



1 **Global patterns and drivers of soil dissolved organic carbon concentrations**

2

3

Tianjing Ren and Andong Cai

4

5

6 Institute of Environment and Sustainable Development in Agriculture, Chinese Academy of Agricultural Sciences,

7 Beijing, 100081, China

8

9

10 ***Corresponding author:** caandong@caas.cn

11 Tel: +86-10-82105615, Fax: +86-10-82105615

12

13



14 **Abstract**

15 Dissolved organic carbon (DOC) is the most active carbon pool in soils, which plays critical roles in soil carbon
16 cycling, plant productivity, and global climate change. An accurate assessment of the quantity of DOC in the soil is
17 essential for the detailed elucidation of ecosystem functions and services. Nevertheless, the global driving factors
18 and distribution of soil DOC remain inadequately quantified due to the scarcity of large-scale data. Here, a
19 comprehensive global database of 12807 soil DOC concentrations derived from 975 target papers in the literature
20 was compiled. Detailed geographic locations, climate, and soil properties were also recorded as predictors of soil
21 DOC. Machine learning techniques were employed to assess the relative importance of various predictors in the
22 determination of soil DOC concentrations, which were subsequently extended for their prediction on a global scale.
23 The worldwide soil DOC concentration spanned a wide range (0.04 to 7859 mg kg⁻¹), averaging 222.78 mg kg⁻¹.
24 The 12 selected variables (including soil properties, month, climate, and ecosystem) explained 65% of the variance
25 in soil DOC concentrations. Elevation, soil clay, and soil organic carbon were three of the most important predictors.
26 Global soil DOC concentration increased from the equator to the poles. The soil DOC stocks in the topsoil layer (0-
27 30 cm) amounted to 12.17 Pg, with significant variations observed across different continents. These results are
28 instrumental for informing strategies on soil management practices, ecosystem services, and the mitigation of
29 climate change. Furthermore, our database can be combined with other carbon pools to explore the total soil carbon
30 turnover and constrain Earth carbon models. The dataset is publicly available at
31 <https://doi.org/10.6084/m9.figshare.26379898> (Ren and Cai, 2024).

32



33 1. Introduction

34 With global changes over the last few decades, terrestrial ecosystems, which serve as the fundamental safeguard for
35 biodiversity and carbon sink on Earth, are becoming increasingly vital toward mitigating global climate warming
36 (IPCC, 2014). Cumulatively, soil carbon pools constitute the largest carbon reservoirs of terrestrial ecosystems,
37 which are three to four times greater than that of the ambient atmospheric carbon pool (Lal, 2004). Even minor
38 fluctuations in soil carbon can have significant impacts on biogeochemical cycles and the global C balance.
39 Dissolved organic carbon (DOC), which consists of simple organic acids and complex macromolecular substances,
40 is recognized as the most active carbon pool in the soil. Currently, the portion of organic carbon that is water-soluble
41 and can filter through 0.45 μm microporous filter membrane is referred to as DOC (Kalbitz et al., 2000; Zsolnay,
42 2003). Although soil DOC typically accounts for only $< 2\%$ of soil carbon pool, it provides a substantial source of
43 carbon and energy for soil microorganisms, while playing a key role in soil carbon sequestration, transport, and
44 stabilization mechanisms (Nakhavali et al., 2021; Ren et al., 2024). The lateral transport of DOC is crucial for
45 linking terrestrial and aquatic ecosystems and plays a key role in the evaluation of terrestrial carbon budgets
46 (Kindler et al., 2011; Sanderman & Amundson, 2008). Thus, an accurate assessment of soil DOC concentrations is
47 vital due to its unique properties and roles, given its broad variations that can span up to three orders of magnitude
48 (Nakhavali et al., 2020; Ren et al., 2024). Despite the significant variations in soil DOC concentrations, their global
49 distribution has not yet been systematically quantified. Bridging this knowledge gap is essential for more accurately
50 depicting the carbon cycle in Earth system models.

51 The soil DOC concentration depends on the dynamic balance between its sources (e.g., leachates from
52 decomposing plant litter, plant root secretions, and microbial decomposition products) and losses (migration and
53 microbial decomposition) (Bolan et al., 2011). Therefore, any factors that affect this dynamic balance would also
54 influence the soil DOC concentrations. Extensive research has demonstrated that the soil DOC concentration is the
55 outcome of climate, vegetation type, as well as soil properties (Chen et al., 2021; Guo et al., 2020; Smreczak &
56 Ukalska-Jaruga, 2021). Each factor plays a distinct role in shaping soil DOC dynamics. For example, the climate,
57 which is characterized by the annual mean temperature and precipitation, is typically recognized as a primary
58 driving factor that influences the soil DOC concentrations (Kalbitz et al., 2000; Neff & Asner, 2001). Temperature
59 and precipitation directly influence soil DOC concentrations by affecting microbial activities, organic matter



60 decomposition rates, its solubility and mobility, and indirectly modulate DOC dynamics by manipulating vegetation
61 growth and soil structures (Andersson & Nilsson, 2001; Kalbitz et al., 2000). The type of vegetation impacts soil
62 DOC concentrations mainly by affecting the input quantity and quality of organic matter (Guo et al., 2020). Together,
63 climate and vegetation types have profound effects on soil biological, chemical, and physical properties, which are
64 closely interconnected with the creation and decomposition of soil DOC (Camino - Serrano et al., 2014). The
65 relationships between soil DOC concentrations and environmental factors have been revealed based on local and
66 regional scales. However, the relative importance of environmental factors that predict soil DOC concentrations on a
67 global scale is still lacking, which impedes the development of effective strategies for the management of soil
68 carbon and mitigation of climate change.

69 Accurate mapping of the soil DOC is essential for addressing pressing global challenges, including climate
70 warming, food security, and eutrophication in aquatic systems (Guo et al., 2020; Langeveld et al., 2020). To the best
71 of our knowledge, there are few global maps of the spatial distribution of soil DOC (Guo et al., 2020; Langeveld et
72 al., 2020). However, these maps have subject to considerable uncertainties due to the limited data employed and the
73 low interpretation rate. Firstly, there is a lack of valid observational data for Africa, South America, Eastern Europe,
74 and Central Asia. Secondly, Guo (Guo et al., 2020) explained only 31% of the variations in the soil DOC using
75 linear regression equations, while Langeveld (Langeveld et al., 2020) explained only 36%. In contrast to linear
76 regression, machine learning has been extensively applied in research due to its capacities to automate feature
77 extraction, handle large datasets, and recognize complex patterns, which offers significant advantages in terms of
78 predictive accuracy and adaptive learning.

79 To address these challenges, we developed a comprehensive database of global soil DOC concentrations,
80 comprising 12,807 samples from 975 published studies. Utilizing Random Forest algorithms, we quantified the
81 relative importance of environmental factors, and further, predicted the soil DOC concentrations on a global scale.
82 The special aims of this study were: (1) What are the global patterns of soil DOC concentrations? (2) What are the
83 primary factors that control soil DOC concentrations on a global scale? (3) How large is total global soil DOC
84 storage?

85 **2. Material and method**

86 **2.1 Data sources and processing**



87 Publication search for this study was performed using Google Scholar (<https://scholar.google.com>), the Web of
88 Science (<http://apps.webofknowledge.com>), and the China Knowledge Resource Integrated Database
89 (<http://www.cnki.net/>) using the following search terms: (dissolved organic carbon OR dissolved organic matter OR
90 "DOC" OR "DOM") AND soil, up to December 2022. The specific data flow through the different phases for the
91 selected papers is shown in Fig. S1. To ensure a standardized and bias-minimized dataset, the following inclusion
92 criteria were applied: (1) Data must be from terrestrial ecosystems, excluding oceans and rivers; (2) Only the topsoil
93 layer data (0-30 cm) were used; (3) Duplicate results from different articles were recorded only once; (4) Soils
94 included agricultural soils that were affected by human activities through tilling and fertilization etc., but did not
95 cover industrial and urban soils. Data presented solely in figures were extracted using the digitizer function of
96 Origin 2019 software.

97 Based on these criteria, a total of 12807 observations of soil DOC were compiled from 975 publications.
98 Additional data included specifics of the experimental sites (longitude, latitude, and altitude), climatic conditions
99 (mean annual temperature (MAT) and mean annual precipitation (MAP)), biomes (e.g., wetland, forest, shrubland,
100 tundra, grassland, and cropland) and soil physical and chemical properties (e.g., soil organic carbon, texture, and
101 pH) (Table 1). These environmental factors are used as predictors. When those environmental factors were missing
102 within the original publication, the missing data were extracted from the grid dataset according to geographic
103 coordinates of observed site (Table S1). This study sites spanned a wide range of latitudes (-64.81° to 78.85°) and
104 longitudes (-159.66° to 175.95°) (Table 1). This database encompassed a large gradient of climate regimes, with
105 MAT from -11.16 to 28.00°C and MAP from 30 to 4200 mm.

106 2.2 Data standardization

107 In our database, the DOC concentrations were quantified using a mix of physical and chemical techniques. Physical
108 methods included soil solution collection using lysimeters or ceramic suction. Chemical methods employed various
109 solvents like distilled water, potassium chloride (KCl), or potassium sulfate (K_2SO_4) as described by Li et al. (2018).
110 Over 74.32% of the DOC was determined using chemical techniques, which highlighted their reliability. For
111 consistency, the DOC values derived from physical approaches was converted to chemical method values using the
112 following equation:

$$113 \text{ DOC}_{\text{soil}} = (\text{DOC}_{\text{solution}} \times V \times 1000) / W \times [1 / (V \times (1 - W) \times \text{BD} \times 1000000)] \quad (1)$$



114 where, DOC_{soil} represents soil DOC concentration determined by chemical methods (mg g^{-1}); $\text{DOC}_{\text{solution}}$ is the
115 concentration measured by physical methods (mg L^{-1}); W denotes the volumetric soil moisture ($\text{m}^3 \text{m}^{-3}$); V is the
116 volume of the soil column for solution extraction (m^3); and BD is the soil bulk density (g cm^{-3}). The factor 1000
117 converts m^3 to L, and 1000000 converts m^3 to cm^3 following the protocol established by Guo (Guo et al., 2020). This
118 standardization allowed for a consistent comparison and analysis of the DOC data across various studies.

119 **2.3 Predictive modeling**

120 The driving factors of soil DOC concentrations were divided into four categories (climate, ecosystem, soil properties,
121 and observation time). The soil properties included physical (clay, sand, bulk density, and depth), chemical (SOC,
122 pH), biological (microbial biomass carbon) attributes. The observation time was represented by month. Climate
123 referred to MAT, MAP, and elevation. Ecosystems encompassed wetland, forest, shrubland, tundra, grassland, and
124 cropland. In predictive models, correlated predictors may substitute for each other, such that their importance will be
125 shared, which results in an estimated importance that is less than the true value. Consequently, the soil total nitrogen,
126 silt, and aridity index were not included as they were correlated with the soil organic carbon, sand, and MAP,
127 respectively (Fig. S2). Further, some variables were not included due to rarely report in target paper.

128 To develop and optimize a predictive model for soil DOC an array of regression methods was employed, which
129 encompassed three linear and four nonlinear approaches (Table S2). The linear regression methods included a least
130 absolute shrinkage and selection operator (LEAPS), elastic net (ENET), and standard linear modeling (LM) to
131 identify the most important predictor variable in a regression model, while minimizing the risk of overfitting. The
132 nonlinear regression methods included the random forest (RF) algorithm, boosted tree (BOOSTED), bagged tree
133 (Bagged), and cubist (CUBIST) models. Each model was equipped with intrinsic feature selection processes and
134 was fine-tuned to improve accuracy and control complexity. During the optimization phase, various actions were
135 implemented; LEAPS models were educated to accommodate the highest count of variables. To discipline the
136 models, the penalty for feature condensation (diminishing the role of less impactful variables in the resultant linear
137 formula) varied between 0 and 0.1, incremented by 0.01. RF models' growth was capped at a maximum of 1,000
138 trees, and the model's predictors were restricted to a third of the possible maximum, ensuring a balance between
139 complexity and manageability. BOOSTED models underwent training with a tree count ranging from ten to a
140 hundred, where each tree had a node range of one to seven. They incorporated a shrinkage rate of either 0.01 or 0.1,
141 and a maximum size limit set to five, optimizing the models' learning process. CUBIST model utilized a sequence of



142 neighboring values from 1 to 9 with increments of 2, alongside community sizes spanning 1 to 100, to refine its
143 predictive accuracy. In every instance, the models were evaluated using Monte Carlo cross-validation, which
144 involved 100 iterations of data resampling with an 80/20 split between training and validation datasets, ensuring an
145 accurate estimation of model uncertainty and safeguarding against over-fitting. The root mean square error and R^2
146 values were calculated to evaluate model accuracy and residual variance, which served as criteria for ranking model
147 performance (Table S2). The relative RMSE, a measure of the estimation uncertainty for soil DOC, was determined
148 by dividing the error's magnitude by the overall average soil DOC value. The nonlinear models ($R^2 = 0.42-0.65$; root
149 mean square error (RMSE) = 250-332) outperformed the linear models ($R^2 = 0.101-0.108$; RMSE = 410-427) (Table
150 S2). The RF model distinguished itself with the lowest RMSE within a standard deviation range, and the model was
151 then selected for subsequent analyses focusing on variable importance (Fig. S3). Consequently, the relative
152 importance of driving the soil DOC and the global map of soil DOC were the averaged values of the RF model
153 results.

154 To evaluate the impacts of independent variables on the soil DOC, a variable importance analysis was conducted
155 using permutation variable importance measurements (Fig. 2). This analysis was performed utilizing the variable
156 importance tool integrated into the R packages for the RF model that exhibited the highest accuracy and predictive
157 quality. In essence, this method assessed prediction errors within the model by calculating mean square errors for
158 each regression tree. The models' variable importance scores assessed the influence of predictor variables on the
159 outcomes. For enhanced comparability of all model inputs, the independent environmental variables were scaled to a
160 0 to 100% range, reflecting their proportional contribution to the model's predictions.

161 Partial dependence analyses were employed to test the relationships between the predicted soil DOC and
162 independent variables across the entire spectrum of potential values considered in the RF model (Fig. 3). In essence,
163 this approach provided insights into the global relationships between the independent variables and predicted
164 outcomes. The focus was set solely on the effects of the targeted independent variables by eliminating the influences
165 of other independent variables. Partial dependence analyses, along with their graphical representations known as
166 partial dependence plots, provided insight into the average marginal effect of one or more independent variables on a
167 machine learning model's predictions within a defined value scope, offering a more nuanced view than assessing the
168 overall relative importance of an independent variable. For instance, partial dependence plots can expose whether
169 the connection between a predicted variable and an independent control is linear, monotonic, or complex. The



170 curvature and inflection points of the partial dependence plot curve help us to decipher and pinpoint areas where an
171 independent variable exerts a notably strong and immediate effect on the forecasted outcome. Additionally, it can
172 indicate where the variable's influence is more subtle, potentially mediated through its effects on other independent
173 variables. To facilitate the interpretation of the partial dependence plots, the x-axis for the standardized value was
174 reported, which ensured a clear progression from low to high values in all curves.

175 **2.4 Global soil DOC mapping**

176 The global distribution of the soil DOC and the relative uncertainties of predictions were generated (Figs. 4, S5).
177 These maps were derived by utilizing our DOC dataset in conjunction with the RF model, which incorporated the
178 global climate, vegetation, and soil-rasterized datasets (Table S1). We generated factor maps from the key input
179 variables, focusing on the 12 distinct variables associated with each raster cell. Subsequently, the factor maps were
180 employed to derive a spatially detailed global map of soil DOC. For global scale mapping, the driving factors were
181 initially processed at a 0.05° resolution to calculate the soil DOC values. Areas that did not meet the following
182 criteria were excluded from our prediction: (1) absence of data for any essential predictors, (2) soil order and biomes
183 not aligning with the previously discussed aggregated land use systems, or (3) locations in climate zones outside the
184 scope of our model's focus. To evaluate the uncertainty associated with map creation due to data resampling and any
185 unexplained variability unaccounted for by the independent variables, we analyzed finer resolution (5 km^2) grids in
186 regions where driving factors were accessible at this detailed level. This analysis illuminated the overall uncertainty
187 inherent in our global soil DOC estimation. A map representing the relative prediction uncertainty was crafted,
188 showcasing the standard deviation in relation to the mean of the predictions. The standard deviation, indicative of
189 the dispersion in potential predictions, was derived from the decision tree model's structure after 500 iterations of the
190 model.

191

192 **3. Results**

193 **3.1 Soil DOC concentrations in different ecosystems globally**

194 A total of 12,807 soil DOC observations were compiled from 975 publications, which spanned six continents as well
195 as major biomes and terrestrial ecosystems (Fig. 1), and the database conformed to a normal distribution (Fig. 1b).
196 The global soil DOC concentrations varied between 0.04 and 7859 mg kg^{-1} . The global average, median, and
197 standard deviation were 222.78, 101.01, and $445.78 \text{ mg kg}^{-1}$, respectively (Table 2). The concentrations of soil DOC



198 varied across different ecosystems. Tundra had the highest soil DOC concentration ($470.78 \text{ mg kg}^{-1}$), while
199 shrubland had the lowest ($160.24 \text{ mg kg}^{-1}$). The average soil DOC concentrations for grassland, forest, wetland, and
200 cropland were 327.77 , 256.18 , 218.53 , and $165.98 \text{ mg kg}^{-1}$, respectively (Table 2).

201

202 **3.2 Model performance and drivers of soil DOC concentrations**

203 Random forest model accounted for 65% of the variability in soil DOC concentrations across all sites, with the
204 lowest RMSE compared with other models (Fig. 2, Table S2). The most important categories of predictors for soil
205 DOC concentrations were climate and soil properties, with elevation and the soil clay content emerging as the most
206 significant. Although less influential, other predictors were nonetheless considered, with soil organic carbon and soil
207 pH having the most notable effects (Fig. 2a). Although the mean annual precipitation and temperature, microbial
208 biomass carbon, bulk density, sand, depth, month, and ecosystem affected soil DOC concentrations, their relative
209 contributions were lower than aforementioned four predictors (Fig. 2). Partial dependence analysis showed similar
210 results to Pearson correlation analysis (Fig. S4) and indicated that there was a positive correlation between the soil
211 DOC and both the elevation and soil organic carbon (Fig. 3g). Conversely, the soil DOC was negatively correlated
212 with mean annual temperature and soil pH (Fig. 3h).

213

214 **3.3 Global soil DOC patterns**

215 Our predicted global soil DOC mapping implied that there was a significant spatial heterogeneity of soil DOC
216 concentrations (Fig. 4a). This revealed a latitudinal pattern that soil DOC concentrations increased from the equator
217 to poles (Fig. 4b). High soil DOC concentrations were found in high-altitude plateaus and mountain ranges at low
218 latitude (e.g., Andes, African Highlands, West Indies) (Fig. 4a). The global average soil DOC concentration was
219 $237.56 \text{ mg kg}^{-1}$ (Table 3), while the soil DOC stock in the topsoil (0-30 cm) was 12.17 Pg .

220 Asia had the highest soil DOC concentration ($274.43 \text{ mg kg}^{-1}$) followed by North America ($263.63 \text{ mg kg}^{-1}$). Next
221 were Europe and South America (227.34 and $215.81 \text{ mg kg}^{-1}$, respectively), with Oceania and Africa having the
222 lowest soil DOC concentrations (198.13 and $186.35 \text{ mg kg}^{-1}$, respectively). For predicted soil DOC stocks, Asia and
223 North America remained in first and second place (4.8 and 2.45 Pg , respectively). Despite the marginal predicted
224 soil DOC concentrations in Africa, its predicted soil DOC stocks ranked third (2.07 Pg) due to its vast area. South



225 America was in fourth place with a predicted soil DOC stock of 1.37 Pg. Finally, Europe and Oceania showed the
226 lowest predicted soil DOC stocks (0.88 and 0.59 Pg, respectively).

227

228 **4 Discussions**

229 **4.1 Variations in soil DOC between ecosystems**

230 Given the substantial number of measurements included in our study (12,807 observations), the range of soil DOC
231 concentrations (0.04-7859 mg kg⁻¹) was broader than that reported by Guo (3,869 observations) (Guo et al., 2020).

232 Our reported global average soil DOC concentration was 222.78 mg kg⁻¹ (Table 2), in contrast to Guo's reported
233 average of only 77.39 mg kg⁻¹. For different ecosystems, the soil DOC concentrations of wetlands, tundra, and

234 shrublands in our study aligned with those of previous research (Guo et al., 2020), which was primarily due to the

235 relatively lower number of observations for these ecosystems in comparison with others, with tundra comprising only

236 1% of our database (Guo et al., 2020). However, significant differences were found in forests, grasslands, and

237 croplands compared with Guo's data. For instance, our average soil DOC concentration for croplands was 165.98

238 mg kg⁻¹, while Guo reported only 60.58 mg kg⁻¹. This discrepancy was due to Guo's database including only 13%

239 cropland observations, whereas our cropland observations are approximately ten times larger (Guo et al., 2020).

240 However, our results consistently indicated that DOC concentrations in forest soils were lower than in grasslands,

241 with tundra showing the highest DOC levels (Table 2) (Guo et al., 2020). This was due to the higher lignin content

242 in forests, which reduces the quality of plant litter, hinders microbial decomposition, and releases less DOC (Wang

243 et al., 2015). For tundra, besides low microbial activities in permafrost due to low temperatures, anaerobic

244 conditions from soil oversaturation severely limit microbial activities and growth, reduce decomposition rates, and

245 increase the DOC (Boddy et al., 2008; Petrone, 2005). Despite the frequent addition of nutrients in croplands, the

246 DOC concentrations remained lower than expected. Intensive anthropogenic activities, such as management

247 practices and frequent harvesting induced the significant loss of soil organic matter, which translated to reduced

248 DOC (Guo et al., 2020; Li et al., 2019; Ren et al., 2024). In summary, our study built on preceding work by

249 incorporating a more extensive dataset that better represented the heterogeneous conditions found globally.

250

251 **4.2 Effects of climate and controlled soil properties on soil DOC concentrations**



252 The two most critical predictors of soil DOC concentrations were climate and soil properties, with elevation and soil
253 clay content being the two most significant factors (Fig. 3). As the elevation gradient increase, temperatures
254 generally decrease, which can constrain microbial metabolic rates and reduce the decomposition of organic matter,
255 which leads to additional organic carbon being retained in the soil as DOC (Li et al., 2023; Nottingham et al., 2019;
256 Wei et al., 2024). Typically, high-altitude regions host specific vegetation types with longer growth cycles and more
257 litterfall (Pesántez et al., 2018; Wei et al., 2024). These plant residues decompose to SOC, a portion of which
258 converts to DOC. Consequently, differences in the vegetation type and productivity also influence the soil DOC
259 concentrations (Camino - Serrano et al., 2014; Rahbek et al., 2019). We also found that forest and grassland sites
260 above 2000 m (which constituted 73% of the high DOC observations) were significant contributors. High-altitude
261 regions often experience distinct precipitation patterns and soil moisture conditions compared with lower elevations
262 (Li et al., 2023). Higher precipitation and lower evaporation rates may result in the greater dissolution and leaching
263 of organic matter, thereby increasing DOC concentrations in the soil (He et al., 2021; Lu et al., 2019). High-altitude
264 areas are generally less frequented by humans, which may assist in the preservation of the DOC in the soil through
265 the prevention of disturbances and losses. Our results also indicated that soils in low-latitude plateaus and mountain
266 ranges (e.g., Tibetan Plateau, Andes, African Highlands, and West Indies) exhibited higher DOC concentrations (Fig.
267 4a). The impacts of the soil clay content on DOC concentrations are complex, which occurred primarily through
268 adsorption, water retention, microbial activities, and organic matter protection mechanisms (Kaiser & Zech, 2000;
269 Singh et al., 2017). Generally, a high clay content tends to stimulate the accumulation of soil DOC through the
270 adsorption and stabilization of organic matter (Gmach et al., 2019; Kalbitz et al., 2000). Furthermore, the effects of
271 SOC and soil pH on DOC should not be overlooked (Fig. 2a). SOC serves as the main source of DOC, where higher
272 SOC generally implies that more DOC can be released into the soil through microbial metabolism (Kalbitz et al.,
273 2000; Neff & Asner, 2001). Variations in the soil pH can affect the charge of soil colloids, thereby altering their
274 adsorption-desorption mechanisms for DOC, which affects its solubility in the soil (Andersson & Nilsson, 2001;
275 Cheng et al., 2020; Kaiser et al., 2005). In summary, the soil DOC concentration is the result of interactions between
276 the soil and climate, biological, chemical, physical processes, and human influences at various spatial and temporal
277 scales, with each factor playing a unique role in shaping DOC dynamics.

278

279 **4.3 Global patterns of soil DOC**



280 Using our soil DOC concentration dataset, we quantified the soil DOC concentrations (0-30 cm) in terrestrial
281 ecosystems, identified their key driving factors, and made global predictions. Global DOC stocks in the topsoil are
282 estimated at 12.17 Pg C, accounting for 0.775% of the global soil organic carbon, which is significantly higher than
283 previous estimates (Guo et al., 2020). Our predictions indicated that soil DOC concentrations decreased significantly
284 with lower latitudes, particularly in the Northern Hemisphere. Previous global maps of soil DOC concentrations
285 failed to capture this latitudinal trend, which was likely due to their limited spatial coverage (Guo et al., 2020;
286 Langeveld et al., 2020). Our predicted map shows that the soil DOC concentrations increased with latitude. This
287 trend was attributed to lower temperatures, specific vegetation types, higher soil moisture, and reduced human
288 activities at higher latitudes (Camino - Serrano et al., 2014; Lapierre et al., 2015). However, there was substantial
289 heterogeneity at regional and local scales. For instance, despite being at similar latitudes, soil DOC concentrations in
290 Northern Europe were significantly lower than in Siberia, which we surmised was primarily due to differences
291 between the maritime climate of Northern Europe and the cold subarctic climate of Siberia. Regional variations in
292 soil DOC concentrations might be related to topographic condition. Higher soil DOC concentrations on the Tibetan
293 Plateau compared to Eastern China might result from the high elevation and low MAT in the plateau (Fig. 4a). In
294 contrast, lower DOC levels in Arctic regions was reported, which might have been due to their omission of DOC
295 concentration in the soil and dry or frozen soil (Langeveld et al., 2020). The predictive model offered higher
296 accuracy in estimating the global soil DOC storage (Fig. 3). This advantage stemmed from our comprehensive
297 dataset, which included DOC concentrations in both dry soil and soil solutions, which provided a robust data
298 foundation for global soil DOC predictions. Additionally, we employed the optimal model for predicting the global
299 soil DOC by comparing various linear and non-linear models.

300

301 **4.4 Limitations and predictive uncertainties**

302 Although we compiled a comprehensive global soil DOC concentration dataset, identified key drivers, and made a
303 global prediction, our study had certain limitations. First, certain ecosystems remained underrepresented; for
304 instance, tundra accounted for only 1% of our database, while shrublands, grasslands, and wetlands collectively
305 constituted only 21%. This underrepresentation may reduce the accuracy of predictions for different ecosystems.
306 Second, although we considered the subsoil at the beginning of dataset, we did not explore this further due to the
307 limited availability of data and considerations of predictive accuracy. We intend to continue expanding the subsoil



308 DOC database in future work. Third, there was a deficiency in some predictive variables; although we had extracted
309 missing data through gridded datasets, this inevitably introduced uncertainty in predictions, particularly for soil
310 variables. Fourth, despite employing advanced machine learning methods with multiple predictors to predict the
311 global soil DOC, 35% of soil DOC concentration variability remains unexplained. However, these limitations also
312 highlighted areas for future soil DOC research.

313

314 **5 Data availability**

315 The global soil DOC in this study and raw dataset of driving factors can be downloaded at
316 <https://doi.org/10.6084/m9.figshare.26379898> (Ren and Cai, 2024).

317

318 **6 Conclusions**

319 Through the development of a comprehensive soil DOC dataset, we quantified soil DOC concentrations in terrestrial
320 ecosystems, identified their driving factors, and made global predictions. Subsequent to comparing multiple
321 predictive models, we selected the Random Forest model as the best performer for mapping soil DOC
322 concentrations. The results indicated that tundra exhibited the highest DOC concentrations, while shrubland and
323 cropland soils had relatively lower concentrations. Climate factors (elevation) and soil properties (clay content, SOC,
324 pH) jointly regulated the DOC variations. The predicted that the soil DOC concentration increased significantly
325 from the equator to the poles, and estimated the DOC stocks in the topsoil of terrestrial ecosystems was 12.17 Pg.
326 The global soil DOC database we created will serve as a critical resource for future research, while enhancing our
327 understanding of the roles of soil in the global carbon cycle. This database provides valuable data support for
328 climate change research, ecosystem management, agricultural sustainability, environmental policymaking, and the
329 improvement of biogeochemical models. This will aid in addressing soil degradation, improving food security, and
330 tackling global environmental challenges.

331

332 **Author contributions**

333 Andong Cai designed this study. Tianjing Ren collected the data. Tianjing Ren and Andong Cai discussed analyzing
334 methods. Andong Cai conducted the analysis. Tianjing Ren drafted the manuscript. All authors discussed the results
335 and contributed to the manuscript.



336 **Competing interests**

337 The contact author has declared that neither they have any competing interests.

338

339 **Acknowledgements**

340 We would like to thank Frank Boehm at NanoApps Consulting2341York Ave. Vancouver, BC, Canada for his
341 assistance with English language and grammatical editing.

342

343 **Financial support**

344 This work was financially supported by the National Key Research and Development Program of China
345 (2022YFD2300500).

346

347 **References**

348 Andersson, S., & Nilsson, S. I. (2001). Influence of pH and temperature on microbial activity, substrate availability
349 of soil-solution bacteria and leaching of dissolved organic carbon in a mor humus. *Soil Biology and*
350 *Biochemistry*, 33(9), 1181-1191. [https://doi.org/10.1016/S0038-0717\(01\)00022-0](https://doi.org/10.1016/S0038-0717(01)00022-0)

351 Boddy, E., Roberts, P., Hill, P. W., Farrar, J., & Jones, D. L. (2008). Turnover of low molecular weight dissolved
352 organic C (DOC) and microbial C exhibit different temperature sensitivities in Arctic tundra soils. *Soil Biology*
353 *and Biochemistry*, 40(7), 1557-1566. <https://doi.org/10.1016/j.soilbio.2008.01.030>

354 Bolan, N. S., Adriano, D. C., Kunhikrishnan, A., James, T., McDowell, R., & Senesi, N. (2011). Dissolved organic
355 matter: biogeochemistry, dynamics, and environmental significance in soils. *Advances in agronomy*, 110, 1-75.
356 <https://doi.org/10.1016/B978-0-12-385531-2.00001-3>

357 Camino-Serrano, M., Gielen, B., Luyssaert, S., Ciais, P., Vicca, S., Guenet, B., Vos, B. D., Cools, N., Ahrens, B.,
358 Altaf Arain, M., Borken, W., Clarke, N., Clarkson, B., Cummins, T., Don, A., Pannatier, E. G., Laudon, H.,
359 Moore, T., Nieminen, T. M., . . . Janssens, I. (2014). Linking variability in soil solution dissolved organic carbon
360 to climate, soil type, and vegetation type. *Global Biogeochemical Cycles*, 28(5), 497-509.
361 <https://doi.org/10.1002/2013gb004726>



- 362 Chen, H., Kong, W., Shi, Q., Wang, F., He, C., Wu, J., Lin, Q., Zhang, X., Zhu, Y. G., Liang, C., & Luo, Y. (2021).
363 Patterns and drivers of the degradability of dissolved organic matter in dryland soils on the Tibetan Plateau.
364 *Journal of Applied Ecology*. <https://doi.org/10.1111/1365-2664.14105>
- 365 Cheng, X., Hou, H., Li, R., Zheng, C., & Liu, H. (2020). Adsorption behavior of tetracycline on the soil and
366 molecular insight into the effect of dissolved organic matter on the adsorption. *Journal of Soils and Sediments*,
367 20(4), 1846-1857. <https://doi.org/10.1007/s11368-019-02553-7>
- 368 Gmach, M. R., Cherubin, M. R., Kaiser, K., & Cerri, C. E. P. (2019). Processes that influence dissolved organic
369 matter in the soil: a review. *Scientia Agricola*, 77. <https://doi.org/10.1590/1678-992X-2018-0164>
- 370 Guo, Z., Wang, Y., Wan, Z., Zuo, Y., He, L., Li, D., Yuan, F., Wang, N., Liu, J., & Song, Y. (2020). Soil dissolved
371 organic carbon in terrestrial ecosystems: Global budget, spatial distribution and controls. *Global Ecology and*
372 *Biogeography*, 29(12), 2159-2175. <https://doi.org/10.1111/geb.13186>
- 373 He, X., Augusto, L., Goll, D. S., Ringeval, B., Wang, Y., Helfenstein, J., Huang, Y., Yu, K., Wang, Z., Yang, Y., &
374 Hou, E. (2021). Global patterns and drivers of soil total phosphorus concentration. *Earth System Science Data*,
375 13(12), 5831-5846. <https://doi.org/10.5194/essd-13-5831-2021>
- 376 IPCC. (2014). *Climate Change 2014: Synthesis Report*. In Core Writing Team, R. K. Pachauri & L. A. Meyer (Eds.),
377 *Contribution of Working Groups I, II and III to the Fifth Assessment Report of the Intergovernmental Panel on*
378 *Climate Change*. Geneva, Switzerland: IPCC. (2). <https://epic.awi.de/id/eprint/37530/>
- 379 Kaiser, K., Guggenberger, G., Haumaier, L., & Zech, W. (2005). Dissolved organic matter sorption on sub soils and
380 minerals studied by ¹³C-NMR and DRIFT spectroscopy. *European Journal of Soil Science*, 48(2), 301-310.
381 <https://doi.org/10.1111/j.1365-2389.1997.tb00550.x>
- 382 Kaiser, K., & Zech, W. (2000). Dissolved organic matter sorption by mineral constituents of subsoil clay fractions.
383 *Journal of Plant Nutrition and Soil Science*, 163(5), 531-535. [https://doi.org/10.1002/1522-2624\(200010\)163:5%3C531::AID-JPLN531%3E3.0.CO;2-N](https://doi.org/10.1002/1522-2624(200010)163:5%3C531::AID-JPLN531%3E3.0.CO;2-N)
- 385 Kalbitz, K., Solinger, S., Park, J.-H., Michalzik, B., & Matzner, E. (2000). Controls on the dynamics of dissolved
386 organic matter in soils: a review. *Soil Science*, 165(4), 277-304. [https://doi.org/10.1097/00010694-200004000-](https://doi.org/10.1097/00010694-200004000-00001)
387 00001



- 388 Kindler, R., Siemens, J., Kaiser, K., Walmsley, D. C., Bernhofer, C., Buchmann, N., Cellier, P., Eugster, W., Gleixner,
389 G., & Grünwald, T. (2011). Dissolved carbon leaching from soil is a crucial component of the net ecosystem
390 carbon balance. *Global Change Biology*, 17(2), 1167-1185. <https://doi.org/10.1111/j.1365-2486.2010.02282.x>
- 391 Lal, R. (2004). Soil carbon sequestration impacts on global climate change and food security. *science*, 304(5677),
392 1623-1627. <https://doi.org/10.1126/science.1097396>
- 393 Langeveld, J., Bouwman, A. F., van Hoek, W. J., Vilmin, L., Beusen, A. H. W., Mogollón, J. M., & Middelburg, J. J.
394 (2020). Estimating dissolved carbon concentrations in global soils: a global database and model. *Sn Applied*
395 *Sciences*, 2(10). <https://doi.org/10.1007/s42452-020-03290-0>
- 396 Lapierre, J. F., Seekell, D. A., & Del Giorgio, P. A. (2015). Climate and landscape influence on indicators of lake
397 carbon cycling through spatial patterns in dissolved organic carbon. *Global Change Biology*, 21(12), 4425-4435.
398 <https://doi.org/10.1111/gcb.13031>
- 399 Li, mengfan, Wang, J., Guo, D., Yang, R., & Fu, H. (2019). Effect of land management practices on the
400 concentration of dissolved organic matter in soil: A meta-analysis. *Geoderma*, 344, 74-81.
401 <https://doi.org/10.1016/j.geoderma.2019.03.004>
- 402 Li, J., Wu, B., Zhang, D., & Cheng, X. (2023). Elevational variation in soil phosphorus pools and controlling factors
403 in alpine areas of Southwest China. *Geoderma*, 431. <https://doi.org/10.1016/j.geoderma.2023.116361>
404 <https://doi.org/10.1016/j.geoderma.2023.116361>
- 405 Li, S., Zheng, X., Liu, C., Yao, Z., Zhang, W., & Han, S. (2018). Influences of observation method, season, soil
406 depth, land use and management practice on soil dissolvable organic carbon concentrations: A meta-analysis.
407 *Science of the Total Environment*, 631-632, 105-114. <https://doi.org/10.1016/j.scitotenv.2018.02.238>
- 408 Lu, S., Xu, Y., Fu, X., Xiao, H., Ding, W., & Zhang, Y. (2019). Patterns and drivers of soil respiration and vegetation
409 at different altitudes in Southern China. *Applied Ecology & Environmental Research*, 17(2).
410 https://doi.org/10.15666/aeer/1702_30973106
- 411 Nakhavali, M., Lauerwald, R., Regnier, P., Guenet, B., Chadburn, S., & Friedlingstein, P. (2020). Leaching of
412 dissolved organic carbon from mineral soils plays a significant role in the terrestrial carbon balance. *Glob Chang*
413 *Biol*. <https://doi.org/10.1111/gcb.15460>



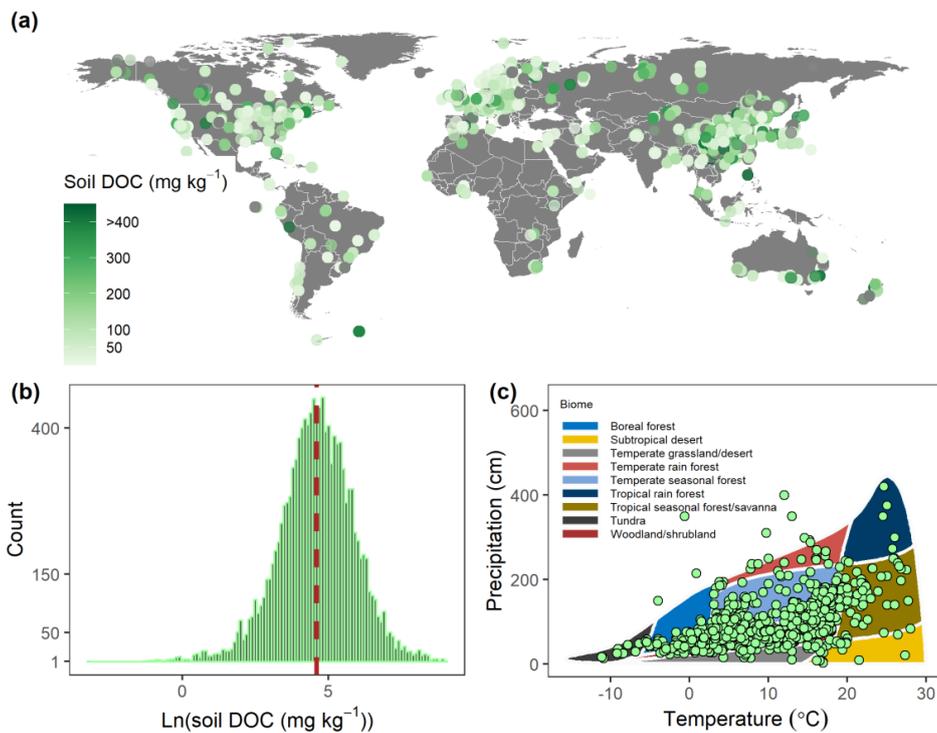
- 414 Nakhavali, M., Lauerwald, R., Regnier, P., Guenet, B., Chadburn, S., & Friedlingstein, P. (2021). Leaching of
415 dissolved organic carbon from mineral soils plays a significant role in the terrestrial carbon balance. *Global*
416 *Change Biology*, 27(5), 1083-1096. <https://doi.org/10.1111/gcb.15460>
- 417 Neff, J. C., & Asner, G. P. (2001). Dissolved organic carbon in terrestrial ecosystems: synthesis and a model.
418 *Ecosystems*, 4(1), 29-48. <https://doi.org/10.1007/s100210000058>
- 419 Nottingham, A. T., Baath, E., Reischke, S., Salinas, N., & Meir, P. (2019). Adaptation of soil microbial growth to
420 temperature: Using a tropical elevation gradient to predict future changes. *Glob Chang Biol*, 25(3), 827-838.
421 <https://doi.org/10.1111/gcb.14502>
- 422 Pesántez, J., Mosquera, G. M., Crespo, P., Breuer, L., & Windhorst, D. (2018). Effect of land cover and
423 hydro-meteorological controls on soil water DOC concentrations in a high-elevation tropical environment.
424 *Hydrological Processes*, 32(17), 2624-2635. <https://doi.org/10.1002/hyp.13224>
- 425 Petrone, K. C. (2005). Export of carbon, nitrogen and major solutes from a boreal forest watershed: The influence of
426 fire and permafrost. University of Alaska Fairbanks.
- 427 Rahbek, C., Borregaard, M. K., Colwell, R. K., Dalsgaard, B., Holt, B. G., Morueta-Holme, N., Nogues-Bravo, D.,
428 Whittaker, R. J., & Fjeldså, J. (2019). Humboldt's enigma: What causes global patterns of mountain biodiversity?
429 *Science*, 365(6458), 1108-1113. <https://www.science.org/doi/10.1126/science.aax0149>
- 430 Ren, T., Ukalska-Jaruga, A., Smreczak, B., & Cai, A. (2024). Dissolved organic carbon in cropland soils: A global
431 meta-analysis of management effects. *Agriculture, Ecosystems & Environment*, 371, 109080.
432 <https://doi.org/10.1016/j.agee.2024.109080>
- 433 Ren, T., & Cai, A. (2024). Global patterns and drivers of soil dissolved organic carbon concentrations. figshare [data
434 set]. <https://doi.org/https://doi.org/10.6084/m9.figshare.26379898>
- 435 Sanderman, J., & Amundson, R. (2008). A comparative study of dissolved organic carbon transport and stabilization
436 in California forest and grassland soils. *Biogeochemistry*, 89, 309-327. [https://doi.org/10.1007/s10533-008-](https://doi.org/10.1007/s10533-008-9221-8)
437 [9221-8](https://doi.org/10.1007/s10533-008-9221-8)
- 438 Singh, M., Sarkar, B., Hussain, S., Ok, Y. S., Bolan, N. S., & Churchman, G. J. (2017). Influence of physico-
439 chemical properties of soil clay fractions on the retention of dissolved organic carbon. *Environmental*
440 *Geochemistry and Health*, 39, 1335-1350. <https://doi.org/10.1007/s10653-017-9939-0>



- 441 Smreczak, B., & Ukalska-Jaruga, A. (2021). Dissolved organic matter in agricultural soils. *Soil Science Annual*.
442 <https://doi.org/10.37501/soilsa/132234>
- 443 Wang, Y.-L., Chang-Ming, Y., Li-Min, Z., & Heng-Zhao, C. (2015). Spatial distribution and fluorescence properties
444 of soil dissolved organic carbon across a riparian buffer wetland in Chongming Island, China. *Pedosphere*, 25(2),
445 220-229. [https://doi.org/10.1016/S1002-0160\(15\)60007-8](https://doi.org/10.1016/S1002-0160(15)60007-8)
- 446 Wei, D., Tao, J., Wang, Z., Zhao, H., Zhao, W., & Wang, X. (2024). Elevation-dependent pattern of net CO₂ uptake
447 across China. *Nature Communications*, 15(1), 2489. <https://doi.org/10.1038/s41467-024-46930-4>
- 448 Zsolnay, Á. (2003). Dissolved organic matter: artefacts, definitions, and functions. *Geoderma*, 113(3-4), 187-209.
449 [https://doi.org/10.1016/s0016-7061\(02\)00361-0](https://doi.org/10.1016/s0016-7061(02)00361-0)
450
451



452 **Figure 1** Global distribution of soil dissolved organic carbon (DOC) concentration according to our site-level
453 dataset. The dataset contains 12807 sets of data **(a, b)**, which covers major terrestrial biomes **(c)**. The dashed red line
454 within the subplot **(b)** signifies the average soil DOC concentration, which is 223 mg kg⁻¹.



455

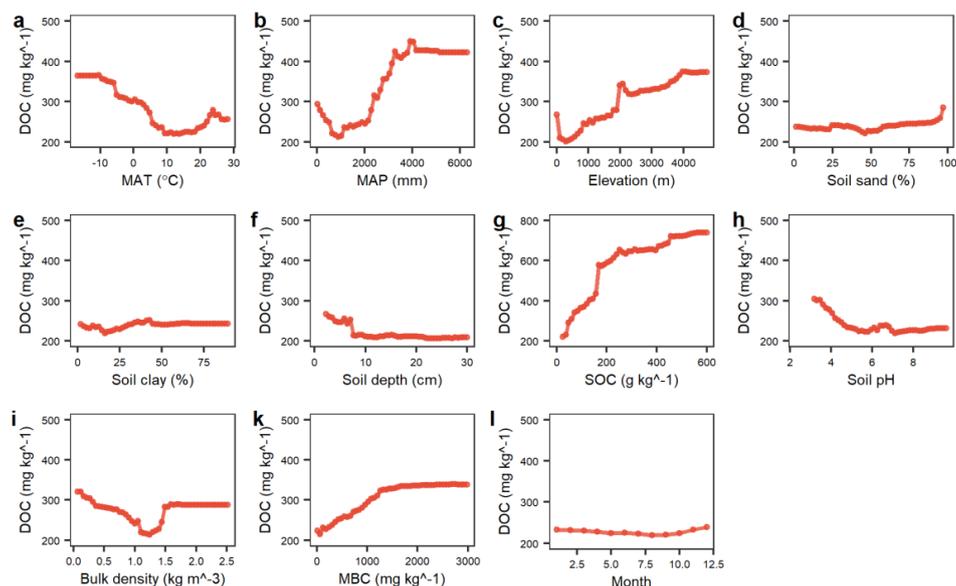
456

457

458



459 **Figure 2** Partial dependence of predictors from random forest algorithm. Soil dissolved organic carbon (DOC)
460 concentration in relation to mean annual temperature (MAT), mean annual precipitation (MAP), elevation, soil sand
461 content, soil clay content, soil depth, soil organic carbon (SOC) content, soil pH, bulk density, microbial biomass
462 carbon content (MBC), and month (a, b, c, d, e, f, g, h, i, k, l, respectively).



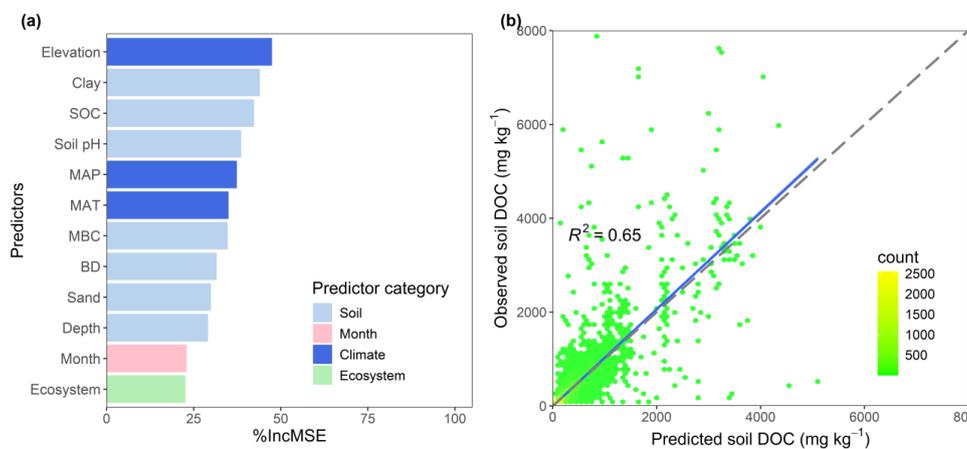
463

464

465



466 **Figure 3** Result of the random forest model predicting soil dissolved organic carbon (DOC) concentration. (a) The
467 relative importance of predictors in the random forest model. (b) Predicted vs. observed soil DOC concentration.
468 The dashed line indicates the 1:1 line and the blue line indicates the regression line between predicted and observed
469 values.

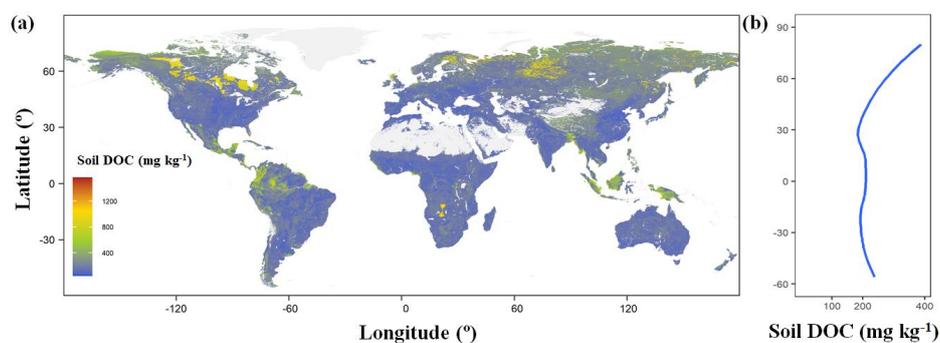


470

471



472 **Figure 4** Prediction of soil dissolved organic carbon (DOC) concentration in global ecosystems. (a) Global map of
473 predicted soil DOC concentration. (b) Latitudinal patterns of soil DOC concentration. Blue line indicates the locally
474 weighted regressions between latitude and soil DOC concentration in the predicted global map. Values in the
475 predicted map reflect soil DOC concentration within a grid cell resolution of $0.05^\circ \times 0.05^\circ$. A value in the grid is the
476 averaged from the result of random forest model.



477

478



479 **Table 1.** Variables information of soil dissolved organic carbon dataset in global terrestrial ecosystems. n/a refers to
 480 values that are not applicable.

Variables	Description	Unit	Number	Range	Mean
No.	Unique identification number of each record	n/a	12807	1 to 12807	6404
Latitude	Latitude of study site	°	12807	-64.81 to 78.85	34.89
Longitude	Longitude of study site	°	12807	-159.66 to 175.95	107.05
MAT	Mean annual temperature	°C	9948	-11.16 to 28.00	11.84
MAP	Mean annual precipitation	mm	10325	30 to 4200	1071
Elevation	Altitude of study site	m	5578	4 to 4730	881
Ecosystems	Community by the dominant plant species		7	n/a	n/a
Soil sand	Soil sand content	%	4062	1 to 98	45
Soil silt	Soil silt content	%	4025	1 to 95	33
Soil clay	Soil clay content	%	4316	0 to 89	22
Soil depth	Mean depth of soil sample	cm	12807	0.53 to 30.00	11.36
SOC	Soil organic carbon	g kg ⁻¹	9136	0.23 to 598.50	38.74
TN	Soil total nitrogen	g kg ⁻¹	7089	0.00 to 33.30	2.57
Soil pH	Measure by 1:2.5 H ₂ O,	n/a	8266	2.30 to 9.59	6.16
BD	Soil bulk density	kg m ⁻³	4380	0.07 to 2.52	1.29
MBC	Soil microbial biomass carbon	mg kg ⁻¹	4218	5.93 to 2986	413
Date	Observation month of DOC	month	12807	1 to 12	6.50
DOC _{phy}	Measure by physical method	mg kg ⁻¹	3289	0.28 to 3181	155.99
DOC _{che}	Measure by chemical process	mg kg ⁻¹	9518	0.04 to 7859	245.83
DOC	Soil dissolved organic carbon	mg kg ⁻¹	12807	0.04 to 7859	222.78

481

482



483 **Table 2.** Global soil dissolved organic carbon concentration (mg kg^{-1}) for major ecosystems. 25% and 75% represent
484 the 25th and 75th percentiles of one group, respectively. SD, Standard deviation; SE, Standard error.

Ecosystems	Mean	SD	SE	Skewness	Kurtosis	25%	Median	75%
Wetland	218.53	340.35	10.23	5.15	39.41	46.40	107.11	266.51
Forest	256.18	531.72	7.62	7.09	69.72	47.60	115.51	246.55
Shrubland	160.24	131.51	6.70	3.40	22.58	76.53	127.84	205.50
Tundra	470.78	721.70	63.30	4.67	29.59	86.91	241.09	577.00
Grassland	327.77	674.43	19.53	4.16	18.03	54.62	126.48	303.63
Cropland	165.98	272.51	3.81	6.53	73.25	40.51	83.00	178.81
Global	222.78	445.78	3.93	7.16	73.67	45.86	101.01	226.47

485



486 **Table 3.** Analysis of the predicted global map of soil dissolved organic carbon. The area-weighted average soil
487 dissolved organic carbon concentration was calculated based on our predicted map. Converting soil dissolved
488 organic carbon concentration to soil dissolved organic carbon content and stock used the soil bulk density and land
489 area.

Continent	Soil DOC concentration (mg kg ⁻¹)	Soil DOC content (g m ⁻²)	Soil DOC stock (Pg)
Asia	274.43	107.79	4.80
North America	263.63	99.37	2.45
Europe	227.34	86.76	0.88
South America	215.81	77.05	1.37
Oceania	198.13	76.92	0.59
Africa	186.35	68.04	2.07
Global	237.56	89.80	12.17

490

491