Stevenson, D. S., and Valdes, P. J.: Enhanced chemistry-climate
- 520 feedbacks in past greenhouse worlds, Proc. Natl. Acad. Sci., 108, 9770-9775,
- 521 doi:10.1073/pnas.1102409108, 2011.
- 522
- Beven, K.J. and Kirkby, M.J.: A physically based variable contributing area model of basin
 hydrology, Hydrol. Sci. Bull., 24, 43-69, doi:10.1080/02626667909491834,1979
- 525
- Boucot, A.J., Chen X., and Scotese, C.R,.: Phanerozoic Paleoclimate: An Atlas of Lithologic
 Indicators of Climate, SEPM Concepts in Sedimentology and Paleontology, (Digital
 Version), No. 11, ISBN 978-1-56576-281-7, Society for Sedimentary Geology, Tulsa, OK,
 478 pp., 2013.
- 530
- Cao, M., Marshal S., and Gregson, K.: Global carbon exchange and methane emissions from
 natural wetlands: Application of a process-based model, Journal of Geophysical Research,
 101, 14399-14414, doi.org/10.1029/96JD00219, 1996
- Carmicheal, M.J., Gordon, N.I., Badger, M.P.S, Naafs, B.D.A., Behrooz, L., Remmelzwaal,
 S., Monteiro, F.M., Rohrssen, M., Farnsworth, A., Buss, H.L., Dickson, A.J., Valdes, P.J.,
 Lunt, D.J., and Pancost, R.D.: Hydrological and associated biogeochemical consequences of
 rapid global warming during the Paleocene-Eocene Thermal Maximum, Global and Planetary
- 538 Change, 157, 114-138, doi:10.1016/j.gloplacha.2017.07.014, 2017.
- 539
- Global Land Cover 2000 database. European Commission, Joint Research Centre,
 http://forobs.jrc.ec.europa.eu/products/glc2000/glc2000.php, 2003, accessed 2005.
- 542
- 543 IPCC, 2013: Climate Change 2013: The Physical Science Basis. Contribution of Working
 544 Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change
- [Stocker, T.F., D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia,
 V. Bex and P.M. Midgley (eds.)]. Cambridge University Press, Cambridge, United Kingdom
- 547 and New York, NY, USA, 1535 pp.
- 548
- Lehner, B. and Doll, P.: Development and validation of a global database of lakes, reservoirs
 and wetlands, Journal of Hydrology, 296, 1-22, doi:10.1016/j.jhydrol.2004.03.028, 2004
- Lu, X. and Zhuang, Q.: Modeling methane emissions from the Alaskan Yukon River basin,
- 553 1986–2005, by coupling a large-scale hydrological model and a process-based methane
- 554 model, J. Geophys. Res, 117, G02010, doi:10.1029/2011JG001843, 2012.

- 555
- Melton, J. R., Wania, R., Hodson, E. L., Poulter, B., Ringeval, B., Spahni, R., Bohn, T., Avis,
- 558 C. A., Beerling, D. J., Chen, G., Eliseev, A. V., Denisov, S. N., Hopcroft, P. O., Lettenmaier,
- D. P., Riley, W. J., Singarayer, J. S., Subin, Z. M., Tian, H., Zürcher, S., Brovkin, V., van
- 560 Bodegom, P. M., Kleinen, T., Yu, Z. C., and Kaplan, J. O.: Present state of global wetland
- 561 extent and wetland methane modelling: conclusions from a model inter- comparison project
- 562 (WETCHIMP), Biogeosciences, 10, 753–788, doi:10.5194/bg-10-753-2013, 2013.
- 563 Papa, F., C. Prigent, F. Aires, C. Jimenez, W. B. Rossow, and E. Matthews. Interannual
- variability of surface water extent at the global scale, 1993–2004, J. Geophys. Res, 115,
 D12111, doi:10.1029/2009JD012674, 2010
- Poulter, B., Bousquet, P., Canadell, J. G., Cias, P., Peregon, A., Saunois, M., Vivek, K. A.,
- 567 Beerling, D., Brovkin, V., Jones, C. D., Joos, F., Gedney, N., Ito, A., Kleinen, T., Koven, C.,
- 568 McDonald, K., Melton, J. R., Peng, C., Peng, S., Prigent, C., Schroder, R., Riley, W., Saito,
- 569 M., Spahni, R., Tian, H., Taylor, L., Viovy, N., Wilton, D., Wiltshire, A., Xu, X., Zhang, B.,
- 570 Zhang, Z., and Zhu, Q.: Global wetland contribution to 2000-2012 atmospheric methane
- 571 growth rate dynamics, Environ. Res. Lett., 12, 094013, doi:10.1088/1748-9326/aa8391, 2017
- Prigent, C., F. Papa, F. Aires, W. B. Rossow, and E. Matthews. Global inundation dynamics
 inferred from multiple satellite observations, 1993–2000. J. Geophys. Res., 112, D12107,
 doi:10.1029/2006JD007847, 2007.
- 575
- 576 Schroeder, R., McDonald, K. C., Chapman, B. D., Jensen K., Podest, E., Tessler Z. D., Bohn,
- 577 T. J., and Zimmermann, R.: Development and Evaluation of a Multi-Year Fractional Surface
- 578 Water Data Set Derived from Active/Passive Microwave Remote Sensing Data, Remote
- 579 Sensing, 7, 16688-16732, doi:10.3390/rs71215843, 2015.
- 580 Sloan, L. C., Walker, J. C. G., Moore Jr, T. C., Rea, D. K., and Zachos, J. C.: Possible
- methane-induced polar warming in the early Eocene, Nature, 357, 320-322,
 doi:10.1038/357320a0 1992.
- 583 Stocker, B. D., Spahni, R. and Joos, F.,: DYTOP: a cost efficient TOPMODEL
- implementation to simulate sub-grid spatio-temporal dynamics of global wetalnds-wetlands
 and peatlands, Geosci. Model Dev., 7, 3089-3110, doi:10.5194/gmd-7-3089-2014, 2014.
- 586
- Valdes P. J., Beerling D. J. and Johnson, C. E.,: The ice age methane budget, Geophys. Res.
 Lett., 32, L02704, doi:10.1029/2004GL021004, 2005.
- 589
- 590 Valdes, P. J., Armstrong, E., Badger, M. P. S., Bradshaw, C. D., Bragg, F., Davies-Barnard,
- 591 T., Day, J. J., Farnsworth, A., Hopcroft, P. O., Kennedy, A. T., Lord, A. S., Lunt, D. J.,
- 592 Marzocchi, A., Parry, L. M., Roberts, W. H. G., Stone, E. J., Tourte, G. J. L., and Williams, J.
- H. T.: The BRIDGE HadCM3 familiy of climate models: HadCM3@Bristol v1.0, Geosci.
- 594 Model Dev., 10, 3715-3743, doi:10.5194/gmd-10-3715-2017, 2017

Wania, R., Ross, I., and Prentice, I. C.: Integrating peatlands and permafrost into a dynamic
global vegetation model: 1. EvalulationEvaluation and sensitivity of physical land surface
processes, Global Biogeochem. Cycles, 23, GB3014, doi:10.1029/2008GB003412, 2009.

- 599
- Wei, Y., Liu, S., Huntzinger, D. N., Michalak, A. M., Viovy, N., Post, W. M., Schwalm, C. 600 R., Schaefer, K., Jacobson, A. R., Lu, C., Tian, H., Ricciuto, D. M., Cook, R. B., Mao, J., and 601 Shi, X.: The North American Carbon Program Multi-scale Synthesis and Terrestrial Model 602 Intercomparison Project - Part 2: Environmental driver data, Geoscientific Model 603 Development, 7(6), 2875-2893, doi:10.5194/gmd-6-2121-2013, 2014 604 605 Westermann, P .: Temperature regulation of methanogenesis in wetlands, Chemosphere, 26, 606 321-328, doi:10.1016/0045-6535(93)90428-8, 1993. 607 608 Whiticar, M. J.: Carbon and hydrogen isotope systematics of bacterial formation and 609 oxidation of methane, Chem. Geol., 161, 291-314, doi:10.1016/S0009-2541(99)00092-3, 610 611 1999. 612 Woodward, F., Smith, T. and Emanual, W.: A global land primary productivity and 613 phytogeography model, Glob. Biogeochem. Cycles, 9, 471-490, 1995 614 615 Zachos, J. C., Dickens, G. R., and Zeebe, R. E.: An early Cenozoic perspective on 616 greenhouse warming and carbon-cycle dynamics, Nature, 451, 279-283, 617 618 doi:10.1038/nature06588, 2008. 619



621 Figure 1: Annual <u>monthly</u> maximum observed FW from the SWAMPS-GLWD data set

(Poulter et- al., 2017), mean of 2000 to 2012. Grey shading indicates bare land, as

predicted by SDGVM, or > 50% cultivated (Global Land Cover 2000 database, 2003).



Figure 2: Seasonal mean FW. Observed interpolated to model grid; (a) Jun–Jul–Aug
and (b) Dec–Jan–Feb. 1NN prediction by SDGVM (c) Jun–Jul–Aug and (d) Dec–Jan–

- and (b) Dec–Jan–Feb. 1NN prediction by SDGVM (c) Jun–Jul–Aug and (d
 Feb. 1NN prediction by LPJ (e) Jun–Jul–Aug and (f) Dec–Jan–Feb.
- 629



631 Figure 3: Monthly zonal variations of FW calculated for the mean 2000-12 climate on a

- **2.5** x 3.75° grid, (a) North of 30° N and (b) South of 30° N.







Figure 4: Monthly zonal variations of wetland CH4 <u>emissions / Tg CH4</u> calculated from
DGVM model data and observed or modelled FW, for the mean 2000-12 climate on a
2.5 x 3.75 ° grid. (a) SDGVM North of 30° N, (b) LPJ north of 30° N, (c) SDGVM South
of 30° N and (d) LPJ south of 30° N.





641 642 Figure 5: Locations of Eocene coal deposits plotted on our Eocene model land mask.

indicates an Eocene coal deposit location (Boucot et al., 2013) 643



Figure 6: Variations of statistics for match between Eocene maxKNN predicted high
FW and coal locations (Boucot et al., 2013). f1 is the mean distance of a coal location to
site with FW > 0.2 for model based on (a) SDGVM and (b) LPJ. f2 is the mean FW of
sites within 2 pixels of a coal location for model based on (c) SDGVM and (d) LPJ data.



650

651 Figure 7: Annual maximum FW calculated by the max3NN method by SDGVM and

652 LPJ for the Eocene climate, compared with coal deposit locations







Figure 8: Monthly variations of total wetland area calculated for the Eocene climate by
SDGVM and LPJ, for (a) all areas north of 30° N, (b) all areas between 30° S and 30° N
and (c) all areas south of 30° S.





660 -Figure 9: Monthly variations of wetland CH4 <u>emissions /Tg CH4</u> calculated from

- 661 predicted FW, for the Eocene climate by SDGVM and LPJ, for (a) all areas north of 30°
- 662 N, (b) all areas between 30° S and 30° N and (c) all areas south of 30° S.



Figure 10: Seasonal mean FW predicted for the Eocene climate by SDGVM and LPJ
using the max3NN (a) SDGVM June–July–August, (b) SDGVM December–January–
February, (c) LPJ June–July–August, (d) LPJ December–January–February



Eocene climate, using evapotranspiration from (a) SDGVM or (b) LPJ.

669

Data set	<u>Time</u>	Climate data source	DGVM used
SDGVM reference	Modern day	<u>CRU-NCEP v4.0</u>	SDGVM
LPJ reference	Modern day	<u>CRU-NCEP v4.0</u>	<u>LPJ</u>
SDGVM modern test	Modern day	Interpolated CRU-NCEP v4.0	SDGVM
LPJ modern test	Modern day	Interpolated CRU-NCEP v4.0	<u>LPJ</u>
SDGVM Eocene test	Early Eocene	HadCM3BL-M2.2	SDGVM
LPJ Eocene test	Early Eocene	HadCM3BL-M2.2	LPJ

Table 1. Summary of reference and test data sets used combining data from dynamic

global vegetation models SDGVM (Woodward et al., 1995; Beerling and Woodward, 2001) and LPJ (Wania et al., 2009) with climate data from CRU-NCEP v4.0 (Wei et al. 2014), for the modern day, and HadCM3BL-M2.2 (Valdes et al. 2017), for the Early

Eocene.

	> 30° N FW	< 30° N FW	Global FW
Observed	1.84	2.11	3.95
Observed	1.47	1.41	2.88
excluding SDGVM bare land			
SDGVM	1.53	1.47	3.00
LPJ	1.95	1.90	3.86

Table 21: Modern day monthly mean FW area (10⁶ km²), for observed data interpolated to the 2.5° x 3.75° grid or calculated by vegetation model.

Model	FW data	> 30° N CH4	< 30° N CH4	Global CH 4
SDGVM	observed	64.32	122.69	187.01
	predicted	57.95	108.63	166.58
LPJ	observed	65.43	68.60	134.03
	predicted	73.11	83.78	156.89
GCP-CH4* WETCHIMP**	observed 0.5° model specific	51±15	126±31	~ 184 190±39

683	* GCP-CH4 (Poulter et al., 2017) results are the mean of 11 different methane emission
684	models with the same observed wetland data as used to produce Figure 1 here. They are
685	quoted as means over specific ranges of years: $2000-2006 = 184.0 \pm 21.1$, $2007-2012 = 2000-2006 = 184.0 \pm 21.1$
686	$183.5 \pm 23.1, 2012 = 185.7 \pm 23.2$. As our results are for a single mean 2000–12 year we
687	therefore only quote an approximate value from this source for comparison.

** WETCHIMP (Melton et al., 2013) results are the mean of 8 different models, 1993-2004,
 each of which used their own definition of wetland extent rather than observed data

690

Table <u>32</u>: Modern day annual total wetland CH₄ emission (Tg CH₄ year⁻¹), calculated

by vegetation model using either observed FW data (interpolated to the $2.5^{\circ} \times 3.75^{\circ}$

693 grid) or model predicted FW, compared with other modelling studies.

FW model	> 30°N	30°S to 30°N	< 30°S	Global
SDGVM	2.82	4.11	1.53	8.48
LPJ	4.84	3.39	2.06	10.29

696 Table <u>43</u>: Eocene monthly mean max3NN modelled FW area / 10⁶ km²

Final response to Referees #1 and #2

We thank both anonymous referees for their comments and recommendations. Below we give our responses, in italics, to each of those and indicate where revisions have been made to the manuscript. The line numbers given refer to the above tracked changes version of the manuscript.

Response to Referee 1

The main confusion I have is on the validation of this approach. It is not convincing that using one reference dataset to train their algorithm, and then evaluate the simulated results with the same reference dataset. It would be necessary to compare with independent inundation products to justify their approach, or the authors need to provide the uncertainty in the estimated inundation using their approach given that there are large uncertainties in wetland extent among existing inundation products (Melton et al., 2013).

There is no training and evaluation in the sense that would normally be understood from a machine learning perspective. For the Eocene results, section 3.2, we clearly have no wetland data with which to train and evaluate our predictions. We simply use the coal deposits as a proxy, comparing those to our wetland predictions to give us the best value of K for the maxKNN approach with this particular data set.

We have added to the text in section 3.2.1 to make this clearer, L362 – 364.

Nor are we using a training set for the modern day test data, section 3.1. These results were included simply to show whether some form of nearest-neighour approach might, in principle, be useful (lines 236-238); we were exploring the potential of this approach. It was a test that if failed would have meant we would not have continued developing a nearest neighbour method; it would have been another unsuccessful attempt along with those briefly discussed in section 2.3. That the method passed this test merely indicated we could explore some form of nearest neighbour method in the context of the Eocene climate.

We have added to the text at the end of section 3.1 to stress this point, L315-317

The logic of this approach is a bit confusing to me. If I understand it correctly, this nearest neighbor-based algorithm implicitly assumes the locations of wetlands should close to each other and inundation is correlated with eight variables the authors proposed. But according to the modern dataset, is there any analysis/evidence prove that this relationship exist

The nearest neighbour approach assumes that sites with similar values of wetland fraction should have some similarity in terms of their values of the 8 climate & vegetation variables we use; or to put it another way, if sites with similar FW show no similarity at all between their values of at least some of those 8 variables, then a nearest neighbour approach will simply not work. There is certainly no simple correlation between FW and those 8 variables in the modern day data, as we briefly explain in our "Initial unsuccessful models" section 2.3; a multiple-linear regression on those 8 variables did not produce a good predictive model of FW. This suggests that any relationship between FW and those 8 variables must be complex.

We have added a sentence to the end of section 2.3 to reflect this, L206 - 209

We have also the revised section 2.4 text should make this clearer, L213 - 227

Fan (2011) suggest that water table depth is a key to simulate wetland distribution - at least it is an important variable to capture the distribution of peatlands in high latitudes as some of the peatlands don't show inundated condition but still emit CH4.

We use soil water content, defined as the amount of water in in the top 1m of soil. This is produced by both vegetation models whereas water table depth is not.

I'm not sure that comparing the simulated wetland distribution with coal deposit can be helpful as the authors have already mentioned some of the limitations using coal deposit. Also, it's hard to tell how good the fit is from reading Figure 7.

Clearly coal deposits are not an ideal proxy for wetland fraction, but they are all we have. Without them we would have had no way of deciding on a value for K in the maxKNN algorithm. Therefore, despite the limitations, they are useful to explore this approach.

It would be great to address a bit more about the background why it's important to develop a dynamic inundation algorithm for deep time paleoclimate simulation and what's the current status of research on this topic.

As explained in the introduction, there is great interest in understanding how the extent of wetlands changed through geological time and what role that could have had on methane cycling. However, there is currently only one model-based approach for deep time paleoclimates (Beerling et al., 2011). The goal of this paper is to explore other methodologies and compare them to this original work, better understanding the potential of the new approaches and the robustness of the previous work.

We have revised the text in a number of places in the introduction, L41-44, 58-59, 62-63, 71-73.

Response to Referee 2

"I am fairly convinced this is a sensible and useful approach, but I must admit to being slightly baffled about the exact methods employed – I found the paper rather unclear in quite a few places. I would encourage the authors to revise the description of the methods to make it clearer. "

In addition to the changes in the introduction recommended by referee 1 we have further added to and revised the text outlining the structure of the paper here, L93 – 110. We have also revised the description of the nearest neighbour method, L213 - 227

Some clarification on how this approach should be employed by the wider modelling community would also be appreciated – can the method be embedded within ESMs to calculate wetland emissions online? Or is it envisaged as an offline only tool? I wasn't clear. If the authors can clarify the methods This has been used as an offline tool with the emphasis in this work on exploring methods for predicting FW. We have added to the conclusion with respect to this point and also added that further work on other paleoclimates and other data would be naturally be of interest, L472-475 and L479-483

L43 ESMs must either prescribe CH4 concentrations as boundary conditions, or "incorporate dynamic methane fluxes from natural sources: : :". If the latter, they must not only simulate the sources but also the sinks of the CH4 (i.e the whole budget) in order to reasonably represent concentrations.

Changed to "... sources and sinks", L47

L55 ': : :no direct observations of wetland extent' – it should be stated that there are however proxies, that you later utilise (i.e. coal deposits).

Amended the text to refer to coal deposits as a proxy, L58-59

L60 ': : :mean monthly temperature drops below 0 _C at some point in the year: : :' I found this slightly confusing. Do you mean if there is one (or more) month in the year below 0 _C, then that grid-cell is classified as producing methane? Clarify.

Changed to "mean monthly temperature drops below 0°C for at least one month of the year" L65-66

L71 So you are using DGVMs to simulate vegetation distributions, rather than using present-day observational datasets. It may be worth saying that the DGVMs have (presumably) been evaluated elsewhere.

We have now named and included the references for the DGVMS used at this point, L78-80

L68 Is it worth briefly defining wetland? Perhaps earlier. E.g. the RAMSAR definition. Is it obvious how such definitions translate into a climate model-specific definition? (Water depth, etc.). What is the basis of the modern day reference data set of FW? Can you say it is 'known' or 'observed' FW?

We have added a sentence here to say how we define wetlands, L80-83 This is the same definition as for our reference data set. However wetland may be defined, what we are predicting is wetland as defined by the reference data.

We have also explained the use of 'observed' to distinguish the reference data from our model results L117-118.

That the reference data is partly derived from satellite observations was already included in the text, in response to a later point we have now also included further detail on the other observations it is derived from, L 123-124

L80 Typo: intercomparson

Corrected, L92

L87/90 Capitalise Nearest Neighbours or not? (Is it a well enough known method to be considered a proper noun? I don't know, but at least be consistent.)

Changed to nearest neighbours, L102

L105 So SWAMPS is based on microwave satellite observations – what is the observational data that GLWD is based on?

Added that it is static inventory of wetland area, L123-124

L127 I didn't fully understand the scaling – are the mean/standard deviation global values?

They are global, added this L149

L130 As previous comment – is the global mean 0?

They are global, added this L152 Also add to later discussion of the test data sets L157

L132 A modern-day test data set...

Added 'data', L154

L134 conducted on -> driven by?

Changed L156

L136 Use the same terminology as l132 to avoid confusion, i.e.: "The paleoclimatic assessment of our model was performed using an early Eocene: ::" -> An early Eocene test data set was made using: ::?

Changed L159-160

L145 It would be useful to provide a summary table of the test/reference data sets to clarify exactly how you are going to evaluate your approach; I didn't find the current explanation completely clear.

Added sentence referring to new Table 1, L171-172 New Table 1 appears at L672 – 676 Subsequent changes to previous tables, increasing all numbers by 1 The tables section L678, 690, 695 In the main text where these are referred to L292, 297, 343, 463 L163 The number of what? Months or grid cells?

It is grid cells, text amended L187

L214 Is Rh an absolute or scaled (0-1) value? If absolute, what are the units? Similarly for GPP in the next equation (I guess it must be absolute value to make sense.)

It is absolute and units of g C m^{-2} month⁻¹ in both cases, text amended L249-250

L218 Is TMP soil, surface, or surface air temperature?

It is air temperature, added L251

L226 Presumably me >= 0? Is there a test for mp >= mo? What are the units of me?

mo is defined as fraction of mp by eq(4), therefore me, by eg (6), has to be >= 0, so no test for mp > mo is done. The units of me, mo and mp are all g CH₄ m⁻²month⁻¹ these have been added earlier where mp and mo defined L220.

L232 'downscaled' – I think the definition of downscaling is to infer something at high resolution from something at low resolution. You seem to be using the word in the opposite sense. I don't think we (scientists) normally use 'upscaled', so I am unsure what to call this (degrading?), but I don't think it is downscaling (also I234).

Changed "downscaled" to "interpolated" L266

L318 I think the term 'maxKNN' appears here for the first time and isn't defined. Is it just KNN with K>1? (As suggested by I316.)

We have changed the section title to "Maximum of K nearest neighbours FW prediction" L356

We have now defined maxKNN later L360

L364 In a similar vein to the last comment - why not just 3NN rather than max3NN?

"3NN" does not indicate it is the maximum of those 3 nearest neighbours, it could imply any relationship to the three nearest neighbours

L383 ': : :both [FW and CH4 emissions] have their highest values in summer months: : :' This is not so clear in Figure 8 for SDGVM. It is clear in Figure 9.

We have changed the text to " ... values in the late spring and summer months .. " L423

L409 ': : :their respective impacts of soil water balance: : :'. Clarify. Is this just a typo of -> on?

Corrected to "on" L450

L409 I got a bit confused here about EVT. It seems EVT is from the vegetation models; but EVT must also be calculated in the underlying climate model – I guess with a much

more simplified vegetation scheme. Is there a large discrepancy between the EVT in the vegetation and climate models? Isn't this a bit of a problem? This decoupling of the simulated water budget between the climate model and the vegetation model should be clearly explained earlier in the methods section, and the implications discussed here.

EVT, as used throughout the paper, is always from the vegetation models. We have not considered EVT from the climate model. So long as it is the same EVT used at all times in our modelling of FW, i.e. same definition for reference and test data sets, this should not be an problem.

We have clarified this in section 2.1, L138-141 as well as in section 3.2.2 L451-452

L423 Global monthly mean FW for the Eocene: : :

Inserted "for the Eocene" L465

L572 Figure 1 caption – Annual monthly maximum: : :

Inserted "monthly" L621

L586 Figure 4 caption and y-axes – clarify these are CH4 emissions – what are the units? (Tg CH4/month?)

Captions and figure titles (there are no y-axes labels on this and similar figures as it makes them too crowded) changed to " CH_4 emissions / Tg CH_4 ", L634-636

We have made similar changes made to Figure 9, L658-660

L628 Incorrect punctuation for list.

Changed ; to : L684

Additional

Deleted repeated spaces between the following words on the given lines

"fraction" and "predicted", L21

"fraction" and "methane", L34

"data" and "at", L93

"models." and "The", L155

"deposits." and "The", L407

Added a reference in the Abstract for the HadCM3BL-M2.2 climate model, L25

Added "a", L75

Deleted ":", L296

Added an acknowledgement to the referees, L505-506

Corrected some citations. Removed extra "." or added missing "," in "et al.," L35, 36, 346, 622

Corrected misspellings of

"wetlands", L584

"Evaluation", L597