# Observational benchmarks inform representation of soil organic carbon dynamics in land surface models

Kamal Nyaupane[1], Umakant Mishra[2*], Feng Tao[3], Kyongmin Yeo[4], William J. Riley[5], Forrest M. Hoffman[6] and Sagar Gautam[2]


[1]Environmental science and Engineering Program, The University of Texas at El Paso, El Paso, TX 79968, United States.

[2] Biomaterials & Biomanufacturing, Sandia National Laboratories, Livermore, CA, 94550, United States.

[3]Ministry of Education Key Laboratory for Earth System Modelling, Department of Earth System Science, Tsinghua University, Beijing, 100084, China

[4]IBM Thomas J. Watson Research Center, Yorktown Heights, NY, 10562, United States.

[5]Earth & Environmental Sciences, Lawrence Berkeley National Laboratory, Berkeley, CA, 94720, United States.

[6]Climate Change Institute, Oak Ridge National Laboratory, Oak Ridge, TN, 37830, United States.


*Corresponding author Email: umishra@sandia.gov

Biogeosciences
Discussions

1  **Abstract**

2  Representing soil organic carbon (SOC) dynamics in Earth system models (ESMs) is a key

3  source of uncertainty in predicting carbon climate feedbacks. Machine learning models can help

4  identify dominant environmental controllers and their functional relationships with SOC stocks.

5  The resulting knowledge can be implemented in ESMs to reduce uncertainty and better predict

6  SOC dynamics over space and time. In this study, we used a large number of SOC field

7  observations (n = 54,000), geospatial datasets of environmental factors (n = 46), and two

8  machine learning approaches (Random Forest (RF) and Generalized Additive Modeling (GAM))

9  to: (1) identify dominant environmental controllers of global and biome-specific SOC stocks, (2)

10  derive functional relationships between environmental controllers and SOC stocks, and (3)

11  compare the identified environmental controllers and predictive relationships with those in

12  Coupled Model Intercomparison Project phase six (CMIP6) models. Our results showed that

13  diurnal temperature, drought index, cation exchange capacity, and precipitation were important

14  observed environmental controllers of SOC stocks. RF model predictions of global-scale SOC

15  stocks were relatively accurate ($R^2 = 0.61$, RMSE = 0.46 kg m$^{-2}$). In contrast, precipitation,

16  temperature, and net primary productivity explained >96% of ESM-modeled SOC stock

17  variability. We also found very different functional relationships between environmental factors

18  and SOC stocks in observations and ESMs. SOC predictions in ESMs may be improved

19  significantly by including additional environmental controls (e.g., cation exchange capacity) and

20  representing the functional relationships of environmental controllers consistent with

21  observations.

22

23  **Keywords**: Environmental controllers, Earth system models, soil organic carbon, net primary

24  productivity, machine learning, model benchmarking

25

26  **1. Introduction**

27  Soil is the largest actively cycling carbon pool in terrestrial ecosystems and stores almost twice

28  the amount of carbon as in the current atmosphere (Lal, 2016). A small change in soil carbon

29  stocks can lead to large changes in the atmospheric $CO_2$ concentration and future climate change

30  trajectories. Soils also play a crucial role in sequestering atmospheric $CO_2$ as soil organic carbon

31  (SOC) (Hinge et al., 2018). Thus, sequestration, protection, and sustainable management of SOC

32  stocks can be a promising climate mitigation strategy (Lal, 2020). Accurate representation of

33  global SOC storage and its environmental controllers are essential for predicting realistic

34  changes of SOC under different land use and climate change scenarios. Yet, no consensus exists

35  among current Earth system models (ESMs) in representing the spatial distributions of global

36  SOC storage and its fate under future climate change scenarios (Friedlingstein et al., 2014.;

37  Arora et al., 2020).

38  Multiple environmental variables, including climatic and topographic factors, land use history,

39  and edaphic properties, have been identified as possible controllers of SOC storage (Georgiou et

40  al., 2021; Mishra et al., 2022). Current ESMs, however, use the effects of only a limited number

41  of environmental factors in representing SOC storage and dynamics. A recent study that

42  compared SOC stocks from multiple ESMs against observations indicated a large knowledge gap

43  in both ESMs and observations (Georgiou et al., 2021). Therefore, it is important to compare

44  ESM simulations against global SOC observational datasets to evaluate model performance and

45  identify key environmental controllers of global SOC storage.

46   Benchmarking ESM simulations with observed data is a common approach for model evaluation

47   (Luo et al., 2012; Todd-Brown et al., 2013; Collier et al., 2018). Through comparing model

48   simulations with observations, model strengths, deficiencies, and needed improvements can be

49   identified. The resulting understanding from SOC benchmarking could lead to new ESM land

50   model structures (by identifying key processes) and new parameterizations (by quantifying key

51   relationships between SOC and environmental variables). Thus, benchmarking analysis of ESMs

52   is an effective tool to reduce uncertainties in predicting SOC dynamics and can provide more

53   realistic information for managing SOC under changing climate conditions (Lauer et al., 2017).

54   Currently ESMs predict SOC stocks primarily with model representations that depend on soil

55   temperature, moisture, and belowground net primary production (Todd-Brown et al., 2013).

56   ESMs capture the positive correlation between NPP and precipitation, resulting in high SOC

57   stocks for areas with high NPP in moist regions (Sun et al., 2016). Higher temperature increases

58   soil respiration, which, in the short-term, reduces SOC storage. In the longer-term, increased soil

59   respiration can release nutrients, leading to increased plant growth, belowground carbon inputs,

60   and thereby SOC stocks; the balance of these factors can take centuries to manifest (Mekonnen

61   et al., 2022). Soil respiration temperature sensitivity is often defined based on $Q_{10}$ or Arrhenius

62   equations in ESMs (Wynn et al., 2006), although low- and high-temperature modifications to

63   these relationships are likely needed (Jiang et al., 2013; Azizi-Rad et al., 2022).

64   In a previous U.S. continental-scale study, we derived empirical non-linear relationships between

65   SOC and environmental factors that produced comparable prediction accuracy to a random forest

66   (RF) machine learning approach (Mishra et al., 2022). We apply a similar approach in this study

67   in both global field observations and ESMs to (1) identify key observed environmental controllers

68   of, and functional relationships with, global SOC stocks and (2) evaluate ESMs with these

69    observational benchmarks. Simulated SOC stocks from three CMIP6 ESMs (i.e., Community

70    Earth System Model (CESM, Hurrell et al., 2013); U.K. Earth System Model (UKESM, Sellar et

71    al., 2019); Beijing Climate Center model (BCC, Xiao-Ge et al., 2019) were benchmarked with

72    50,000 SOC profile observations across the globe. We used a machine learning (i.e., random

73    forest) approach with 46 environmental factors to identify the key environmental controllers of

74    SOC stocks at the global scale. We then applied a generalized additive model (GAM) to derive the

75    predictive relationships between these key environmental factors and SOC stocks in observations

76    and ESM simulations.   Specific objectives of this study were to: (1) identify dominant

77    environmental controllers of SOC stocks in field observations and CMIP6 ESMs, (2) derive

78    observed and ESM-modeled functional relationships between environmental factors and SOC

79    stocks, and (3) analyze these functional relationships to inform needed improvements in ESM

80    representations of SOC dynamics.

81    **2. Materials and Methods**

83    *2.1 Soil organic carbon stock observations*

85    We used two datasets of SOC stocks for the topsoil layer (i.e., $0 - 30$ cm) and the whole soil profile

86    (i.e., $0 - 100$ cm). The World Soil Information Service (WoSIS) compiled SOC profiles across the

87    globe after quality assessment. The 2019 snapshot of the WoSIS dataset contained 111,380 soil

88    profiles with SOC content information (unit: g C g-soil$^{-1}$) at different soil depths (Batjes et al.,

89    2020). We estimated the SOC stock (g C m$^{-2}$) at different soil layers using:

90    $$SOC\ Stock = SOC\ Content \times \left(1 - \frac{G}{100}\right) \times BD \times D \qquad (1)$$

91    where G is the coarse fragment fraction (%); BD is the bulk density of soil (g m$^{-3}$); and D is the

92    soil layer depth (m).

93    When the measured bulk density value was absent from the dataset, we used a pedo-transfer

94    function (Yigini et al., 2018) to estimate the soil bulk density:

95    $BD = \alpha + \beta \times exp(-\gamma \times OM)$                                      (2)

96    Where OM is organic matter, equivalent to SOC×1.724, with SOC content in percent (%); α, $\beta$,

97    and $\gamma$ are fitting parameters. We found α = 0.32, $\beta$ = 1.30, and $\gamma$ = 0.0089 after fitting WoSIS data

98    to this equation.

99    Another dataset we used in this study was compiled from Mishra et al. (2021). This dataset

100   contained 2,546 soil profiles with SOC stock (g C m$^{-3}$) information from permafrost regions in

101   North America, northern Eurasia, and the Qinghai-Tibet Plateau. In total, we used 113,926 soil

102   profile observations from these two data sources. SOC stocks of different soil layers were then

103   summed to SOC stocks in 0 – 30 cm and 0 – 100 cm depth intervals. Because not all these soil

104   profiles covered the whole 0 – 30 cm or 0 – 100 cm intervals, we used a total of 54,000 soil profiles

105   that included SOC stock information for both depth intervals. The geographical distributions of

106   soil profiles used in this study are shown in Figure 1. Because SOC stock values across the globe

107   were highly skewed, we used a natural logarithm transformation in this study.

108
109   ***2.2 Environmental predictors of SOC stocks***
110
111   The storage and cycling of SOC are controlled by multiple environmental factors. In this study,

112   we used observations of 46 environmental variables, which represented major soil forming factors

113   (McBratney et al., 2003.). Twenty-one of the 46 environmental variables were climatic variables,

114   including annual average temperature, precipitation, evapotranspiration, drought severity index,

115   and statistics for different temporal scales (e.g., during the wettest and driest quarter in a year).

116   Thirteen of the 46 variables described soil properties (e.g., clay content, sand content, silt content,

117   soil texture, pH, and cation exchange capacity). Six variables represented topographic factors (e.g.,

118  elevation and soil depth). Six variables represented land use and land cover types. All the

119  categorical variables were converted to integer variables and the environmental variables were

120  resampled to a common 1 km resolution. The environmental factors, their original spatial

121  resolution, and data sources are provided in the supporting information (Table S1).

122

123  ***2.3 Selection of dominant environmental controllers of SOC stocks***
124
125  We used RF to select dominant environmental predictors of SOC stocks within biomes and at

126  global scale in both observations and ESMs. RF is an ensemble learning method, which is an

127  extension of the classical Classification and Regression Trees (CART). Building a collection of

128  uncorrelated CARTs through bootstrapping the samples and applying the random subspace method

129  at each branch of the trees, RF improves the prediction performance (Breiman, 2001; Wiesmeier

130  et al., 2011; Mishra et al., 2020). RF is well known for its strength in modeling highly nonlinear

131  relationships between the predictors and is robust to overfitting (Chagas et al., 2016). Moreover,

132  RF is not very sensitive to the choice of the hyperparameters, which makes RF one of the most

133  popular off-the-shelf model for many classification and regression problems.

134  In this study, we trained the RF model using SOC content as a response variable and environmental

135  factors as predictors. The model performance was evaluated using the coefficient of determination

136  ($R^2$) and root mean square error (RMSE). A 10-fold cross-validation was used to compute $R^2$ and

137  RMSE. Biome-specific analyses were conducted on a subset of the global dataset. For biome

138  classification, we used the IGBP land classes (Loveland and Belward, 1997). The "Random-

139  Forest" package in R was used to train a RF model using all the observed environmental factors in

140  the dataset and to identify dominant environmental controllers of SOC stocks. Prior to fitting into

141  the final model, we performed a potential collinearity test among the environmental variables by

142    calculating pairwise correlations and variance influence factors. Predictors showing a variance

143    influence factor (VIF) value greater than 10 were omitted, leaving 14 uncorrelated environmental

144    predictors of SOC stocks in the observations.

145

146    ### 2.4 Generalized additive model
147

148    Generalized additive model (GAM) is an extension of generalized linear models, which employs

149    spline functions to model nonlinear relationships between predictor and response variables (Arnold

150    et al., 2013).  In GAM, the relationship between predictor and response variable can be modeled as

151    (Hastie and Tibshirani, 1987):

152    $Y = C + \sum_{i=1}^{p} f_i(X_i)$                              (3)

153    Here, Y is the response variable (SOC), $C$ is a constant, $X_i$ are the environmental controller

154    variables, $f_i$ is a spline function for $X_i$, and p is the total number of environmental controllers. We

155    used the "mgcv" package in R to build GAMs for the observations as well as CMIP6 ESMs

156    (Arnold et al., 2013). The performance of GAMs was evaluated by using $R^2$ and RMSE.

157
158    ### 2.5 Earth system model outputs
159
160    We downloaded and aggregated the SOC and environmental controller data from three ESMs that

161    participated in CMIP6: Community Earth System Model (Hurrell et al., 2013.), U.K. Earth System

162    Model (Sellar et al., 2019), and Beijing Climate Center model (Xiao-Ge et al., 2019). These ESMs

163    included most of the environmental factors used by CMIP6 ESMs. ESMs did not report depth-

164    dependent soil carbon projections, making direct comparison with depth-dependent SOC

165    observations difficult. The majority of land models used in ESMs were designed to simulate topsoil

166 carbon for topsoil depth; thus, we assumed that the simulated soil carbon is contained within 1 m

167 of soil profile to simplify comparison with observations.

168

169 3.    **Results**
170
171 *3.1    Descriptive statistics of SOC observations*
172
173 The average global SOC stock in the 0 - 1 m depth interval was 13.5 kg C m$^{-2}$, ranging from 0.14-

174 435.3 kg C m$^{-2}$. Summary statistics of SOC stocks at global scale and within different biomes is

175 presented in Table 1. The standard deviation showed a similar spread in SOC stock values in

176 croplands (n=21820), savannas (n=9807) and grasslands (n=5938). However, in forests (n=12164)

177 and shrublands (n=3769), the standard deviation was higher indicating a large range in SOC stock

178 values. Distributions of total SOC stocks in different biomes are presented in Figure 2. Across

179 different biomes, forests contain the largest organic carbon content globally, with a mean value of

180 15.9 kg C m$^{-2}$ and standard deviation 20.7 kg C m$^{-2}$.

181
182 *3.2    Dominant environmental controllers of SOC stocks in observations and ESMs*
183

184 At the global scale, we found that diurnal temperature, drought severity index, annual

185 temperature, and cation exchange capacity are the dominant environmental controllers of SOC

186 stocks in observations (Figure 3). By including all the environmental controllers, the RF model

187 explained 61% of observed global spatial SOC variation. $R^2$ ranged from 48% in savannas to

188 65% in croplands (Table 2) and the importance of key environmental controllers varied between

189 biomes (Figure 4). In croplands, precipitation, drought, diurnal temperature, and cation exchange

190 capacity were identified as the dominant controllers of SOC stocks. In grasslands, annual

191 temperature, cation exchange capacity, and sand content were the dominant controllers. In

192 forests, cation exchange capacity, precipitation, and temperature were dominant controllers. In

193    shrublands, annual temperature, soil pH, and cation exchange capacity were the most important

194    controllers. In savannas, soil related variables, temperature, and precipitation were the most

195    important controllers. Across all land cover types, we found that cation exchange capacity and

196    seasonal climatic variables were the dominant environmental controllers of SOC stocks.

197         In contrast, the RF model with 8 environmental variable predictors made near-perfect

198    predictions of ESM simulated SOC stocks (average $R^2 = 0.95$, $R^2$ values for UKESM, CESM,

199    and BCC model were 0.99, 0.89, and 0.98, respectively). In contrast to the results obtained from

200    the observed SOC stocks, the dominant controllers of ESM simulated SOC stocks were annual

201    temperature, net primary productivity (NPP), and annual precipitation (Figure 5). In particular,

202    NPP was by far the most dominant predictor of SOC stocks in the UKESM.

203

204    *3.2 Predictive relationships between environmental factors and SOC stocks*

205         Dominant environmental controllers of observed SOC stocks identified by the RF model

206    were used in GAM to derive predictive relationships. We retrieved explicit analytical

207    expressions by fitting the splines derived from GAM in the observation dataset. Notwithstanding

208    its role as the sole carbon source to soil, our results did not show NPP as a strong controller on

209    observed SOC stocks (Figure 6a). In contrast with field observations, all ESMs showed

210    significant dependence (exponential increase) of SOC stocks on NPP. Our results also showed

211    that observed SOC stocks increased almost linearly with observed annual precipitation (Figure

212    6b). In contrast, ESMs show different relationships between SOC and precipitation. We found a

213    nonlinearly increasing SOC with precipitation in CESM, an initial sharply increasing and then

214    decreasing relationship in UKESM, and a decreasing relationship in BCC ESM. On the

215    relationship between SOC storage and soil texture and elevation, ESMs do not capture the

216    observed relationships. Our results indicated that observed SOC stocks decreased with clay

217    content in the interval between 0 and 20%, and then increased with clay content above 20%

218    (Figure 6c). Observed SOC stocks increased with silt content up to 55% and then decreased

219    (Figure 6d).

220         SOC stock functional relationships differed between the three ESMs and in many cases

221    differed with the relationships we derived from observations. In terms of the effects of annual

222    temperature on modeled SOC storage, we found that SOC stocks decreased with annual

223    temperature and were most sensitive to temperature in the range between 0 and 10$^{o}$C (Figure 6e).

224    However, while the three ESMs captured the general negative relationship between SOC storage

225    and temperature, none of them correctly described the varying sensitivity of SOC in different

226    temperature ranges (especially in extreme temperature ranges <0$^{o}$C and >20$^{o}$C). In representing

227    the control of elevation on SOC storage, only UKESM showed consistent patterns with

228    observations, where SOC storage remained stable when the elevation was lower than 2000 m and

229    decreased when the elevation was higher than 2000 m (Figure 6f).

230

231    **Discussion**

232    Previous studies have suggested that the spatial variation of SOC is dependent on multiple

233    environmental factors such as climatic and edaphic variables, geography, and vegetation. Here,

234    we found that climatic variables (i.e., temperature and precipitation) are the most important

235    controllers of global SOC stocks, followed by edaphic variables (i.e., cation exchange capacity),

236    topography (i.e., elevation), and vegetation (i.e., NPP). Using boosted regression trees, Luo et al.

237    (2021) studied edaphic and climatic controls on SOC dynamics at different soil depths and found

238    that soil type and climatic variables are the most important variables in explaining the SOC

239    stocks (Luo et al., 2021). In this study, we found that seasonal climatic variables such as diurnal

240    temperature range and precipitation seasonality are among the most important environmental

241    controllers in explaining the spatial variation of SOC stocks. This result indicates the critical role

242    of seasonal and interannual climatic variables in understanding SOC dynamics.

243        The importance of climatic variables on global SOC storage emerges from close links

244    with processes that affect ecosystem productivity and soil microbial processes. Consistent with

245    our findings, Wiesmeier et al. (2014) reported climatic variables (temperature and precipitation)

246    as significant controllers of SOC stocks up to 1 m depth in German soils under oceanic climate

247    (Wiesmeier et al., 2014). Sreenivas et al. (2014) used RF to predict the SOC variability across

248    semi-arid and humid areas of India in the top 30 cm of soil and found that the top three

249    environmental controllers were land cover, mean temperature of hottest months, and mean

250    annual precipitation (Sreenivas et al., 2016). In our analysis, the overall relative importance of

251    climatic variables was significantly higher than other variables at the global and biome scales.

252        Soil properties were identified as the second most important controllers of global SOC

253    stocks. Soil properties impact various processes that govern soil carbon dynamics. For example,

254    soil properties impact microbial activity, porosity, and oxygen availability in the soil profile,

255    which directly or indirectly control soil water dynamics, plant growth, and SOC stocks.

256    Consistent with our findings, Luo et al. (2021) reported that sand content, silt content, and soil

257    pH were significant controllers of SOC stocks in all soil depths globally.

258        The Palmer drought severity index, which indicates low soil moisture availability, was a

259    dominant controller of global SOC stocks. Consistent with our findings, Li et al. (2021) reported

260    that soil particle size and soil water content were the most influential predictors of SOC variation

261    (Li et al., 2021). Soil drought, indicating more negative soil water potential and low soil

262    hydraulic conductivity, can cause tree mortality (Anderegg et al., 2012). Climate extremes like

263    droughts can impact the structure, composition, and functioning of terrestrial ecosystems and can

264    thereby severely affect the regional carbon cycle (Frank et al., 2015).

265        Cation exchange capacity is a soil property that indicates the active soil surface to which

266    SOC may be adsorbed, and polyvalent metal cations can play a significant role in SOC

267    stabilization by binding organic compounds to mineral surfaces (O'Brien et al., 2015; Solly et

268    al., 2020). O'Brien et al., (2015) found that exchangeable soil $Ca^{2+}$ is a significant predictor of

269    SOC stocks. This relationship is supported by the mechanism that $Ca^{2+}$ and $Mg^{2+}$ promote clay

270    flocculation and bind organic matter to clay surfaces. Solly et al. (2020) reported that SOC and

271    cation exchange capacity are significantly related in both topsoil and subsoil with strong positive

272    relationship.

273        After climatic factors and cation exchange capacity, topography and vegetation (NPP)

274    were important controllers of observed global SOC stocks. Effects of NPP on observed SOC

275    stocks was found to be small (~6% in 0-100 cm soil depth). Similar to our findings, Luo et al.

276    (2021) reported NPP explaining about 10% of the variation of SOC stocks. NPP delivers the

277    primary inputs of carbon to soil and NPP generally increases with moisture, temperature, and

278    $CO_2$ up to a certain limit (Todd-Brown et al., 2013). NPP also depends on the availability of soil

279    nutrients. Most ESMs overestimate the increase in SOC pools in response to NPP increases

280    (Todd-Brown et al., 2013). The effects of NPP on SOC also depend on biome type and soil

281    depths (Luo et al., n.d.; Georgiou et al., 2021). The contribution of NPP on SOC stocks mostly

282    depends on how much NPP ends up in the soil and how it is translocated to different soil depths.

283    Georgiou et al. (2021) reported a saturating relationship of SOC stocks with increasing NPP in a

284    global observational dataset. However, Chen et al., (2018) reported high SOC stocks with

285    increasing productivity and soil water holding capacity (Chen et al., 2018).

286          The three CMIP6 ESMs we analyzed predicted SOC stocks mostly as a function of

287    temperature, precipitation, and NPP. These ESMs simulated positive correlations between SOC

288    stocks and NPP (Figure 5a), resulting in high SOC stocks in areas with high NPP in most regions

289    (Shi et al., 2013; Sun et al., 2016). In these ESMs, effects of temperature and precipitation on

290    SOC stocks are driven by soil respiration. Most current ESMs simulate the response of soil

291    respiration to temperature using either a $Q_{10}$ or Arrhenius equation (Wynn et al., 2006), such that

292    a higher temperature causes more soil respiration, and, all else equal, eventually reduces SOC

293    stocks (Figure 5b). Our results showed diverse control of precipitation on SOC stocks in

294    different ESMs. Todd-Brown et al. (2013) showed that ESM soil respiration either increases

295    monotonically with precipitation, or first increases to a plateau under optimal precipitation and

296    then decreases with further increasing precipitation. Consistent with those results, the ESMs we

297    analyzed in this study showed different dependence of SOC storage on annual precipitation.

298          In this study, we found that, in comparison to the patterns that emerged from

299    observations, ESMs have distinctively different emergent relationships between environmental

300    factors and SOC stocks. These results could either result from unrealistic parameterization or

301    missing critical processes in model representation. Our results show that observed global SOC

302    stocks are controlled not only by temperature, precipitation, and NPP. Effects of other

303    environmental factors, such as drought severity index and cation exchange capacity should also

304    be considered in future representations of SOC dynamics in ESMs. It is also imperative to

305    compare observational data and ESM simulations to improve model structures and

306    parameterization.

307

308

**5. Conclusion**

310 Our results document disagreement between environmental controllers of SOC stocks in

311 observations and ESM land models. Specifically, NPP, annual temperature, and annual

312 precipitation have dominant control in modeled SOC stocks. In contrast, diurnal temperature,

313 drought index, annual temperature, cation exchange capacity, and other soil related variables are

314 the dominant controllers of observed SOC stocks. Using field observations and data for

315 environmental factors, machine learning techniques predict about 60% of the variability in

316 observed global SOC stocks, while in ESMs, only a few environmental factors predict about

317 95% of the variability in predicted SOC stocks. Comparisons of derived functional relationships

318 between the environmental factors and SOC stocks in observations and ESM models also show

319 discrepancies. These discrepancies indicate the importance of efforts to benchmark ESM land

320 models and to improve the mechanistic representations that are affected by the observed

321 dominant environmental controllers. Such an effort could decrease disagreements between

322 observed and modeled SOC stocks.

323

Biogeosciences

Discussions

Open Access

EGU

**References**

341    Anderegg, W. R. L., Berry, J. A., Smith, D. D., Sperry, J. S., Anderegg, L. D. L., and Field, C.
342    B.: The roles of hydraulic and carbon stress in a widespread climate-induced forest die-off, Proc
343    Natl Acad Sci U S A, 109, 233–237, https://doi.org/10.1073/PNAS.1107891109, 2012.

345    Arnold, D., Wagner, P., and Baayen, R. B.: Using generalized additive models and random
346    forests to model prosodic prominence in German, isca-speech.org, 2013.

348    Arora, V., Katavouta, A., Williams, R., Jones, C. D., Brovkin, V., Friedlingstein, P., Schwinger,
349    J., Bopp, L., Boucher, O., Cadule, P., Chamberlain, M., Christian, J., Delire, C., Fisher, R.,
350    Hajima, T., Ilyina, T., Joetzjer, E., Kawamiya, M., Koven, C., Krasting, J., Law, R., Lawrence,
351    D., Lenton, A., Lindsay, K., Pongratz, J., Raddatz, T., Séférian, R., Tachiiri, K., Tjiputra, J.,
352    Wiltshire, A., Wu, T., and Ziehn, T.: Carbon–concentration and carbon–climate feedbacks in
353    CMIP6 models and their comparison to CMIP5 models, Biogeosciences, 17, 4173–4222,
354    https://doi.org/10.5194/bg-17-4173-2020, 2020.

356    Azizi-Rad, M., Guggenberger, G., Ma, Y., and Sierra, C. A.: Sensitivity of soil respiration rate
357    with respect to temperature, moisture and oxygen under freezing and thawing, Soil Biology and
358    Biochemistry, 165, https://doi.org/10.1016/j.soilbio.2021.108488, 2022.

360    Batjes, N., Ribeiro E., and Oostrum, A.: Standardised soil profile data to support global mapping
361    and modelling (WoSIS snapshot 2019), Earth Syst. Sci. Data, 12, 299–320,
362    https://doi.org/10.5194/essd-12-299-2020, 2020.

364   Breiman, L.: Random forests, Mach Learn, 45, 5–32, https://doi.org/10.1023/A:1010933404324,
365   2001.
366
367   Chagas, C. da S., Junior, W de C., Bhering, S. B., and Filho, B. C.: Spatial prediction of soil
368   surface texture in a semiarid region using random forest and multiple linear regressions, Catena,
369   139, 232-240, https://doi.org/10.1016/j.catena.2016.01.001, 2016.
370
371   Chen, S., Wang, W., Xu, W., Wang, Y., Wan, H., Chen, D., Tang, Z., Tang, X., Zhou, G., Xie,
372   Z., Zhou, D., Shangguan, Z., Huang, J., He, J. S., Wang, Y., Sheng, J., Tang, L., Li, X., Dong,
373   M., Wu, Y., Wang, Q., Wang, Z., Wu, J., Stuart Chapin, F., and Bai, Y.: Plant diversity enhances
374   productivity and soil carbon storage, Proc Natl Acad Sci U S A, 115, 4027–4032,
375   https://doi.org/10.1073/PNAS.1700298114, 2018.
376
377   Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J.,
378   Mu, M., and Randerson, J. T.: The International Land Model Benchmarking (ILAMB) system:
379   design, theory, and implementation, Wiley Online Library, 10, 2731–2754,
380   https://doi.org/10.1029/2018MS001354, 2018.
381
382   Frank, S., Schmid, E., Havlík, P., Schneider, U.A., Böttcher, H., Balkovič, J. and Obersteiner,
383   M.: The dynamic soil organic carbon mitigation potential of European cropland, Global
384   Environmental Change, 35, 269-278, https://doi.org/10.1016/j.gloenvcha.2015.08.004, 2015.
385
386   Friedlingstein, P., Meinshausen, M., Arora, V.K., Jones, C.D., Anav, A., Liddicoat, S.K. and
387   Knutti, R.: Uncertainties in CMIP5 climate projections due to carbon cycle feedbacks, Journal of
388   Climate, 27, 511-526, https://doi.org/10.1175/JCLI-D-12-00579.1, 2014.
389
390   Georgiou, K., Malhotra, A., Wieder, W. R., Ennis, J. H., Hartman, M. D., Sulman, B. N., Berhe,
391   A. A., Grandy, A. S., Kyker-Snowman, E., Lajtha, K., Moore, J. A. M., Pierson, D., and Jackson,
392   R. B.: Divergent controls of soil organic carbon between observations and process-based models,
393   Biogeochemistry, 156, 5–17, https://doi.org/10.1007/S10533-021-00819-2, 2021.
394
395   Hastie, T. and Tibshirani, R.: Generalized additive models: Some applications, J Am Stat Assoc,
396   82, 371–386, https://doi.org/10.1080/01621459.1987.10478440, 1987.
397
398   Hinge, G., Surampalli, R. Y., and Goyal, M. K.: Prediction of soil organic carbon stock using
399   digital mapping approach in humid India, Environ Earth Sci, 77, https://doi.org/10.1007/S12665-
400   018-7374-X, 2018.
401
402   Hurrell, J.W., Holland, M.M., Gent, P.R., Ghan, S., Kay, J.E., Kushner, P.J., Lamarque, J.F.,
403   Large, W.G., Lawrence, D., Lindsay, K. and Lipscomb, W.H.: The community earth system
404   model: a framework for collaborative research, Bulletin of the American Meteorological Society,
405   94, 1339-1360, https://doi.org/10.1175/BAMS-D-12-00121.1, 2013.
406
407   Jiang, H., Deng, Q., Zhou, G., Hui, D., Zhang, D., Liu, S., Chu, G., and Li, J.: Responses of soil
408   respiration and its temperature/moisture sensitivity to precipitation in three subtropical forests in
409   southern China, Biogeosciences, 10, 3963–3982, https://doi.org/10.5194/bg-10-3963-2013, 2013.

Biogeosciences
Discussions

410
411   Lal, R.: Soil health and carbon management, Food Energy Secur, 5, 212–222,
412   https://doi.org/10.1002/fes3.96, 2016.
413
414   Lal, R.: Managing soils for negative feedback to climate change and positive impact on food and
415   nutritional security, Soil Sci Plant Nutr, 66, 1–9,
416   https://doi.org/10.1080/00380768.2020.1718548, 2020.
417
418   Lauer, A., Eyring, V., Righi, M., Buchwitz, M., Defourny, P., Evaldsson, M., Friedlingstein, P.,
419   de Jeu, R., de Leeuw, G., Loew, A. and Merchant, C.J.: Benchmarking CMIP5 models with a
420   subset of ESA CCI Phase 2 data using the ESMValTool, Remote Sensing of Environment, 203,
421   9-39, https://doi.org/10.1016/j.rse.2017.01.007, 2017.
422
423   Li, S., Liu, Y., Lyu, S., Wang, S., Pan, Y., and Qin, Y.: Change in soil organic carbon and its
424   climate drivers over the Tibetan Plateau in CMIP5 earth system models, Theor Appl Climatol,
425   145, 187–196, https://doi.org/10.1007/S00704-021-03631-Y, 2021.
426
427   Loveland, T. R. and Belward, A. S.: The igbp-dis global 1km land cover data set, discover: First
428   results, Int J Remote Sens, 18, 3289–3295, https://doi.org/10.1080/014311697217099, 1997.
429
430   Luo, Y. Q., Randerson, J., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., Ciais, P.,
431   Dalmonech, D., Fisher, J., Fisher, R., Friedlingstein, P., Hibbard, K., Hoffman, F., Huntzinger,
432   D., Jones, C. D., Koven, C., Lawrence, D., Li, D. J., Mahecha, M., Niu, S. L., Norby, R., Piao, S.
433   L., Qi, X., Peylin, P., Prentice, I. C., Riley, W., Reichstein, M., Schwalm, C., Wang, Y. P., Xia,
434   J. Y., Zaehle, S., and Zhou, X. H.: A framework for benchmarking land models,
435   bg.copernicus.org, 9, 1899–1944, https://doi.org/10.5194/bgd-9-1899-2012, 2012.
436
437   Luo, Z., Viscarra-Rossel, R.A. and Qian, T.: Similar importance of edaphic and climatic factors
438   for controlling soil organic carbon stocks of the world, Biogeosciences, 18, 2063-2073.,
439   https://doi.org/10.5194/bg-18-2063-2021, 2021.
440
441   McBratney, A.B., Santos, M.M. and Minasny, B.: On digital soil mapping, Geoderma, 117, 3-52,
442   https://doi.org/10.1016/S0016-7061(03)00223-4, 2003.
443
444   Mekonnen, Z., Riley, W., … J. R.-E., and 2022, undefined: Wildfire exacerbates high-latitude
445   soil carbon losses from climate warming, iopscience.iop.org, https://doi.org/10.1088/1748-
446   9326/ac8be6, 2022.
447
448   Mishra, U., Gautam, S., Riley, W. J., and Hoffman, F. M.: Ensemble Machine Learning
449   Approach Improves Predicted Spatial Variation of Surface Soil Organic Carbon Stocks in Data-
450   Limited Northern Circumpolar Region, Front Big Data, 3,
451   https://doi.org/10.3389/FDATA.2020.528441/FULL, 2020.
452
453   Mishra, U., Yeo, K., Adhikari, K., Riley, W. J., Hoffman, F. M., Hudson, C., and Gautam, S.:
454   Empirical relationships between environmental factors and soil organic carbon produce

455    comparable prediction accuracy to machine learning, Wiley Online Library, 86, 1611–1624,
456    https://doi.org/10.1002/saj2.20453, 2022.
457
458    O'Brien, S.L., Jastrow, J.D., Grimley, D.A. and Gonzalez-Meler, M.A.: Edaphic controls on soil
459    organic carbon stocks in restored grasslands, Geoderma, 251, 117-123,
460    https://doi.org/10.1016/j.geoderma.2015.03.023, 2015.
461
462    Sellar, A. A., Jones, C. G., Mulcahy, J. P., Tang, Y., Yool, A., Wiltshire, A., O'Connor, F. M.,
463    Stringer, M., Hill, R., Palmieri, J., Woodward, S., de Mora, L., Kuhlbrodt, T., Rumbold, S. T.,
464    Kelley, D. I., Ellis, R., Johnson, C. E., Walton, J., Abraham, N. L., Andrews, M. B., Andrews,
465    T., Archibald, A. T., Berthou, S., Burke, E., Blockley, E., Carslaw, K., Dalvi, M., Edwards, J.,
466    Folberth, G. A., Gedney, N., Griffiths, P. T., Harper, A. B., Hendry, M. A., Hewitt, A. J.,
467    Johnson, B., Jones, A., Jones, C. D., Keeble, J., Liddicoat, S., Morgenstern, O., Parker, R. J.,
468    Predoi, V., Robertson, E., Siahaan, A., Smith, R. S., Swaminathan, R., Woodhouse, M. T., Zeng,
469    G., and Zerroukat, M.: UKESM1: Description and Evaluation of the U.K. Earth System Model, J
470    Adv Model Earth Syst, 11, 4513–4558, https://doi.org/10.1029/2019MS001739, 2019.
471
472    Shi, X., Mao, J., Thornton, P.E. and Huang, M.: Spatiotemporal patterns of evapotranspiration in
473    response to multiple environmental factors simulated by the Community Land Model,
474    Environmental Research Letters, 8, 024012, https://doi.org/10.1088/1748-9326/8/2/024012,
475    2013.
476
477    Solly, E. F., Weber, V., Zimmermann, S., Walthert, L., Hagedorn, F., and Schmidt, M. W. I.: A
478    Critical Evaluation of the Relationship Between the Effective Cation Exchange Capacity and
479    Soil Organic Carbon Content in Swiss Forest Soils, Frontiers in Forests and Global Change, 3,
480    https://doi.org/10.3389/FFGC.2020.00098/FULL, 2020.
481
482    Sreenivas, K., Dadhwal, V.K., Kumar, S., Harsha, G.S., Mitran, T., Sujatha, G., Suresh, G.J.R.,
483    Fyzee, M.A. and Ravisankar, T.: Digital mapping of soil organic and inorganic carbon status in
484    India, Geoderma, 269, 160-173, https://doi.org/10.1016/j.geoderma.2016.02.002 , 2016.
485
486    Sun, Y., Piao, S., Huang, M., Ciais, P., Zeng, Z., Cheng, L., Li, X., Zhang, X., Mao, J., Peng, S.,
487    Poulter, B., Shi, X., Wang, X., Wang, Y. P., and Zeng, H.: Global patterns and climate drivers of
488    water-use efficiency in terrestrial ecosystems deduced from satellite-based datasets and carbon
489    cycle models, Global Ecology and Biogeography, 25, 311–323,
490    https://doi.org/10.1111/GEB.12411, 2016.
491
492    Todd-Brown, K. E. O., Randerson, J. T., Post, W. M., Hoffman, F. M., Tarnocai, C., Schuur, E.
493    A. G., and Allison, S. D.: Causes of variation in soil carbon simulations from CMIP5 Earth
494    system models and comparison with observations, Biogeosciences, 10, 1717–1736,
495    https://doi.org/10.5194/BG-10-1717-2013, 2013.
496
497    Wiesmeier, M., Barthold, F., Blank, B., and Kögel-Knabner, I.: Digital mapping of soil organic
498    matter stocks using Random Forest modeling in a semi-arid steppe ecosystem, Plant Soil, 340,
499    7–24, https://doi.org/10.1007/S11104-010-0425-Z, 2011.
500

501    Wiesmeier, M., Barthold, F., Spörlein, P., Geuß, U., Hangen, E., Reischl, A., Schilling, B.,
502    Angst, G., von Lützow, M. and Kögel-Knabner, I.: Estimation of total organic carbon storage
503    and its driving factors in soils of Bavaria (southeast Germany), Geoderma Regional, 1, 67-78,
504    https://doi.org/10.1016/j.geodrs.2014.09.001, 2014.
505
506    Wynn, J. G., Bird, M. I., Vellen, L., Grand-Clement, E., Carter, J., and Berry, S. L.: Continental-
507    scale measurement of the soil organic carbon pool with climatic, edaphic, and biotic controls,
508    Wiley Online Library, 20, https://doi.org/10.1029/2005GB002576, 2006.
509
510    Xiao-Ge, X.I.N., Tong-Wen, W.U., Jie ZHANG, F.Z., Wei-Ping, L.I., Yan-Wu ZHANG,
511    Y.X.L., Yong-Jie, F.A.N.G., Wei-Hua, J.I.E., Li ZHANG, M.D., Xue-Li, S.H.I., Jiang-Long, L.I.
512    and Min, C.H.U.: Introduction of BCC models and its participation in CMIP6, Advances in
513    Climate Change Research, 15, 533, https://doi.org/10.12006/j.issn.1673-1719.2019.039, 2019.
514
515    Yigini, Y., Olmedo, G., Reiter, S., Baritz, R., and Viatkin, K.: Soil organic carbon mapping:
516    cookbook, 2018.
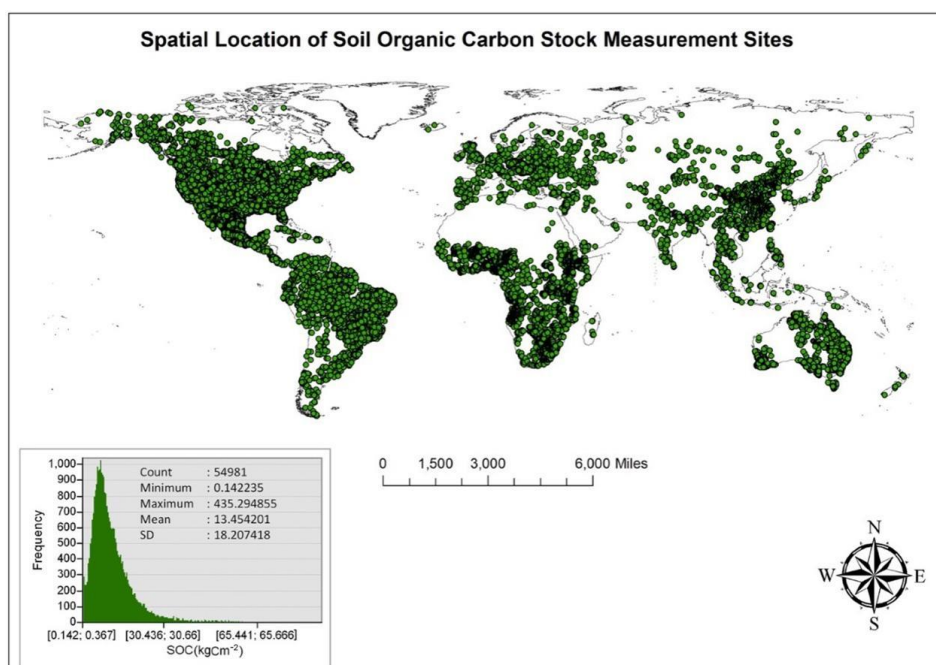517

**Figures and Tables**



Figure 1. Spatial and statistical distributions of 54,000 soil organic carbon profiles used in this study.
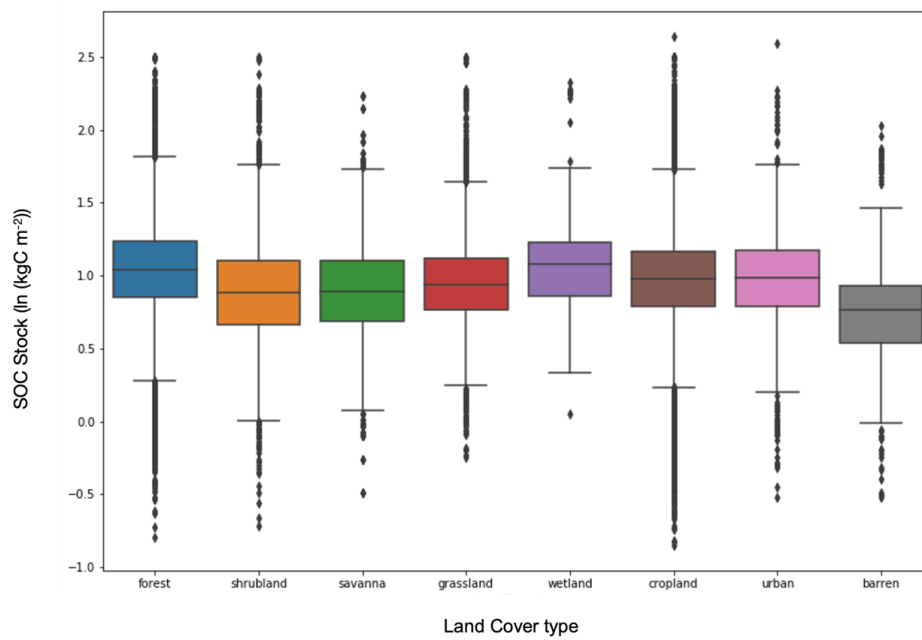
Figure 2: Boxplot of soil organic carbon content (logarithmic scale) for each biome or land cover type analyzed in this study. The horizontal line in the middle of the boxes is the median while their lower and upper limits correspond to the first and third quartiles.
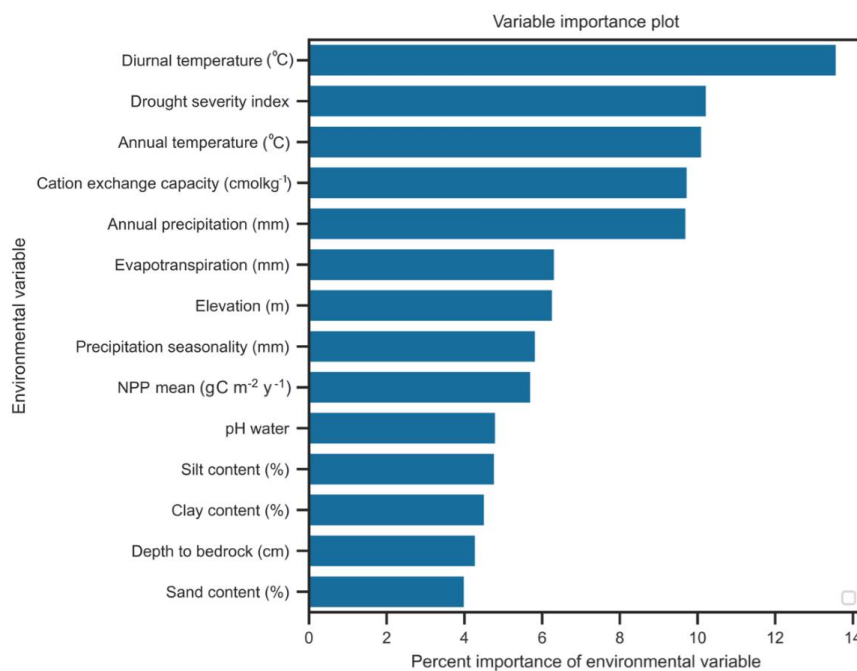
Figure 3: Importance of different environmental factors to predict the global soil organic carbon stocks in observations.
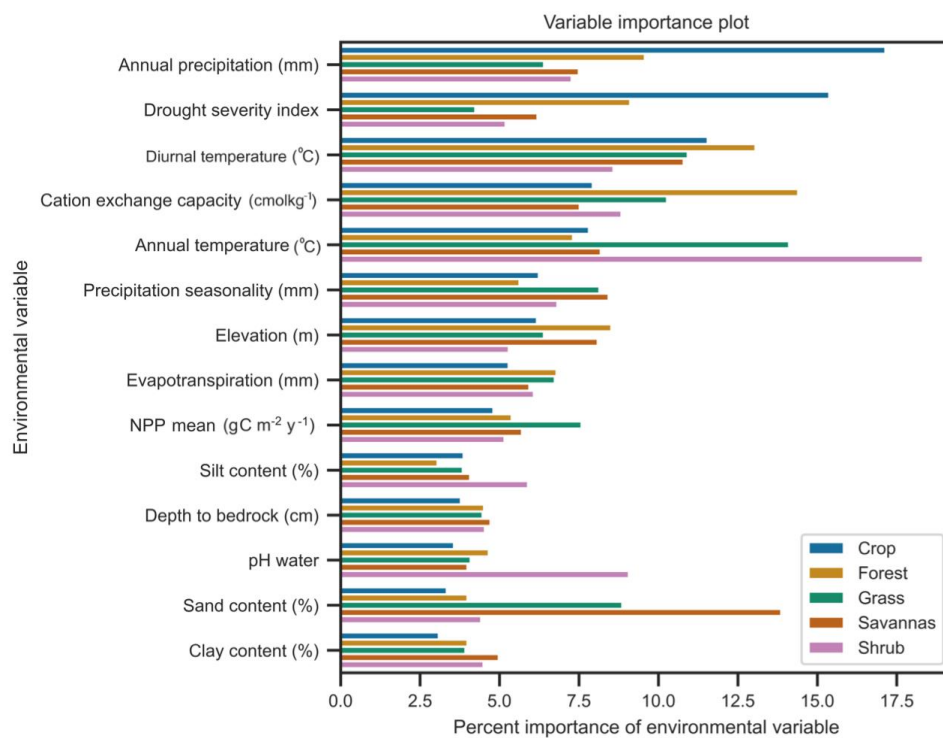
Figure 4: Strengths and importance of environmental controllers of observed SOC stocks within different biomes.
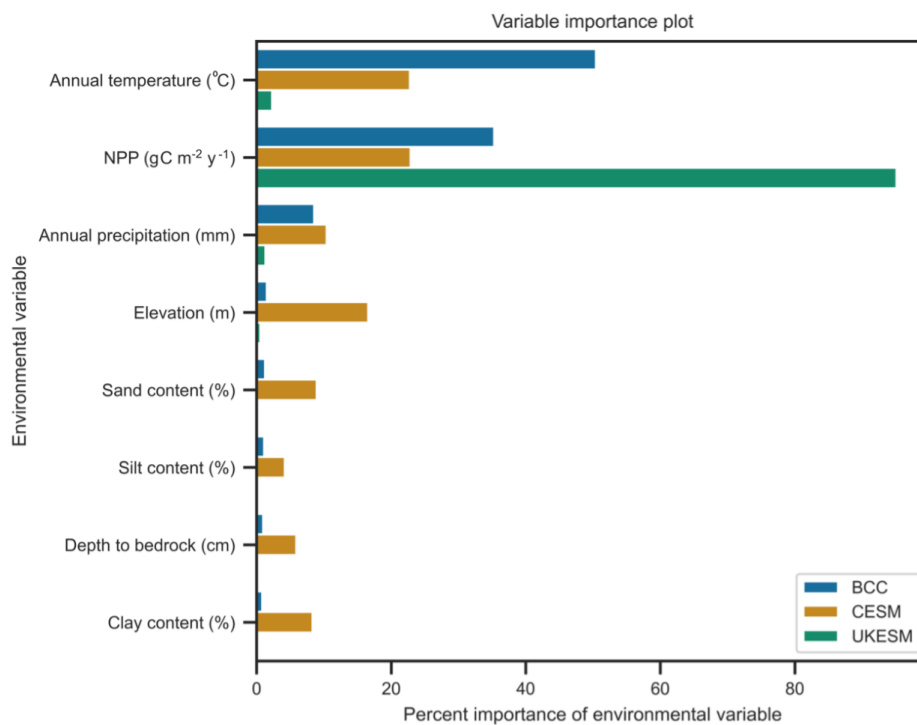
Figure 5: Importance of different environmental factors on global soil organic carbon stocks in three CMIP6 earth system models.
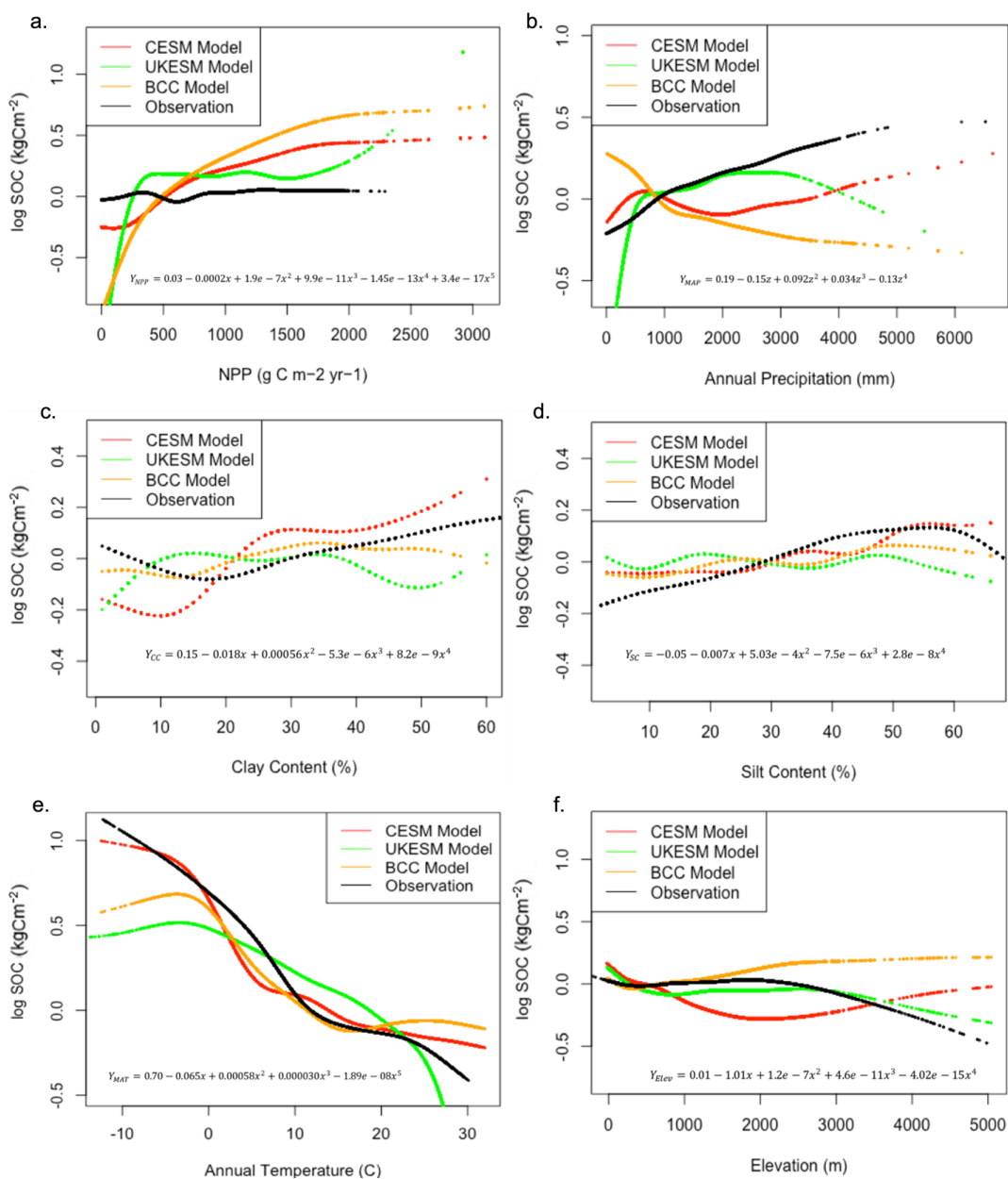
Figure 6: Predictive relationships between environmental factors and soil organic carbon stocks in observations (black line) and CMIP6 earth system models (different colors).

Table 1: Descriptive statistics of global soil organic carbon stocks at 0-100 cm depth interval.

| Location | Depth (cm) | Minimum (kgC m$^{-2}$) | Maximum (kgC m$^{-2}$) | Mean (kgC m$^{-2}$) | Median (kgC m$^{-2}$) | Standard Deviation (kgC m$^{-2}$) |
|---|---|---|---|---|---|---|
| Global | 0-100 | 0.14 | 435.3 | 13.5 | 9.5 | 18.2 |
| Cropland | 0-100 | 0.14 | 435.3 | 12.75 | 9.5 | 16.0 |
| Grassland | 0-100 | 0.56 | 315.9 | 12.1 | 8.7 | 16.8 |
| Forest | 0-100 | 0.16 | 314.4 | 15.9 | 10.9 | 20.7 |
| Shrubland | 0-100 | 0.19 | 312.5 | 13.6 | 7.6 | 25.6 |
| Savannas | 0-100 | 0.32 | 309.1 | 12.6 | 9.2 | 15.2 |

Table 2: Prediction accuracies of Random Forest models across biomes and at global scale in predicting SOC stocks.

| Location | Depth (cm) | R square (RF) | RMSE |
|---|---|---|---|
| Global | 0-100 | 0.61 | 0.46 |
| Cropland | 0-100 | 0.65 | 0.51 |
| Grassland | 0-100 | 0.57 | 0.46 |
| Forest | 0-100 | 0.59 | 0.52 |
| Shrubland | 0-100 | 0.64 | 0.54 |
| Savannas | 0-100 | 0.48 | 0.52 |