

1 An ensemble estimate of Australian soil organic carbon 2 using machine learning and process-based modelling

3
4 Lingfei Wang^{1,2}, Gab Abramowitz^{1,2}, Ying-Ping Wang³, Andy Pitman^{1,2} and Raphael A.
5 Viscarra Rossel⁴

6
7 ¹ ARC Center of Excellence for Climate Extremes, Sydney NSW 2052, Australia

8 ² Climate Change Research Center, University of New South Wales, Sydney NSW 2052, Australia

9 ³ CSIRO Environment, Private Bag 10, Clayton South VIC 3169, Australia

10 ⁴ Soil & Landscape Science, School of Molecular & Life Sciences, Faculty of Science & Engineering,
11 Curtin University, GPO Box U1987, Perth WA 6845, Australia.

12
13 Correspondence to: Lingfei Wang (lingfei.wang@unsw.edu.au)

14 Abstract

15
16 Spatially explicit prediction of soil organic carbon (SOC) serves as a crucial foundation for
17 effective land management strategies aimed at mitigating soil degradation and assessing carbon
18 sequestration potential. Here, using more than 1000 in-situ observations, we trained two
19 machine learning models (random forest, and K-means coupled with multiple linear
20 regression), and one process-based model (the vertically resolved Microbial-MIneral Carbon
21 Stabilization (MIMICS)) to predict SOC stocks of the top 30 cm of soil in Australia. Parameters
22 of MIMICS were optimized for different site groupings, using two distinct approaches, plant
23 functional types (MIMICS-PFT), and the most influential environmental factors (MIMICS-
24 ENV). All models showed good performance in SOC predictions with R^2 greater than 0.8
25 during out-of-sample validation with random forest being the most accurate, and SOC in forests
26 is more predictable than that in non-forest soils excluding croplands. The performance of
27 continental-scale SOC predictions by MIMICS-ENV is better than that by MIMICS-PFT
28 especially in non-forest soils. Digital maps of terrestrial SOC stocks generated using all the
29 models showed similar spatial distribution with higher values in southeast and southwest
30 Australia, but the magnitude of estimated SOC stocks varied. The mean ensemble estimate of
31 SOC stocks was 30.3 t ha^{-1} with K-means coupled with multiple linear regression generating
32 the highest estimate (mean SOC stocks at 38.15 t ha^{-1}) and MIMICS-PFT generating the lowest
33 estimate (mean SOC stocks at 24.29 t ha^{-1}). We suggest that enhancing process-based models
34 to incorporate newly identified drivers that significantly influence SOC variations in different
35 environments could be key to reducing the discrepancies in these estimates. Our findings
36 underscore the considerable uncertainty in SOC estimates derived from different modelling
37 approaches and emphasize the importance of rigorous out-of-sample validation before applying
38 any one approach in Australia.

41 1. Introduction

42
43 Globally, the soil is the largest biogeochemically active terrestrial carbon pool, storing more
44 organic carbon than plants and the atmosphere combined (Jackson et al., 2017). The turnover
45 of soil organic carbon (SOC) is a key function in plant growth, maintenance of soil water and
46 nutrients, soil structure stabilization and other biogeochemical processes (Lefèvre et al., 2017).
47 Soil can act as either a carbon sink or carbon source depending on the balance of carbon input
48 through plant litter and root exudates and output through respiration and leaching (Terrer et al.,
49 2021; Panchal et al., 2022). Even a small change in SOC stocks, in any direction, could
50 significantly affect the atmospheric concentration of CO₂ and thereby climate change
51 (Stockmann et al., 2013).

52
53 Given the importance of SOC, there is now a large and growing interest in estimating spatially
54 explicit SOC content and stocks. SOC supports critically important soil-derived ecosystem
55 services, and the amount of SOC indicates the degree of land and soil degradation (Lorenz et
56 al., 2019). SOC content below a certain limit will lead to the decline of microbial diversity,
57 water holding capacity and soil productivity (Stockmann et al., 2015). Additionally, with
58 growing concerns about increasing anthropogenic CO₂ emissions, soil carbon sequestration has
59 emerged as a potential strategy for climate change mitigation (Smith, 2016; Rumpel et al.,
60 2018). Protection of existing SOC and rebuilding depleted stocks through land management are
61 potential strategies in mitigating climate change (Bossio et al., 2020). However, effective SOC
62 management requires accurate knowledge of its existing distribution. Reliable estimates of SOC
63 stocks and their spatial variation serve as a reference point for assessing how close soil is to its
64 maximum SOC storage capacity and its potential to sequester additional carbon (Six et al.,
65 2002; Georgiou et al., 2022). Precise estimation of contemporary SOC stocks also provides a
66 baseline map that can be used to calibrate and initialize dynamic-mechanistic models, enabling
67 the study of how SOC will respond to climate and land-use change (Minasny et al., 2013;
68 Viscarra Rossel et al., 2014). It is, for example, a prerequisite for accurately predicting future
69 carbon–climate feedback in Earth system models (ESMs) (Todd-Brown et al., 2013).

70
71 Accurately assessing SOC storage is challenging due to the complexity of carbon formation
72 and degradation processes in space and time (Keskin et al., 2019). Soil exists as a continuum
73 containing organic compounds at different stages of decomposition (Lehmann and Kleber,
74 2015). Soil formation can be described by a function of climate, organisms, relief, parent
75 material and time (Jenny, 1941). These factors are widely used in SOC studies for digital soil
76 mapping (McBratney et al., 2003; Viscarra Rossel et al., 2015; Liang et al., 2019). However,
77 the relationship between SOC storage and these driving variables is complex and spatially
78 variable (Mishra and Riley, 2015; Viscarra Rossel et al., 2019; Adhikari et al., 2020) leading to
79 substantial challenges and inherent uncertainties in SOC predictions.

80
81 Mechanistic process-based models and empirical models (including machine learning models)
82 are two widely employed approaches used to predict SOC stocks and their spatial distribution.

83 Conventional process-based models assume first-order kinetics for SOC decomposition,
84 wherein the rate of C decomposition is dependent on temperature and moisture but independent
85 of microbial biomass, and equilibrium SOC stock is proportional to carbon input and mean
86 residence time (Abs and Ferrière, 2020; Wang et al., 2021). ESMS coupled with conventional
87 SOC models cannot accurately simulate spatial pattern of contemporary soil carbon and show
88 large divergence in projected SOC dynamics under future climate change (Todd-Brown et al.,
89 2013; Todd-Brown et al., 2014). In addition to the biases introduced by errors in model
90 parameters and the lack of independent model validation based on observed time series data,
91 the uncertainties in predicted SOC by ESMS can also result from the lack of explicit
92 representation of soil microbial activities and metabolic traits (Wieder et al., 2015; Le Neo et
93 al., 2023). Numerous microbial models have been developed in the past few decades to improve
94 model performance of SOC predictions (Chandel et al., 2023), but these models have rarely
95 been incorporated into large-scale modelling frameworks due to the difficulty of constraining
96 parameters relating to microbial activities and the lack of rigorous validation (Todd-Brown et
97 al., 2013; Luo et al., 2016). Process-based SOC models are constructed based on our
98 understanding of the major processes governing SOC dynamics (e.g., carbon input,
99 decomposition, and loss). However, the disagreement in projections of carbon dynamics by
100 different models highlights the need to improve our knowledge of SOC cycling (Luo et al.,
101 2016). Machine learning models without any process-level assumptions provide a tool to
102 identify the most influential controls on SOC variations. Machine learning models can represent
103 non-linear and non-smooth relationships between predictor and response variables as well as
104 interactions between different predictors (Heung et al., 2016). Various machine learning
105 algorithms have been successfully used in digital soil mapping to predict high-resolution
106 spatially explicit SOC concentration/stocks (Lamichhane et al., 2019).

107
108 Several modelling studies of soil carbon stocks have been conducted in Australia. Wang et al.
109 (2018a) trained boosted regression trees and random forest models using field observations and
110 applied the trained random forest model to map the spatial distribution of SOC at two soil depths
111 (0-5 cm and 0-30 cm) for the semi-arid rangelands of eastern Australia. Continentally, Viscarra
112 Rossel et al. (2014) trained the CUBIST model, a form of piecewise linear decision tree, using
113 more than five thousand observations to produce a high resolution (90 m × 90 m) baseline map
114 of SOC stocks of Australian terrestrial systems and its uncertainty of the top 30 cm soils. Based
115 on the baseline map, Walden et al. (2023) derived spatially explicit estimates of Australian SOC
116 stocks and uncertainty including additional data from forests from southeastern Australia and
117 coastal marine (or blue carbon) ecosystems. SOC content at multiple soil depths along with
118 associated uncertainties were also estimated using different machine learning algorithms
119 (Viscarra Rossel et al., 2015; Wadoux et al., 2023). Moreover, the distribution of different soil
120 carbon compositions (i.e., the particulate, mineral-associated and pyrogenic organic carbon
121 fractions) and the importance of environmental factors on their variations were also studied
122 using machine learning (Viscarra Rossel et al., 2019). However, despite the progress made in
123 SOC modelling, significant uncertainties persist in SOC estimates due to the inherent
124 complexities of SOC variations and the lack of appropriately sampled SOC observations. All

125 these continental estimates were generated using empirical modelling approaches or first-order
126 biogeochemical models without explicitly representing the important role of soil microbes in
127 SOC stabilization (Grace et al., 2006; Lee et al., 2021). Estimates from mechanistic SOC
128 models with explicit representation of microbial metabolism are missing despite offering the
129 potential to better constrain SOC dynamics under future climate change scenarios in a way that
130 empirical approaches cannot.

131
132 Our primary objective in this paper is to assess the predictability of SOC concentration
133 (excluding cropland soils) in Australia and generate a range of estimates of terrestrial SOC
134 stocks, employing both process-based and empirical modelling, and examine why these
135 estimates might differ. First, we discern the significance of environmental predictors, both at
136 continental and biome scales. We then evaluate the performance of random forests, K-means
137 with multiple linear regression and the vertically resolved MIMICS model with different
138 parametrization approaches. Finally, we compare the spatial estimates of SOC stocks using
139 these different approaches across Australia, and discuss their differences and potential
140 application to future SOC projection.

141

142 2. Materials and Methods

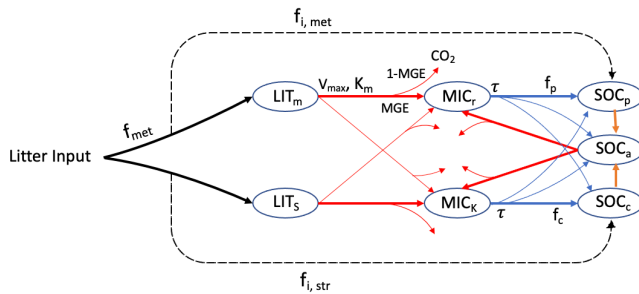
143 2.1. Model descriptions

144 2.1.1. Vertically resolved MIMICS

145

146 The MIMICS model (Wieder et al., 2015; Zhang et al., 2020) explicitly considers relationships
147 between litter quality, functional trade-offs in microbial physiology, and the physical protection
148 of microbial by-products in forming stable soil organic matter. There are two litter pools:
149 metabolic (LIT_m) and structural (LIT_s) litter (Figure 1), and the partitioning of litter input into
150 metabolic and structural pools is determined by the chemical properties of the litter. Litter and
151 SOC turnover are governed by two microbial functional types that exhibit copiotrophic (i.e., r-
152 selected, MIC_r) and oligotrophic (i.e., K-selected, MIC_k) growth strategies. The MIC_r is
153 assumed to have higher growth and turnover rates, and a preference for consuming labile litter
154 (LIT_m), while MIC_k is characterized by lower growth and turnover rates, and a greater
155 competitive advantage when consuming low-quality litter (LIT_s) and chemically recalcitrant
156 SOC. SOC in MIMICS is divided into three pools: physically protected (SOC_p),
157 (bio)chemically recalcitrant (SOC_c) and available (SOC_a) carbon (Figure 1).

158



159 **Figure 1.** SOC pools and fluxes represented in MIMICS (adapted from Wieder et al., (2015)). Litter inputs
 160 are partitioned into metabolic and structural litter pools (LIT_m and LIT_s) based on litter quality (f_{met}).
 161 Decomposition of litter and available SOC pool (SOC_a) are governed by temperature sensitive Michaelis-
 162 Menten kinetics (V_{max} (maximum reaction velocity) and K_m (half saturation constant)), shown by red lines.
 163 Microbial growth efficiency (MGE) determines the partitioning of C fluxes entering microbial biomass pools
 164 vs. heterotrophic respiration. Turnover of microbial biomass (τ , blue) depends on microbial functional types
 165 (MIC_r and MIC_k), and is partitioned into available, physically protected and chemically recalcitrant SOC
 166 pools (SOC_a , SOC_p and SOC_c , respectively).
 167

168
 169 The decomposition of litter pools and SOC pools follows temperature-sensitive Michaelis-
 170 Menten kinetics. Microbial growth efficiency (MGE) determines the partitioning of carbon
 171 fluxes entering microbial biomass pools (MIC_r and MIC_k) versus heterotrophic respiration.
 172 Access of microbial enzymes to available substrates depends on soil texture. The equations of
 173 MIMICS are from Wieder et al. (2015), except that the density-dependent microbial turnover
 174 was introduced to MIMICS to minimize an unrealistic oscillation (Zhang et al., 2020). To better
 175 simulate carbon turnover at different soil depths, vertical transport of soil carbon was introduced
 176 into MIMICS considering carbon transported through bioturbation and diffusion among
 177 adjacent soil layers (Wang et al., 2021).
 178

179 Vertically resolved MIMICS is run using a daily time step. The soil was divided into 15 layers,
 180 each of 10 cm thickness. All the sites in this study are assumed to be at steady state (i.e., no
 181 interannual variation of SOC). Historical climate, litterfall input and soil properties were all
 182 assumed to be similar to the average conditions. At each site, the initial pool fractions were
 183 0.03, 0.03, 0.14, 0.47 and 0.33 for MIC_r , MIC_k , SOC_p , SOC_c and SOC_a , respectively. All pools
 184 were then spun up to finally achieve steady state with the maximal difference in any pool size
 185 between two successive spins being less than 0.05%.
 186

187 2.1.2. Machine learning

188
 189 Two machine learning algorithms were applied in this study to predict SOC. First, random forest
 190 (RF) is a tree-based ensemble learning method that works by building a set of regression trees
 191 and averaging results (Breiman, 2001). Within the training procedure, the RF algorithm
 192 produces multiple trees. Each regression tree in the forest is independently constructed based
 193 on a unique bootstrap sample (with replacement) from the original training data set. The

194 response, as well as the predictor variables are either categorical (classification trees) or
195 numeric (regression trees). Bootstrap sampling makes RF less sensitive to overfitting and
196 allows for robust error estimation based on the remaining test set, the so-called Out-Of-Bag
197 (OOB) sample (Wiesmeier et al., 2014). We used the “ranger” package R (version 4.2.0) for RF
198 computation. We trained the RF model with different numbers (100, 200, 300, 400 and 500) of
199 trees and observed that the model's performance remained similar regardless of the number of
200 trees used. The number of regression trees generated in the forest (num.trees) was finally set as
201 200, and the number of predictors randomly selected at each node (mtry) was set as default,
202 which was 2.

203
204 Multiple linear regression (MLR) is widely used in SOC studies but found to be less effective
205 than machine learning algorithms (Lamichhane et al., 2019). Here, instead of applying MLR
206 directly with all environmental factors as predictors, our approach involved a preliminary step
207 where we partitioned all observations into distinct clusters using K-means, an unsupervised
208 machine learning algorithm. K-means aims to divide the data into a predefined number of
209 clusters (k), with the objective of maximizing the similarity among data within each cluster.
210 The underlying assumption here was that sites sharing similar environmental conditions would
211 exhibit comparable SOC concentration. In cases where certain clusters had fewer observations
212 than five times the number of predictors, we augmented these clusters by incorporating
213 observations from other clusters. This augmentation process was guided by the Euclidean
214 distance between the observation and the cluster centre, ensuring a more robust construction of
215 the linear regression model. To determine the number of clusters, we applied the coupled K-
216 means and MLR with varying number of clusters. The selection of the optimal number of
217 clusters was based on the criterion of producing the smallest root mean square error during
218 independent out-of-sample validation.

219 220 2.2. Relative importance of environmental variables for SOC prediction

221
222 RF-based measures of variable importance have gained widespread popularity as tools for
223 evaluating the contributions made by predictor variables within a fitted random forest model
224 (Debeer and Strobl, 2020). In the context of this study, we employed permutation variable
225 importance (PVI) within the random forest framework to gauge the significance of predictors
226 (see Section 2.4) in predicting SOC concentration.

227
228 The PVI entails measuring the reduction in a RF model's performance score upon random
229 shuffling of a single variable values. By doing so, the inherent relationship between the variable
230 and the SOC concentration is disrupted. Consequently, the disparity in prediction accuracy
231 observed in a RF model before and after such shuffling serves as a quantitative representation
232 of the significance of the particular predictor in predicting SOC concentration. The greater the
233 importance of the predictor, the higher its corresponding PVI value becomes.

234
235

236 2.3. Parameter optimization

237
238 MIMICS parameters were derived from Zhang et al. (2020) and Wang et al. (2021), except that
239 five parameters (Table 1) which directly control the organic carbon decomposition were
240 optimized. An effective global optimization algorithm called the shuffled complex evolution
241 (SCE-UA, version 2.2) method (Duan et al., 1993) was applied for parameter optimization by
242 minimizing the sum of squared residuals between the observed and modelled values.

243
244 Vertically resolved MIMICS simulated SOC concentration for 15 soil layers with a uniform
245 layer thickness of 10 cm. As observations only provide one measurement for the top 30 cm soil,
246 we computed the average of the modelled values spanning the 0-10 cm, 10-20 cm, and 20-30
247 cm soil layers. This average was then adopted as the modelled SOC concentration for top 30
248 cm soil, serving as the basis for evaluating the difference between observations and simulations.

249
250 **Table 1.** The optimized model parameters (dimensionless) and their value range.

Parameter	Definition	Range
a_v	A scaling factor for V_{\max}	0-30
a_k	A scaling factor for K_m	0-20
xdesorp	A scaling factor for SOC desorption rate	0-3
xbeta	An exponent of the biomass density dependent mortality rate of microbes	1.05-2
xdiffsoc	A scaling factor for SOC diffusion coefficient in soil	0-30

251
252 Parameters in MIMICS were optimized for different groups divided based on two approaches.
253 The first approach involved categorizing all observations into four groups based on plant
254 functional type (PFT). The second approach used the most influential abiotic variables as
255 predictors (as outlined in Section 2.2) and divided all observations into 6 clusters using the K-
256 means algorithm. The determination of the optimal number of clusters was achieved through
257 the minimization of the sum of the within-cluster-sum-of-squares-of-all-clusters (WCSSE), a
258 process facilitated by the "ClusterR" package in R (version 4.2.0). This clustering aimed to
259 ensure the highest possible similarity among the environmental factors within each cluster. It
260 was anticipated that SOC ranges within each cluster would be narrow due to the high similarity
261 of environmental predictors.

262

263 2.4. Data

264 2.4.1. Predictors of spatial variations of observed SOC concentration

265
266 MIMICS requires gridded mean annual temperature (MAT), carbon input and clay content as
267 driving variables for a spatial simulation. Gridded mean annual precipitation (MAP) and
268 vegetation types were also used during calibration and when understanding the drivers and
269 spatial variability of SOC. Details of gridded data can be found in Table 2.

270

271 Gridded daily maximum temperature, minimum temperature, and precipitation at 0.05°
272 resolution were obtained from the SILO database (Jeffrey et al., 2001) of Australian climate
273 data. Mean daily temperature was approximated as the average of maximum and minimum

274 daily temperature. MAT was calculated from mean daily temperature from 1991 to 2020, and
275 MAP was calculated from daily precipitation from 1991 to 2020.

276
277 Carbon input was represented by NPP. Gridded mean annual NPP at 500 m was calculated
278 based on annual NPP from 2001 to 2020 obtained from MODIS (MOD17A3HGF V6.1)
279 (Running and Zhao, 2021). NPP was partitioned to above-/belowground part by multiplying by
280 the root/shoot ratio for different vegetation types (Mokany et al., 2006). Here we did not account
281 for the fraction of NPP that is appropriated by human activities.

282
283 The distribution of vegetation types at 3'' resolution was obtained from National Vegetation
284 Information System (NVIS, version 6.0, <https://www.dcceew.gov.au/environment/land/native-vegetation/national-vegetation-information-system>). Pixels of non-vegetated regions were
285 removed and 28 types from NVIS were aggregated to just 4 PFT: forest, woodland, shrubland
286 and grassland.

287
288
289 Soil bulk density and clay content were obtained from Soil and Landscape Grid National Soil
290 Attributes Maps (SLGA – Release 2) (Grundy et al., 2015; Viscarra Rossel et al., 2015). Soil
291 properties were predicted based on machine learning at depths 0-5 cm, 5-15 cm, 15-30 cm, 30-
292 60 cm, 60-100 cm, and 100-200 cm in SLGA. Bulk density and clay content were estimated for
293 top 30 cm soil as weighted average of first 3 layers in SLGA.

294
295 The initial spatial resolution of the gridded data was maintained when extracting the required
296 environmental factors for each SOC observation. All data were then resampled to 0.05°
297 resolution using bilinear interpolation for estimation of terrestrial SOC stocks at continental
298 scale.

299
300 **Table 2.** Information of gridded data used in this study.

	Source	Spatial Scale	Temporal Scale	Unit	Time Period
Maximum Temperature	SILO	~5 km	daily	°C	1991-2020
Minimum Temperature	SILO	~5 km	daily	°C	1991-2020
Precipitation	SILO	~5 km	daily	mm	1991-2020
NPP	MODIS	500 m	annually	g C/m ² /year	2001-2020
Vegetation Types	NVIS	100 m	/	/	/
Soil Bulk Density	SLGA	~90 m	/	kg/m ³	/
Soil Clay Content	SLGA	~90 m	/	%	/

301
302 2.4.2. Soil organic carbon observations

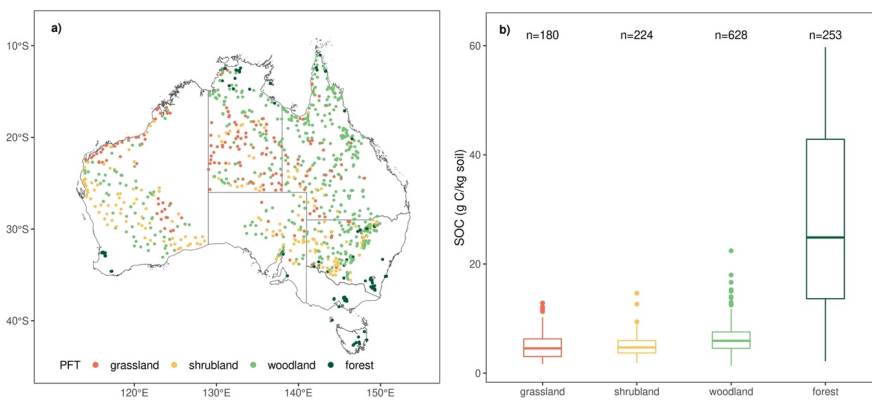
303
304 SOC observations for top 30 cm soil in Australia were collected from two datasets. The first
305 dataset is described in Viscarra Rossel et al. (2014) and Viscarra Rossel et al. (2019). We
306 removed the observations collected from croplands based on the land-use record in the dataset
307 and removed those from unvegetated regions based on NVIS vegetation map (see above). A
308 total of 1070 site observations with only 38 from forest soils were retained. SOC stocks were
309 reported in t ha⁻¹. To better represent SOC distribution in forest, we obtained additional forest
310 SOC observations from a second dataset, the Biomes of Australian Soil Environments (BASE)

311 described in Bissett et al. (2016). Here, SOC (%) was reported for 0-10 and 20-30 cm. We
312 estimated SOC for 0-30 cm soil following the method described in Viscarra Rossel et al. (2014).

313
314 To compare the observations with MIMICS outputs, we then converted both simulated SOC
315 (mg/cm^3) and observed SOC (t/ha) in the first dataset (Viscarra Rossel et al. 2014) to SOC
316 concentration ($\text{g C}/\text{kg soil}$) using spatially explicit soil bulk density (BD) from SLGA. The unit
317 conversion will not affect the results of MIMICS. Soil clay content is extracted from SLGA.

318
319 The spatial distribution of SOC observations from different PFT is shown in Figure 2a. SOC
320 concentration in top 30 cm is positively skewed, ranging from 1.36 to 59.73 $\text{g C}/\text{kg soil}$ with
321 mean value at 9.97 $\text{g C}/\text{kg soil}$ and median value at 6.11 $\text{g C}/\text{kg soil}$. SOC concentration in
322 grassland, shrubland and woodland show similar distribution patterns (Figure 2b), while SOC
323 concentration in forest is more variable with a standard deviation at 15.92 $\text{g C}/\text{kg soil}$.

324
325



326
327 **Figure 2.** a) Spatial distribution of 1285 soil organic carbon observations used in this study and the plant
328 functional types which they belong to; b) boxplots of SOC concentration distributions for each plant
329 functional type. For boxplots, centre lines represent the median value, and upper and lower box boundaries
330 represent third and first quartile. Whiskers extend to the smallest and largest values within 1.5 times the
331 interquartile range.
332

333
334
335

336 2.5. Model evaluation

337
338 For machine learning models, 70% of the observations were randomly selected as training data
339 to train the models and the remaining 30% used as test data to validate the predictions of SOC
340 concentration. For vertically resolved MIMICS, parameters were optimized for each PFT or
341 environmental group (see Section 2.3 above), and we again randomly selected 70% of

342 observations in each group to train the model and used the remaining 30% for validation. To
343 cross-validate, the procedure was repeated 10 times.

344
345 The performance of models was evaluated using four metrics. Mean Absolute Error (MAE)
346 indicates how close the average predictions are to average observations. Root Mean Square
347 Error (RMSE) measures the overall accuracy combining mean, standard deviation differences
348 (across sites) and (spatial) correlation. Coefficient of determination (R^2) measures the
349 percentage of variation explained by the model. Lin's Concordance Correlation Coefficient
350 (LCCC) (Lawrence and Lin, 1989) measures the level of agreement between predictions and
351 observations following the 1:1 line. A good model will have MAE and RMSE close to 0 and R^2
352 and LCCC close to 1.

353
354 2.6. Estimation of terrestrial SOC stocks

355
356 SOC concentrations were used to train the models, and we then estimated terrestrial SOC stocks
357 and their continental-scale spatial distribution in top 30 cm soil utilizing the four models
358 validated within this study. SOC stock ($t\ ha^{-1}$) is calculated using SOC concentration (g C/kg
359 soil), bulk density (BD, kg/m^3) and soil depth (m),

360
361
$$SOC_{stock} = SOC_{concentration} \times BD \times depth/100 \quad (2)$$

362

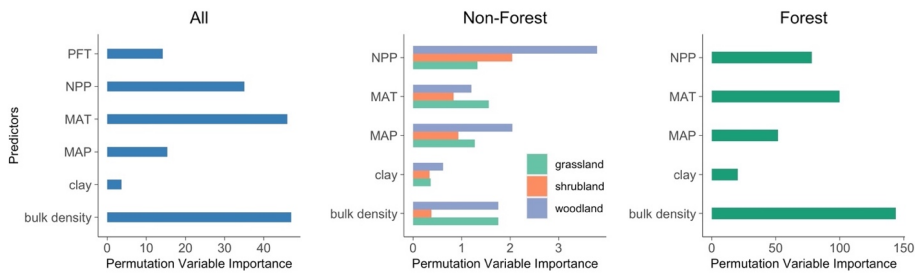
363 In the cases of MIMICS-PFT and MIMICS-ENV, the initial step involved grouping all pixels
364 into four distinct plant functional groups or six environmental clusters. Since cross-validation
365 was performed, the machine learning and process-based models were evaluated using test data,
366 and the models with the optimal performance were subsequently employed at each pixel to
367 estimate terrestrial SOC stocks. The map of ensemble estimate of SOC stocks was produced as
368 the average of four model estimates at each pixel.

369 3. Results

370
371 3.1. Relative importance of environmental predictors of SOC concentration

372
373 Using the PVI in random forest, we identified the significance of environmental factors in
374 predicting SOC concentration. At the continental scale, soil bulk density contributes most to the
375 prediction of SOC concentration, following by MAT, NPP and MAP (Figure 3). Soil clay
376 content and plant functional type exhibit relatively lesser significance in this regard.

377
378 The relative predictor importance for forests and grasslands aligns with the importance at
379 continental scale. In shrubland and woodland, NPP and MAP emerge as the pivotal factors.
380 Collectively, across both continental and regional scales, soil bulk density, MAT, and MAP are
381 the three most influential abiotic factors.

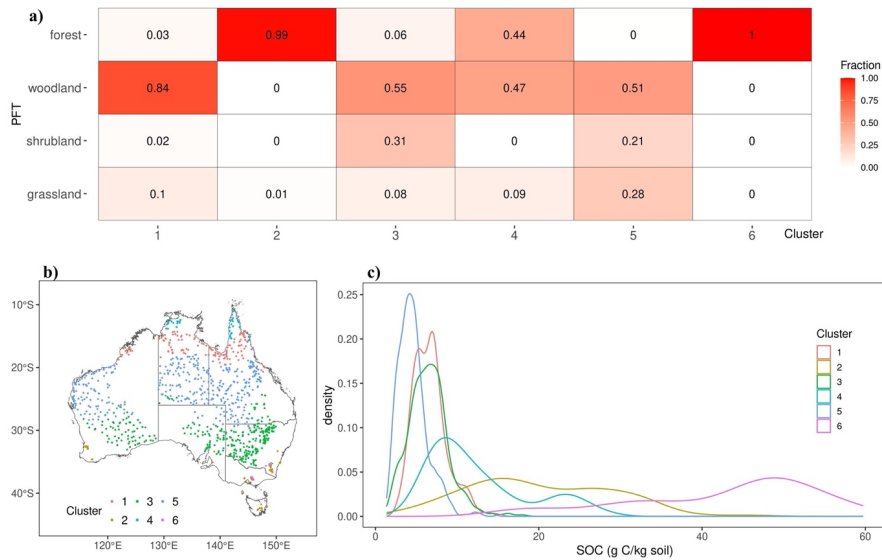


382 **Figure 3.** Importance of predictors on SOC concentration for different plant functional types.
 383

384 3.2. Data clustering based on environmental factors
 385

386 To develop the calibration groups for MIMICS-ENV, we partitioned the top three important
 387 abiotic factors, which are soil bulk density, MAT and MAP, into six distinct clusters using K-
 388 means (see Section 2.3). The resulting characteristics and spatial distribution of SOC belonging
 389 to these six clusters are illustrated in Figure 4.

390
 391 Notably, a substantial majority of forests were assigned to clusters 2 and 6 (Figure 4a), while
 392 woodland, shrubland, and grassland observations were distributed across the remaining four
 393 clusters. Among these clusters, cluster 5 exhibits the lowest SOC concentration, while SOC of
 394 cluster 1 and 3 display a comparable pattern but spread across different biomes. Conversely,
 395 distribution of SOC concentration in clusters 2, 4, and 6 shows more pronounced variability
 396 (Figure 4c).
 397

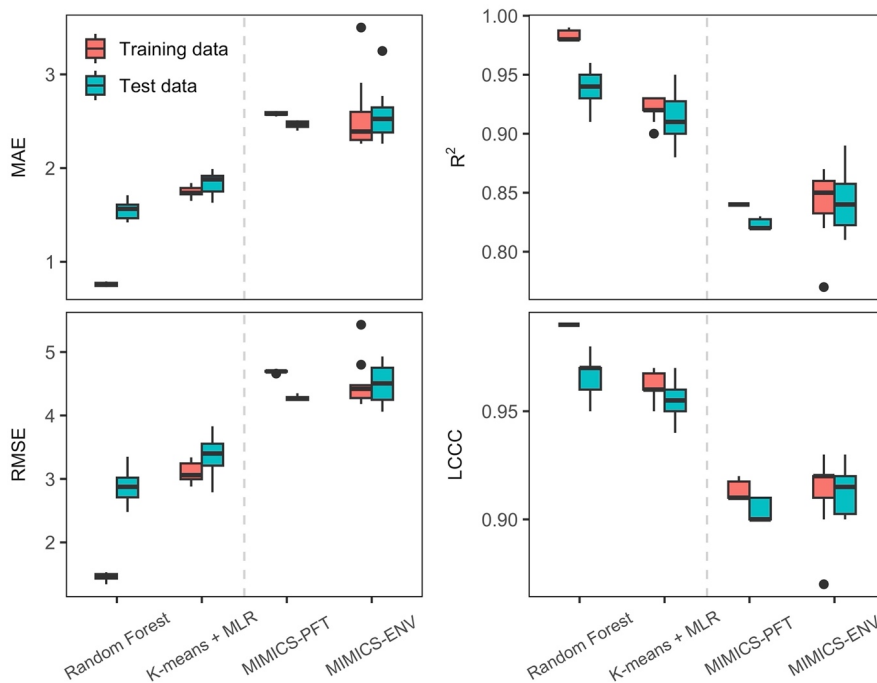


398
 399 **Figure 4.** a) Fraction of different PFT in each cluster divided based on environmental factors; b) spatial
 400 distribution of SOC observations from different environmental clusters and c) density plot of observed SOC
 401 concentration for different clusters.

402
 403 **3.3. Evaluation of model performance**

404
 405 All models employed in this study (RF, K-means + MLR, MIMICS-PFT and MIMICS-ENV)
 406 predicted SOC concentration well for both training data and test data (Figure 5). As anticipated,
 407 sample data versus in-sample training or calibration data. When using test data, the mean value
 408 of R^2 for all models ranges from 0.82 to 0.94, mean LCCC ranges from 0.90 to 0.97, mean
 409 RMSE ranges from 2.88 to 4.51 g C/kg soil, and mean MAE ranges from 1.55 to 2.57 g C/kg
 410 soil.

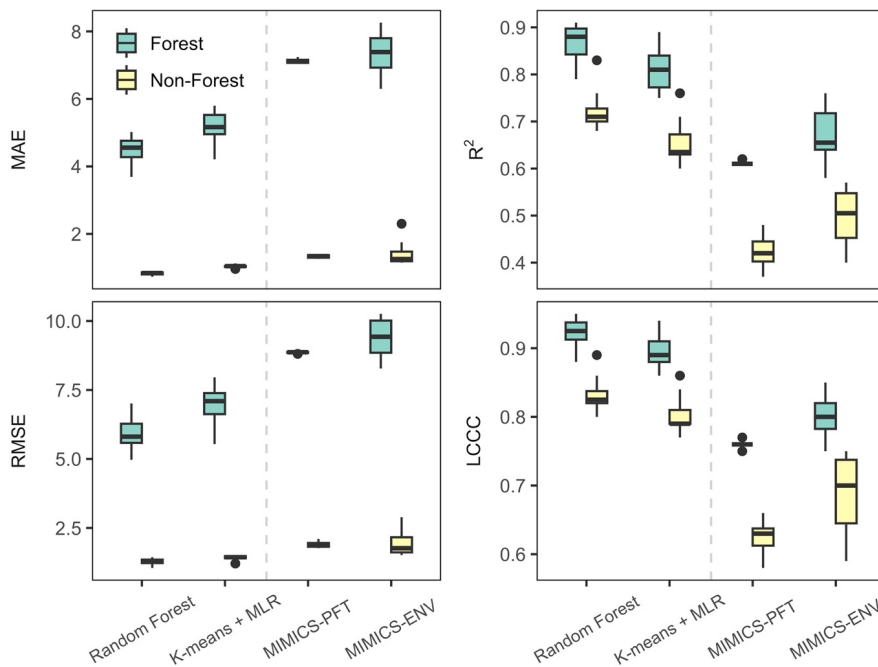
411
 412 The machine learning models outperformed MIMICS in predicting SOC concentration,
 413 regardless of the optimisation approach taken. Particularly, the RF model demonstrated the most
 414 accurate predictions characterized by higher R^2 and LCCC values and lower RMSE and MAE
 415 values for both training and test data. While MIMICS-ENV displayed performance similar to
 416 that of MIMICS-PFT in SOC concentration predictions based on RMSE and MAE, the former
 417 exhibited slightly superior median R^2 and LCCC values but with a higher variability (Figure 5).



418
 419 **Figure 5.** Performance metrics of SOC concentration predictions. Units for MAE and RMSE are g C/kg soil.
 420 Centre line represents median value, and upper and lower box boundaries represent third and first quartile of
 421 metrics from cross-validation. Whiskers extend to the smallest and largest values within 1.5 times the
 422 interquartile range.

423
 424 SOC concentration in forest soil exhibited significantly higher predictability than those in non-
 425 forest (woodland, shrubland and grassland) soil, evidenced by higher R^2 (ranging from 0.58 to
 426 0.91) and LCCC (ranging from 0.75 to 0.95) for test data (Figure 6). Machine learning models
 427 surpassed MIMICS in predicting SOC concentration for both forest and non-forest soils.
 428 Notably, MIMICS-ENV outperformed MIMICS-PFT in SOC concentration predictions,
 429 particularly in non-forest soils.

430
 431



432
 433 **Figure 6.** Performance metrics of SOC concentration predictions for forest and non-forest (woodland,
 435 shrubland and grassland) soils in test (out-of-sample) data. Unit for MAE and RMSE is g C/kg soil. Centre
 436 line represents median value, and upper and lower box boundaries represent third and first quartile of metrics
 437 from cross-validation. Whiskers extend to the smallest and largest values within 1.5 times the interquartile
 438 range.

439
 440 **3.4. Estimations of terrestrial SOC stocks**

441
 442 Using the best fitted models after cross-validation (see Section 2.6 for details), we estimated
 443 the total amount of SOC stocks in the top 30 cm for the whole Australia continent at a spatial
 444 resolution of 0.05° by 0.05°. The optimized parameters used for MIMICS-PFT and MIMICS-
 445 ENV at continental scale are shown in Table 3.

446
 447
 448
 449
 450
 451
 452
 453
 454

455 Table 3. Optimized parameter ranges of MIMICS for cross-validation. Values in brackets were used for
 456 estimating SOC stocks at continental scale. See Table 1 for further explanations of each parameter.

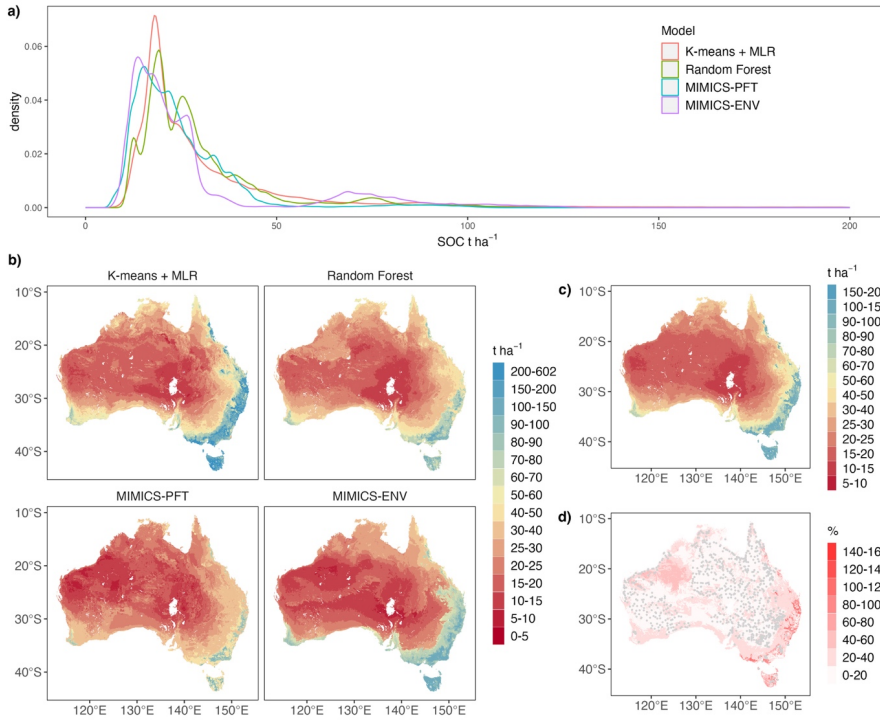
Model	PFT/Cluster	a_v	a_k	xdesorp	xbeta	xdiffsoc
MIMICS- PFT	Grassland	4.36-18.11 (5.45)	4.42-19.11 (5.62)	1.90-3.0 (2.97)	1.06-1.42 (1.06)	16.21-29.90 (29.3)
	Shrubland	12.15-17.91 (12.46)	14.46-18.87 (16.80)	1.54-2.92 (2.58)	1.14-1.27 (1.24)	20.21-29.96 (29.73)
	Woodland	8.41-17.01 (10.92)	9.35-16.99 (12.73)	1.12-1.23 (1.10)	1.12-1.23 (1.18)	20.17-29.96 (23.91)
	Forest	3.15-8.56 (4.70)	12.61-19.69 (13.53)	0.39-3.0 (1.36)	1.42-1.88 (1.35)	11.55-27.70 (10.20)
MIMICS- ENV	Cluster 1	5.23-13.82 (10.189)	6.08-17.80 (11.93)	1.62-2.85 (1.84)	1.07-1.20 (1.07)	0.00-29.81 (28.80)
	Cluster 2	3.56-10.76 (7.60)	7.36-18.24 (15.70)	1.01-2.94 (2.07)	1.05-1.07 (1.05)	3.61-12.75 (6.91)
	Cluster 3	8.31-10.52 (8.48)	15.98-19.91 (19.66)	1.84-2.83 (2.25)	1.36-1.52 (1.52)	10.83-29.45 (26.25)
	Cluster 4	2.47-5.52 (5.10)	6.44-16.80 (13.52)	0.54-1.78 (0.92)	1.21-1.74 (1.42)	14.75-28.91 (20.37)
	Cluster 5	12.24-20.57 (19.55)	10.90-17.56 (17.56)	2.89-3.0 (2.98)	1.05-1.06 (1.05)	25.32-29.83 (25.75)
	Cluster 6	3.25-7.18 (6.40)	7.73-18.23 (15.86)	1.91-2.97 (2.73)	1.05-1.09 (1.09)	6.19-28.57 (15.47)

457
 458 Descriptive statistics of predicted terrestrial SOC stocks at 0-30 cm soil depth are shown in
 459 Table 4. Forests have the largest mean SOC stocks ranging from 70.3 to 113.9 t ha⁻¹ by all
 460 models, and shrubland is estimated to have the lowest mean SOC stocks. The distributions of
 461 predicted continental SOC stocks by all models are positively skewed with most estimated SOC
 462 stocks less than 50 t ha⁻¹ (Figure 7a), and SOC stocks at peak density predicted by MIMICS-
 463 ENV and MIMICS-PFT are smaller than those predicted by the two machine learning
 464 approaches.

465
 466 As expected, all models consistently projected larger SOC stocks in the southeast region,
 467 southwest corner and Tasmania, and consistently indicated lower SOC stocks in central and
 468 western Australia (Figure 7b). Among the models, K-means coupled with multiple linear
 469 regression consistently provided the highest SOC estimations across all vegetation types, while
 470 MIMICS-PFT model consistently yielded the lowest mean SOC stocks.

471
 472 The ensemble estimate of SOC stocks (Figure 7c) shows a similar distribution pattern as that
 473 generated by single model. SOC stocks of the ensemble range from 10.0 to 180.4 t ha⁻¹ with an
 474 average value of 30.3 t ha⁻¹. Coefficient of variation calculated as the ratio of standard deviation
 475 to mean across the four estimates (Figure 7d) is positively correlated with the ensemble mean
 476 estimate. That is, soils with higher SOC stocks exhibit greater variability in SOC predictions
 477 among different models. Note also that the variability of estimates tends to be smaller in areas
 478 with denser numbers of observations (Figure 7d).

479



480

481

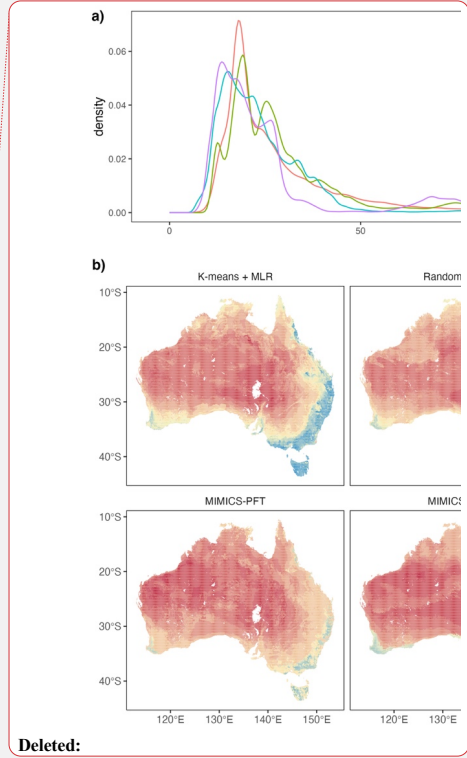
482 **Figure 7.** Estimated Australian terrestrial SOC stocks ($t\ ha^{-1}$) for top 30 cm soil and ensemble statistical
 483 characteristics: a) density plot of estimated terrestrial SOC stocks by all models, noting that only stocks less
 484 than $200\ t\ ha^{-1}$ are shown for better comparison of the distribution; b) estimated SOC stocks by each model;
 485 c) estimated SOC stocks of the ensemble; d) coefficient of variation of the ensemble estimates of SOC stocks.
 486 Grey points represent locations of SOC observations.

487

488 **Table 4.** Descriptive statistics of estimated terrestrial SOC stocks ($t\ ha^{-1}$) at 0-30 cm soil. Min. and Max. are
 489 minimum and maximum value, respectively. 1st Qu and 3rd Qu are first and third quartile, respectively.

490

	PFT	Min.	1 st Qu	median	mean	3 rd Qu	Max.
K-means + MLR	grassland	4.2	17.9	21.2	41.5	42.5	601.1
	shrubland	7.2	16.4	19.3	23.6	24.4	472.2
	woodland	7.1	20.1	26.1	33.3	33.7	483.1
	forest	18.0	51.3	95.2	113.9	153.4	474.0
	all	4.2	18.1	23.6	38.2	36.7	601.1
Random Forest	grassland	10.4	18.5	26.0	30.4	37.2	125.3
	shrubland	10.3	17.0	19.6	21.4	24.4	104.4
	woodland	10.5	20.3	25.8	28.2	32.4	122.1
	forest	29.3	55.0	82.3	78.4	97.0	161.7
	all	10.3	18.9	25.0	29.8	33.7	161.7
	grassland	10.8	16.4	24.1	25.1	33.3	58.7



MIMICS- PFT	shrubland	6.5	12.2	15.5	16.5	20.6	56.5
	woodland	7.8	17.4	21.2	22.1	25.9	61.4
	forest	17.9	44.5	77.4	70.3	88.5	109.9
	all	6.5	15.7	21.2	24.3	28.9	109.9
MIMICS- ENV	grassland	6.8	13.7	18.7	29.9	27.6	124.0
	shrubland	6.7	13.4	16.7	18.3	20.2	131.9
	woodland	8.1	18.0	24.0	27.5	28.0	131.6
	forest	15.8	35.7	90.4	79.4	106.5	134.1
	all	6.7	15.0	20.2	28.9	27.5	134.1
Ensemble	grassland	11.4	17.1	21.1	31.7	36.3	180.4
	shrubland	10.0	15.2	17.3	20.0	21.7	170.4
	woodland	11.0	18.8	24.4	27.8	30.0	168.0
	forest	22.0	46.8	93.1	85.5	112.7	166.3
	all	10.0	17.2	22.2	30.3	31.5	180.4

492

493 4. Discussion

494 4.1. Relative importance of predictors on SOC variations

495

496 Extensive research has been conducted to discern the factors that govern SOC
497 concentration/stocks. Among the commonly employed predictors for SOC spatial variations,
498 climate, organisms, topography, parent material, and soil properties are prominent (Wiesmeier
499 et al., 2019). Within this study, we conducted a comparative assessment of the significance of
500 key variables, namely MAT, MAP, NPP, soil clay content and bulk density, in driving variations
501 in SOC in Australia. Although the number of predictors utilized in our approach is fewer than
502 that employed in most digital mapping methodologies, our models show good performance in
503 predicting SOC in Australia (Figure 5 and 6) and its strength lies in the potential for a more
504 direct comparison between empirical and process-based models.

505

506 Consistent with the result by Hobbey et al., (2015) on the soils from eastern Australia, this study
507 identified soil bulk density as an important predictor of SOC concentration at continental scale
508 (Figure 3). However, the relationship is largely interactive between soil bulk density and soil
509 carbon concentration (Murphy, 2015). Higher concentrations of soil organic matter facilitate
510 soil aggregation formation and increase soil porosity, which results in lower bulk density.
511 Meanwhile, a soil with reduced bulk density exhibits higher permeability for water and oxygen,
512 which enhances plant root growth and SOC dynamics. Physically, the bulk density of organic
513 matter is less than 1 g cm^{-3} , much lower than soil mineral solids with a density of 2.66 g cm^{-3} ,
514 therefore lower bulk density soils usually have higher SOC concentration (Marshall and
515 Holmes, 1988).

516

517 Across the Australia continent, MAT emerges as the second most influential factor governing
518 SOC variations, followed by NPP, MAP, and clay content. This sequence of significance
519 diverges from the findings of Walden et al. (2023), where the order of importance was observed
520 as $\text{NPP} > \text{clay content} > \text{MAP} > \text{MAT}$ on a continental scale in Australia. The number of
521 predictors used in their study is much higher than that in our study, which may affect the
522 contribution of given predictors in SOC variation (Guo et al., 2019). This discrepancy might

523 however be attributable to the utilization of observations encompassing both terrestrial and blue
524 carbon ecosystems in their study. Clay emerges as key driver mainly in the groups where aquatic
525 plants (e.g., seagrass, tidal marsh) appeared. The more extensive dataset encompassing the
526 eastern coastline, characterized by greater variability and abundance of NPP input, potentially
527 elevates NPP to a dominant role in influencing SOC variations within their study.

528
529 For SOC in different vegetation types (Figure 3), soil bulk density and MAT are more important
530 than other factors in forest, and all factors except clay content showed similar importance in
531 predicting SOC concentration in grassland. NPP and MAP dominate the SOC variations in
532 woodland and shrubland. Climate conditions as represented by MAT and MAP exert their
533 impact on SOC in all vegetation types. It was proposed that the primary climatic determinant
534 of SOC variation hinges on the primary constraint affecting SOC production and turnover
535 (Hobley et al., 2016). In this study, most shrublands and woodlands are distributed in arid and
536 semi-arid regions characterized by limited precipitation, which leads to water stress in surface
537 soil, limiting plant productivity and reducing soil C input (Hobley et al., 2015). Consequently,
538 MAP and NPP exhibited relatively higher influence on SOC variations in soils under these
539 vegetation types. In contrast, forest SOC observations are mainly distributed in areas with
540 relatively lower temperatures, therefore experience constrained microbial metabolism, leading
541 to reduced decomposition rates and the high accumulation of SOC (Wynn et al., 2006).
542 Consequently, MAT emerges as a key factor influencing SOC variations in forests. Furthermore,
543 it is noteworthy that soil bulk density plays a crucial role in determining SOC distribution within
544 forest, where it is found to be significantly lower compared to other vegetation types. This lower
545 soil bulk density likely improves oxygen availability to soil microbial communities, and
546 facilitates the formation of microaggregates to enhance the preservation of SOC within the soil
547 matrix (Bronick and Lal, 2005). Consequently, it effectively contributes to elevated SOC
548 concentration levels in forested areas.

549
550 PFT is the only categorical predictor for SOC concentration in this study. SOC is mainly derived
551 from plant C input through above-/belowground tissues, and SOC turnover and storage are
552 influenced by plant traits like plant growth rate and chemical and physical composition (De
553 Deyn et al., 2008; Faucon et al., 2017). With shared representation of similar plant traits, PFT
554 is widely used in process-based models (Poulter et al., 2015; Famiglietti et al., 2023). It was
555 found that the vertical distribution of SOC is highly related to PFT due to the different root
556 distribution and above- and belowground allocation (Jobbágy and Jackson, 2000). However,
557 our study is limited by the absence of SOC observations at multiple soil depths, restricting the
558 analysis to the spatial distribution of SOC at 30 cm soil depth. The influence of PFT on SOC
559 concentration at this particular depth appears relatively insignificant (Figure 3), casting doubt
560 on the effectiveness of optimizing parameters of process-based models for individual PFT
561 (Cranko Page et al., 2023). Considering this, employing the top 3 influential abiotic predictors,
562 soil bulk density, MAT, and MAP, we partitioned all observations into six distinct clusters using
563 K-means. It was anticipated that SOC ranges within each cluster would be narrow due to the
564 high similarity of these three predictors within each group. However, the distribution of SOC

565 in clusters 2, 4, and 6 exhibited considerable variability (Figure 4). Given that these clusters are
566 predominantly composed of forests, it becomes apparent that these three abiotic factors alone
567 are insufficient to fully characterize the intricacies of forest SOC concentration. It was found
568 that elevation and evapotranspiration also drive the variation of forest SOC in Australia (Walden
569 et al., 2023), and taking them into account might potentially increase the predictability of forest
570 SOC.

571 4.2. Model evaluation and comparison with other studies

572
573 Although the predictors used for machine learning models are not exactly same as the inputs of
574 MIMICS, the missing factors (e.g., MAP) were used for parameter optimization of MIMICS-
575 ENV, making the predictions dependent on similar information and so comparable to some
576 extent. Besides, our study presented clear evaluation metrics for out-of-sample validation,
577 enabling a more robust assessment of model performance when applied to new datasets.

578
579 Based on the performance metrics of test data, the machine learning models performed
580 remarkably well (Figure 5). The R^2 suggested that both machine learning models can explain
581 more than 90% of SOC variability across sites, and random forest did the best job with greatest
582 R^2 and LCCC, and lowest MAE and RMSE. Random forest algorithms were widely adopted in
583 predicting spatial-temporal SOC dynamics and produced moderately good performance
584 regionally and globally. For example, Wang et al. (2022) applied random forest to estimate SOC
585 stocks in south-eastern Australia and explained 69% of the variation of current SOC stocks.
586 Nyaupane et al. (2023) trained a random forest model using global SOC observations and
587 explained 61% of SOC variation. The good performance of random forest might be attributed
588 to reduced susceptibility to over-fitting and better capacity to manage the hierarchical non-
589 linear relationships that exist between SOC and environmental predictors (Wang et al., 2018b).
590 Other machine learning methods have been applied to predict continental SOC stocks in
591 Australia. For example, Walden et al. (2023) trained regression-tree algorithm CUBIST to
592 predict SOC stocks for top 30 cm soil using the harmonised datasets. The mean LCCC and
593 RMSE for out-of-sample validation in their study was 0.78 and 0.20 respectively when \log_{10}
594 transformed SOC (t ha^{-1}) values were used. Wadoux et al. (2023) applied quantile regression
595 forest to predict SOC stocks at multiple soil depths. The prediction accuracy decreased
596 dramatically for deeper depth intervals with the greatest R^2 (0.53) at 0-5 cm soil. The better
597 results in this study may be attributed to the removal of cropland ecosystems, which are clearly
598 highly managed and so less predictable. Agricultural practices greatly affect SOC stocks in
599 Australia and add the complexity to the relationship between SOC and environmental factors
600 (Luo et al., 2010). Models using environmental predictors without representation of land use
601 management are unlikely to be able to fully capture the SOC dynamics in croplands (Abramoff
602 et al., 2022).

603
604 Although MIMICS was not as accurate as machine learning models in simulating spatial
605 variation of SOC concentration in Australia, it did well at continental scale with mean R^2 at
606 0.82 and 0.84 for MIMICS-PFT and MIMICS-ENV, respectively (Figure 5), much greater than

607 the values (<0.4) obtained by Abramoff et al. (2022) who applied a different microbial explicit
608 model to Australian SOC dataset. Georgiou et al. (2021) found that there was a mismatch
609 between observations and MIMICS in the role of different environmental controls on SOC
610 variability at global scale. In their study, NPP and MAT had the most explanatory power for
611 SOC stocks from MIMICS, while clay content had the most explanatory power for global SOC
612 observations, which limits the predictability of SOC using MIMICS in their study. However, in
613 our study, NPP and MAT rather than clay content played a greater role in observed SOC
614 variations, perhaps contributing to a better performance of MIMICS in Australia. It also means
615 that SOC estimates in our study are highly sensitive to the estimates of NPP. In this study, we
616 used MODIS NPP product (Running and Zhao, 2021) and did not account for the loss of NPP
617 due to human activities, which may likely influence the optimized estimates of some model
618 parameters, and the uncertainties of simulated SOC concentration. Future studies would ideally
619 use multiple NPP products to quantify the impacts of NPP uncertainties in simulating SOC
620 variation in Australia.

621
622 The modest performance of process-based model MIMICS relative to machine learning models
623 could potentially be attributed to the absence of explicit representation of MAP. The
624 augmentation of MAP within parameter optimization in MIMICS-ENV did allow improved
625 performance compared to MIMICS-PFT, particularly within non-forest regions where the
626 importance of MAP rivals or surpasses that of temperature. Precipitation is a determinant of
627 plant productivity, especially in arid and semi-arid regions. Besides, arid regions with limited
628 precipitation are characterized by lower weathering rate limiting the formation of mineral-
629 associated soil carbon (Doetterl et al., 2015). Hence, we assume that introducing the effect of
630 moisture to MIMICS could contribute to more accurate prediction of SOC, as compared with
631 just taking MAP into account for parametrization, especially in arid and semiarid regions.

632
633 All models produced lower MAE and RMSE for non-forest SOC but higher R^2 and LCCC for
634 forest SOC (Figure 6). SOC in forest is more abundant and variable compared to SOC in other
635 vegetation types even when climate conditions are similar, which leads to greater absolute error
636 in the estimated forest SOC than in other vegetation types. However, in terms of the consistency
637 and concordance between the pattern of observations and predictions, all models show higher
638 ability to predict SOC in forest. Forests, given that they are less perturbed ecosystems, might
639 show greater SOC predictability due to the reduced influence of direct anthropogenic
640 disturbances. Grasslands, shrublands, and woodlands, predominantly situated in Australian
641 rangelands may experience extensive grazing and land management. Primarily, grazing reduces
642 soil carbon input by consumption of aboveground biomass and accelerate SOC decomposition
643 through input of nutrient-enriched animal waste. This introduces additional uncertainties to our
644 modelled SOC estimates, since C input is represented solely by NPP without accounting for the
645 impact of grazing and land managements. Moreover, the cascading effects of grazing extend to
646 potential alterations in plant composition and structural attributes, inducing consequential shifts
647 in litter properties that modulate soil carbon decomposition kinetics (Lunt et al., 2007; Bai and
648 Cotrufo, 2022). The disturbances triggered by grazing manifest in soil carbon pools, leading to

649 a state of disequilibrium rather than adhering to the assumption of SOC convergence toward
650 equilibrium, as embraced in this study's framework. Notably, forests, as relatively undisturbed
651 natural ecosystems, demonstrate a better coherence with the equilibrium assumption, rendering
652 their SOC more amenable to prediction through environmental drivers.

653

654 4.3. Spatial prediction of SOC stocks in Australia

655

656 We produced gridded SOC stocks across Australia using the models validated in this study and
657 an ensemble estimate as the average of four models (Figure 7). Among the models, K-means
658 coupled with multiple linear regression produced the largest mean SOC stocks both at
659 continental scale and for all vegetation types. In contrast, RF and MIMICS with different
660 parameterization approaches produced lower SOC stock estimations (Table 4). The mean
661 terrestrial SOC stocks estimated by random forest and MIMICS are comparable with that
662 estimated by Australian baseline map, which was generated using machine learning algorithm,
663 reporting mean SOC stocks at 29.7 t ha⁻¹ with 95% confidence limits of 22.6 and 37.9 t ha⁻¹
664 (Viscarra Rossel et al., 2014). However, SOC stocks might be underestimated by these methods
665 because of the scarcity of data from the most productive temperate forest both in the baseline
666 map (Bennett et al., 2020) and in our study. Parameter optimization process of MIMICS and
667 the training process of random forest are greatly affected by data used to train the model. Most
668 SOC observations in this study were sourced from arid and semiarid regions, characterized by
669 relatively low SOC content. As a result, the models' ability to predict SOC stocks beyond the
670 observed data range is somewhat constrained. PFT was found to be less important than other
671 environmental factors in driving spatial SOC variations (Figure 3), so it was perhaps not
672 surprising that applying parameters optimized for each plant functional type to the regions with
673 same PFT but broader climate conditions led to inferior results than applying parameters
674 optimized for each environmental group.

675

676 The utilization of linear regression in K-means + MLR generated SOC estimates beyond the
677 range of observations, particularly in eastern Australia where environmental conditions deviate
678 from the training data. The mean SOC stocks estimated by K-means + MLR (38.2 t ha⁻¹) are
679 higher than those of the other models employed in this study, and align closely with the mean
680 value 36.2 t ha⁻¹ reported by Walden et al. (2023) who updated the Australian baseline SOC
681 map (Viscarra Rossel et al., 2014) by incorporating additional SOC observations from forests
682 and coastal marine ecosystems. However, caution is required when interpreting extreme values
683 derived from the K-means + MLR, such as the instance of grassland SOC stocks reaching 601
684 t ha⁻¹ (Table 4). These values raise concerns about the reliability of this approach when
685 extrapolating out-of-sample. Though there is a positive relationship between NPP and SOC
686 observations in this study, SOC accumulation cannot continuously increase linearly in the
687 regions where environmental conditions seem highly conducive to SOC formation. The greater
688 amount of carbon input in eastern Australia might trigger the acceleration of microbial
689 decomposition because of a priming effect, and lead to a decreased accumulation of SOC stocks
690 (Ren et al., 2022). The existence of SOC saturation also implies that SOC cannot be

691 accumulated without limit (Georgiou et al., 2022; Viscarra Rossel et al., 2023). In light of these
692 complexities, applying linear regression to predict SOC stocks, especially under the extreme
693 environmental conditions, should be undertaken with care.

694
695 Continentally, higher SOC stocks were estimated for the southwest corner and southeast
696 Australia (Figure 7), aligning with other SOC maps for Australia (Wadoux et al., 2023; Walden
697 et al., 2023). These regions are characterized by lower temperature and higher precipitation,
698 therefore high SOC accumulation appeared because of high carbon input of NPP and low
699 decomposition rate. However, the high variability of SOC estimates among the four models in
700 these regions should be highlighted (Figure 7d), along with the difference of magnitudes
701 between the estimates in this study and other Australian SOC products (Viscarra Rossel et al.,
702 2014; Walden et al., 2023). Despite inherent differences in model structures, the scarcity of
703 observations in these regions likely contributes to the large uncertainties in SOC estimates.
704 Forest has the largest mean SOC stocks ranging from 70.3 to 113.9 t ha⁻¹ estimated by four
705 models in this study. Around 75% of the forest SOC is from soil under Eucalypt open forest,
706 and mean SOC stocks under this type of forest were estimated to be 87.5 t ha⁻¹ (63.8 -119.6 t
707 ha⁻¹ for 95% confidence interval) (Walden et al., 2023). Shrublands are estimated to have the
708 lowest mean SOC stocks, and more than 90% of shrub SOC observations are from soil under
709 Acacia shrubland and Chenopod shrubland, which rank at the bottom of SOC stocks among
710 different vegetation types (Walden et al., 2023). The low SOC in shrubland is probably due to
711 low carbon input because of limited rainfall (MAP < 280 mm). Though the mean SOC stocks
712 in non-forest regions are much smaller than that for forest, the greater area of vegetation cover
713 results in considerable total SOC stocks, highlighting the importance of carbon building and
714 maintaining via improved managements in these areas. Greater variability of SOC estimates
715 among different models appears in the regions where SOC stocks are higher (Figure 7). The
716 sparsity of SOC observations is a primary contributor to the uncertainties associated with SOC
717 estimates in these regions, highlighting the importance on continual collection of data to better
718 constrain models' behaviour. This imperative is especially pronounced in regions covered by
719 forests, as forested soils exhibit substantial SOC stocks, amplifying the significance of abundant
720 and accurate data acquisition in these specific ecosystems.

721 5. Conclusion

722
723 We compared the performance of two machine learning models, and one process-based
724 microbial model employing two parameterization approaches, to explain the spatial variation
725 of SOC concentration in the top 30 cm soil in Australia. We found that climate conditions and
726 NPP contribute more than soil clay content in predicting SOC concentration in Australia.

727
728 Validation results affirm that with appropriate filtering of data (e.g. removing highly managed
729 crop ecosystems) models can predict SOC concentration at a continental scale with reasonably
730 high reliability, achieving explained variances exceeding 80% for out-of-sample test data, with
731 random forest showing highest prediction accuracy. Notably, all models show higher R² in

732 prediction of SOC in forest than in non-forest soils. MIMICS, with parameters optimized for
733 different environmental clusters, performed better in SOC prediction than MIMICS with
734 parameters optimized for different PFT, especially in non-forest regions.

735
736 All models broadly agree on the spatial distribution of SOC stocks, with higher SOC stocks
737 concentrated in the southeast and southwest regions of Australia. However, the variations in
738 estimated values need to be acknowledged, particularly in highly productive regions. Among
739 these estimates, K-means algorithm coupled with multiple linear regression yields the highest
740 mean SOC stocks estimate, while the MIMICS-PFT model generates the lowest estimate.
741 Considerable disagreement of the maximum and minimum SOC stock values predicted by all
742 models exists partly because models are less constrained by observations in these environments,
743 highlighting the need for continued observational campaigns.

744
745 Our investigation has revealed significant disparities in estimated SOC stocks when different
746 methodologies were employed. This highlights the need for a critical re-evaluation of land
747 management strategies that heavily depend on SOC estimates derived from a single approach.
748 The incorporation of an ensemble of SOC estimates is more likely to effectively capture
749 elements of the uncertainty associated with SOC estimations, providing a more robust basis for
750 informing strategies in soil carbon management and climate change mitigation.

751 Code availability

752 Source Code of vertically resolved MIMICS can be accessed at
753 <https://github.com/Wanglingfei170/MIMICS.git>. Codes for data analysis and machine learning
754 can be accessed by contacting the correspondence author.

Deleted: the CSIRO data portal
<https://doi.org/10.25919/843a-w584> (Wang et al., 2021).

756 Data availability

757
758 The SOC observations described in Viscarra Rossel et al. (2014) are not publicly available but
759 are available from Raphael A. Viscarra Rossel (r.viscarra-rossel@curtin.edu.au) on reasonable
760 request. All other data used in this study are publicly accessible and the specific references are
761 provided in Section 2.4.

Deleted: of these databases

762 Author contribution

763
764 Conceptualization: LW, GA, Y-PW, AP; Methodology: LW, GA, Y-PW; Investigation: LW,
765 RAVR; Formal analysis and Visualization: LW; Writing-original draft preparation: LW;
766 Writing-review & editing: LW, GA, Y-PW, AP, RAVR.

767 Competing interests

768
769 The co-author Raphael A. Viscarra Rossel is a member of the editorial board of SOIL.

773 **Acknowledgements**

774
775 LW thanks the China Scholarship Council and the University of New South Wales for financial
776 support during her PhD study. RAVR and Y-PW thank the Australian Research Council's
777 Discovery Projects scheme (project DP210100420) for funding. LW, GA and AP thank the ARC
778 Centre of Excellence for Climate Extremes for supporting this work (CE170100023).

779
780
781

782 **Reference**

- 783
784 Abramoff, R. Z., Guenet, B., Zhang, H., Georgiou, K., Xu, X., Viscarra Rossel, R. A., Yuan, W.
785 and Ciais, P.: Improved global-scale predictions of soil carbon stocks with Millennial Version
786 2. *Soil Biol Biochem*, 164, 108466, <https://doi.org/10.1016/j.soilbio.2021.108466>, 2022.
- 787 Abs, E. and Ferrière, R.: Modeling microbial dynamics and heterotrophic soil respiration: Effect
788 of climate change. *Biogeochemical cycles: ecological drivers and environmental impact*,
789 103-129, <https://doi.org/10.1002/9781119413332.ch5>, 2020.
- 790 Adhikari, K., Mishra, U., Owens, P., Libohova, Z., Wills, S., Riley, W., Hoffman, F. and Smith,
791 D.: Importance and strength of environmental controllers of soil organic carbon changes with
792 scale. *Geoderma*, 375, 114472, <https://doi.org/10.1016/j.geoderma.2020.114472>, 2020.
- 793 Bai, Y. and Cotrufo, M. F.: Grassland soil carbon sequestration: Current understanding,
794 challenges, and solutions. *Science*, 377, 603-608, doi: 10.1126/science.abo2380, 2022.
- 795 Bennett, L. T., Hinko-Najera, N., Aponte, C., Nitschke, C. R., Fairman, T. A., Fedrigo, M. and
796 Kasel, S.: Refining benchmarks for soil organic carbon in Australia's temperate forests.
797 *Geoderma*, 368, 114246, <https://doi.org/10.1016/j.geoderma.2020.114246>, 2020.
- 798 Bissett, A., Fitzgerald, A., Meintjes, T., Mele, P. M., et al.: Introducing BASE: the biomes of
799 Australian soil environments soil microbial diversity database. *GigaScience*, 5, s13742–
800 016–0126–5. <https://doi.org/10.1186/s13742-016-0126-5>, 2016.
- 801 Bossio, D., Cook-Patton, S., Ellis, P., Fargione, J., Sanderman, J., Smith, P., Wood, S., Zomer,
802 R., Von Unger, M. and Emmer, I.: The role of soil carbon in natural climate solutions. *Nat*
803 *Sustain*, 3, 391-398, <https://doi.org/10.1038/s41893-020-0491-z>, 2020.
- 804 Breiman, L.: Random forests. *Machine learning*, 45, 5-32,
805 <https://doi.org/10.1023/A:1010933404324>, 2001.
- 806 Bronick, C. J. and Lal, R.: Soil structure and management: a review. *Geoderma*, 124, 3-22,
807 <https://doi.org/10.1016/j.geoderma.2004.03.005>, 2005.
- 808 Cranko Page, J., Abramowitz, G., De Kauwe, M. G. and Pitman, A. J.: Are plant functional
809 types fit for purpose? *Geophys Res Lett*, 51, e2023GL104962,
810 <https://doi.org/10.1029/2023GL104962>, 2024.
- 811 Chandel, A. K., Jiang, L. and Luo, Y.: Microbial Models for Simulating Soil Carbon Dynamics:
812 A Review. *J Geophys Res-Bioge*, e2023JG007436, <https://doi.org/10.1029/2023JG007436>,
813 2023.

814 De Deyn, G. B., Cornelissen, J. H. and Bardgett, R. D.: Plant functional traits and soil carbon
815 sequestration in contrasting biomes. *Ecol Lett*, 11, 516-531, [https://doi.org/10.1111/j.1461-](https://doi.org/10.1111/j.1461-0248.2008.01164.x)
816 [0248.2008.01164.x](https://doi.org/10.1111/j.1461-0248.2008.01164.x), 2008.

817 Debeer, D. and Strobl, C.: Conditional permutation importance revisited. *BMC bioinformatics*,
818 21, 1-30, <https://doi.org/10.1186/s12859-020-03622-2>, 2020.

819 Doetterl, S., Stevens, A., Six, J., Merckx, R., Van Oost, K., Casanova Pinto, M., Casanova-
820 Katny, A., Muñoz, C., Boudin, M. and Zagal Venegas, E.: Soil carbon storage controlled by
821 interactions between geochemistry and climate. *Nat Geosci*, 8, 780-783,
822 <https://doi.org/10.1038/ngeo2516>, 2015.

823 Duan, Q., Gupta, V. K. and Sorooshian, S.: Shuffled complex evolution approach for effective
824 and efficient global minimization. *J Optim Theory Appl*, 76: 501-521,
825 <https://doi.org/10.1007/BF00939380>, 1993.

826 Famiglietti, C. A., Worden, M., Quetin, G. R., Smallman, T. L., Dayal, U., Bloom, A. A.,
827 Williams, M. and Konings, A. G.: Global net biome CO₂ exchange predicted comparably
828 well using parameter–environment relationships and plant functional types. *Glob Change*
829 *Biol*, 29, 2256-2273, <https://doi.org/10.1111/gcb.16574>, 2023.

830 Faucon, M.-P., Houben, D. and Lambers, H.: Plant functional traits: soil and ecosystem services.
831 *Trends Plant Sci*, 22, 385-394, <https://doi.org/10.1016/j.tplants.2017.01.005>, 2017.

832 Georgiou, K., Malhotra, A., Wieder, W. R., Ennis, J. H., Hartman, M. D., Sulman, B. N., Berhe,
833 A. A., Grandy, A. S., Kyker-Snowman, E. and Lajtha, K.: Divergent controls of soil organic
834 carbon between observations and process-based models. *Biogeochemistry*, 156, 5-17,
835 <https://doi.org/10.1007/s10533-021-00819-2>, 2021.

836 Georgiou, K., Jackson, R. B., Vindušková, O., Abramoff, R. Z., Ahlström, A., Feng, W., Harden,
837 J. W., Pellegrini, A. F., Polley, H. W. and Soong, J. L.: Global stocks and capacity of mineral-
838 associated soil organic carbon. *Nat Commun*, 13, 3797, [https://doi.org/10.1038/s41467-022-](https://doi.org/10.1038/s41467-022-31540-9)
839 [31540-9](https://doi.org/10.1038/s41467-022-31540-9), 2022.

840 Grace, P. R., Post, W. M. and Hennessy, K.: The potential impact of climate change on
841 Australia's soil organic carbon resources. *Carbon Balance Manag*, 1, 1-10,
842 <https://doi.org/10.1186/1750-0680-1-14>, 2006.

843 Grundy, M., Viscarra Rossel, R. A., Searle, R., Wilson, P., Chen, C. and Gregory, L.: Soil and
844 landscape grid of Australia. *Soil Res*, 53, 835-844, <https://doi.org/10.1071/SR15191>, 2015.

845 Guo, Z., Adhikari, K., Chellasamy, M., Greve, M. B., Owens, P. R. and Greve, M. H.: Selection
846 of terrain attributes and its scale dependency on soil organic carbon prediction. *Geoderma*,
847 340, 303-312, <https://doi.org/10.1016/j.geoderma.2019.01.023>, 2019.

848 Heung, B., Ho, H. C., Zhang, J., Knudby, A., Bulmer, C. E. and Schmidt, M. G.: An overview
849 and comparison of machine-learning techniques for classification purposes in digital soil
850 mapping. *Geoderma*, 265, 62-77, <https://doi.org/10.1016/j.geoderma.2015.11.014>, 2016.

851 Hobley, E., Wilson, B., Wilkie, A., Gray, J. and Koen, T.: Drivers of soil organic carbon storage
852 and vertical distribution in Eastern Australia. *Plant Soil*, 390, 111-127,
853 <https://doi.org/10.1007/s11104-015-2380-1>, 2015.

854 Hobley, E. U., Baldock, J. and Wilson, B.: Environmental and human influences on organic
855 carbon fractions down the soil profile. *Agric Ecosyst Environ*, 223, 152-166,
856 <https://doi.org/10.1016/j.agee.2016.03.004>, 2016.

857 Jackson, R. B., Lajtha, K., Crow, S. E., Hugelius, G., Kramer, M. G. and Piñeiro, G.: The
858 ecology of soil carbon: pools, vulnerabilities, and biotic and abiotic controls. *Annual review*
859 *of ecology, evolution, and systematics*, 48, 419-445, [https://doi.org/10.1146/annurev-](https://doi.org/10.1146/annurev-ecolsys-112414-054234)
860 [ecolsys-112414-054234](https://doi.org/10.1146/annurev-ecolsys-112414-054234), 2017.

861 Jeffrey, S. J., Carter, J. O., Moodie, K. B. and Beswick, A. R.: Using spatial interpolation to
862 construct a comprehensive archive of Australian climate data. *Environ Model Softw*, 16, 309-
863 330, [https://doi.org/10.1016/S1364-8152\(01\)00008-1](https://doi.org/10.1016/S1364-8152(01)00008-1), 2001.

864 Jenny, H.: Factors of soil formation: a system of quantitative pedology, *Agron. J.*, 33, 857-858,
865 <https://doi.org/10.2134/agronj1941.00021962003300090016x>, 1941.

866 Jobbágy, E. G. and Jackson, R. B.: The Vertical Distribution of Soil Organic Carbon and Its
867 Relation to Climate and Vegetation. *Ecol Appl*, 10, 423-436, [https://doi.org/10.1890/1051-](https://doi.org/10.1890/1051-0761(2000)010[0423:TVDOSO]2.0.CO;2)
868 [0761\(2000\)010\[0423:TVDOSO\]2.0.CO;2](https://doi.org/10.1890/1051-0761(2000)010[0423:TVDOSO]2.0.CO;2), 2000.

869 Keskin, H., Grunwald, S. and Harris, W. G.: Digital mapping of soil carbon fractions with
870 machine learning. *Geoderma*, 339, 40-58, <https://doi.org/10.1016/j.geoderma.2018.12.037>,
871 2019.

872 Lamichhane, S., Kumar, L. and Wilson, B.: Digital soil mapping algorithms and covariates for
873 soil organic carbon mapping and their implications: A review. *Geoderma*, 352, 395-413,
874 <https://doi.org/10.1016/j.geoderma.2019.05.031>, 2019.

875 Lawrence, I. and Lin, K.: A concordance correlation coefficient to evaluate reproducibility.
876 *Biometrics*, 45, 255-268, <https://doi.org/10.2307/2532051>, 1989.

877 Le Noë, J., Manzoni, S., Abramoff, R., Bolscher T., Bruni, E., Cardinael, R., Ciais, P., Chenu,
878 C., Clivot, H., Derrien, D., Ferchaud, F., Garnier, P., Goll, D., Lashermes, G., Martin, M.,
879 Rasse, D., Rees, F., Sainte-Marie J., Salmon, E., Schiedung, M., Schimel, J., Wieder, W.,
880 Abiven, S., Barre, P., Cecillon, L. and Guenet, B.: Soil organic carbon models need
881 independent time-series validation for reliable prediction. *Commun Earth Environ*, 4, 158,
882 <https://doi.org/10.1038/s43247-023-00830-5>, 2023.

883 Lee, J., Viscarra Rossel, R. A., Zhang, M., Luo, Z. and Wang, Y. P.: Assessing the response of
884 soil carbon in Australia to changing inputs and climate using a consistent modelling
885 framework. *Biogeosciences*, 18, 5185-5202, <https://doi.org/10.5194/bg-18-5185-2021>, 2021.

886 Lefèvre, C., Rekik, F., Alcantara, V. and Wiese, L.: Soil organic carbon: the hidden potential,
887 Food and Agriculture Organization of the United Nations (FAO), [http://www.fao.org/3/a-](http://www.fao.org/3/a-i6937e.pdf)
888 [i6937e.pdf](http://www.fao.org/3/a-i6937e.pdf), 2017.

889 Lehmann, J. and Kleber, M.: The contentious nature of soil organic matter. *Nature*, 528, 60-68,
890 <https://doi.org/10.1038/nature16069>, 2015.

891 Liang, Z., Chen, S., Yang, Y., Zhou, Y. and Shi, Z.: High-resolution three-dimensional mapping
892 of soil organic carbon in China: Effects of SoilGrids products on national modeling. *Sci Total*
893 *Environ*, 685, 480-489, <https://doi.org/10.1016/j.scitotenv.2019.05.332>, 2019.

894 Lorenz, K., Lal, R. and Ehlers, K.: Soil organic carbon stock as an indicator for monitoring land
895 and soil degradation in relation to United Nations' Sustainable Development Goals. *Land*
896 *Degrad Dev*, 30, 824-838, <https://doi.org/10.1002/ldr.3270>, 2019.

897 Lunt, I. D., Eldridge, D. J., Morgan, J. W. and Witt, G. B.: A framework to predict the effects
898 of livestock grazing and grazing exclusion on conservation values in natural ecosystems in
899 Australia. *Australian Journal of Botany*, 55, 401-415, <https://doi.org/10.1071/BT06178>,
900 2007.

901 Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Chappell, A.,
902 Ciais, P., Davidson, E. A. and Finzi, A.: Toward more realistic projections of soil carbon
903 dynamics by Earth system models. *Global Biogeochem Cycles*, 30, 40-56,
904 <https://doi.org/10.1002/2015GB005239>, 2016.

905 Luo, Z., Wang, E. and Sun, O. J.: Soil carbon change and its responses to agricultural practices
906 in Australian agro-ecosystems: a review and synthesis. *Geoderma*, 155, 211-223,
907 <https://doi.org/10.1016/j.geoderma.2009.12.012>. 2010.

908 Marshall, T. J. and Holmes, J. W.: *Soil physics*, 2nd ed., Cambridge University Press, New York,
909 1988.

910 McBratney, A. B., Santos, M. M. and Minasny, B.: On digital soil mapping. *Geoderma*, 117, 3-
911 52, [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.

912 Minasny, B., McBratney, A. B., Malone, B. P. and Wheeler, I.: Digital mapping of soil carbon.
913 *Advances in agronomy*, 118, 1-47, <https://doi.org/10.1016/B978-0-12-405942-9.00001-3>,
914 2013.

915 Mishra, U. and Riley, W.: Scaling impacts on environmental controls and spatial heterogeneity
916 of soil organic carbon stocks. *Biogeosciences*, 12, 3993-4004, <https://doi.org/10.5194/bg-12-3993-2015>, 2015.

918 Mokany, K., Raison, R. J. and Prokushkin, A. S.: Critical analysis of root: shoot ratios in
919 terrestrial biomes. *Glob Change biol*, 12, 84-96, <https://doi.org/10.1111/j.1365-2486.2005.001043.x>, 2006.

921 Murphy, B. W.: Impact of soil organic matter on soil properties – a review with emphasis on
922 Australian soils. *Soil Research*, 53, 605-635, <https://doi.org/10.1071/SR14246>, 2015.

923 Nyaupane, K., Mishra, U., Tao, F., Yeo, K., Riley, W. J., Hoffman, F. M. and Gautam, S.:
924 Observational benchmarks inform representation of soil organic carbon dynamics in land
925 surface models. *Biogeosci Discuss*, 2023, 1-28, <https://doi.org/10.5194/bg-2023-50>, 2023.

926 Panchal, P., Preece, C., Penuelas, J. and Giri, J.: Soil carbon sequestration by root exudates.
927 *Trends Plant Sci*, 27, 749-757, <https://doi.org/10.1016/j.tplants.2022.04.009>, 2022.

928 Poulter, B., MacBean, N., Hartley, A., Khlystova, I., Arino, O., Betts, R., Bontemps, S.,
929 Boettcher, M., Brockmann, C. and Defourny, P.: Plant functional type classification for earth
930 system models: results from the European Space Agency's Land Cover Climate Change
931 Initiative. *Geosci Model Dev*, 8, 2315-2328, <https://doi.org/10.5194/gmd-8-2315-2015>,
932 2015.

933 Ren, C., Mo, F., Zhou, Z., Bastida, F., Delgado-Baquerizo, M., Wang, J., Zhang, X., Luo, Y.,
934 Griffis, T. J. and Han, X.: The global biogeography of soil priming effect intensity. *Global*
935 *Ecol Biogeogr*, 31, 1679-1687, <https://doi.org/10.1111/geb.13524>, 2022.

936 Rossel, R. V., Chen, C., Grundy, M., Searle, R., Clifford, D. and Campbell, P. The Australian
937 three-dimensional soil grid: Australia's contribution to the GlobalSoilMap project. *Soil Res*,
938 53, 845-864, <https://doi.org/10.1071/SR14366>, 2015.

939 Rumpel, C., Amiraslani, F., Koutika, L.-S., Smith, P., Whitehead, D. and Wollenberg, E.: Put
940 more carbon in soils to meet Paris climate pledges, *Nature*, 564, 32-34,
941 <https://doi.org/10.1038/d41586-018-07587-4>, 2018.

942 Six, J., Conant, R. T., Paul, E. A. and Paustian, K.: Stabilization mechanisms of soil organic
943 matter: Implications for C-saturation of soils. *Plant Soil*, 241, 155-176,
944 <https://doi.org/10.1023/A:1016125726789>, 2002.

945 Smith, P.: Soil carbon sequestration and biochar as negative emission technologies. *Glob
946 Change Biol*, 22, 1315-1324, <https://doi.org/10.1111/gcb.13178>, 2016.

947 Stockmann, U., Padarian, J., McBratney, A., Minasny, B., de Brogniez, D., Montanarella, L.,
948 Hong, S. Y., Rawlins, B. G. and Field, D. J.: Global soil organic carbon assessment. *Glob
949 Food Sec*, 6, 9-16, <https://doi.org/10.1016/j.gfs.2015.07.001>, 2015.

950 Stockmann, U., Adams, M. A., Crawford, J. W., Field, D. J., Henakaarchchi, N., Jenkins, M.,
951 Minasny, B., McBratney, A. B., De Courcelles, V. d. R. and Singh, K.: The knowns, known
952 unknowns and unknowns of sequestration of soil organic carbon. *Agric Ecosyst Environ*,
953 164, 80-99, <https://doi.org/10.1016/j.agee.2012.10.001>, 2013.

954 Terrer, C., Phillips, R. P., Hungate, B. A., Rosende, J., Pett-Ridge, J., Craig, M. E., van
955 Groenigen, K. J., Keenan, T. F., Sulman, B. N., Stocker, B. D., Reich, P. B., Pellegrini, A. F.
956 A., Pendall, E., Zhang, H., Evans, R. D., Carrillo, Y., Fisher, J. B., Van Sundert, K., Vicca, S.
957 and Jackson, R. B.: A trade-off between plant and soil carbon storage under elevated CO₂.
958 *Nature*, 591, 599-603, <https://doi.org/10.1038/s41586-021-03306-8>, 2021.

959 Todd-Brown, K., Randerson, J., Hopkins, F., Arora, V., Hajima, T., Jones, C., Shevliakova, E.,
960 Tjiputra, J., Volodin, E. and Wu, T.: Changes in soil organic carbon storage predicted by Earth
961 system models during the 21st century. *Biogeosciences*, 11, 2341-2356,
962 <https://doi.org/10.5194/bg-11-2341-2014>, 2014.

963 Todd-Brown, K. E., Randerson, J. T., Post, W. M., Hoffman, F. M., Tarnocai, C., Schuur, E. A.
964 and Allison, S. D.: Causes of variation in soil carbon simulations from CMIP5 Earth system
965 models and comparison with observations. *Biogeosciences*, 10, 1717-1736,
966 <https://doi.org/10.5194/bg-10-1717-2013>, 2013.

967 Viscarra Rossel, R. A., Webster, R., Bui, E. N. and Baldock, J. A.: Baseline map of organic
968 carbon in Australian soil to support national carbon accounting and monitoring under climate
969 change. *Glob Change Biol*, 20, 2953-2970, <https://doi.org/10.1111/gcb.12569>, 2014.

970 Viscarra Rossel, R. A., Chen, C., Grundy, M. J., Searle, R., Clifford, D. and Campbell, P. H.:
971 The Australian three-dimensional soil grid: Australia's contribution to the GlobalSoilMap
972 project. *Soil Res*, 53, 845-864, <https://doi.org/10.1071/SR14366>, 2015.

973 Viscarra Rossel, R. A., Lee, J., Behrens, T., Luo, Z., Baldock, J. and Richards, A.: Continental-
974 scale soil carbon composition and vulnerability modulated by regional environmental
975 controls. *Nat Geosci*, 12, 547-552, <https://doi.org/10.1038/s41561-019-0373-z>, 2019.

976 Viscarra Rossel, R. A., Webster, R., Zhang M., Shen, Z., Dixon, K., Wang, Y. P., Walden, L.:
977 How much organic carbon could the soil store? The carbon sequestration potential of
978 Australian soil. *Glob Change Biol*, 30, e17053, <https://doi.org/10.1111/gcb.17053>, 2023.

979 Wadoux, A. M. J., Román Dobarco, M., Malone, B., Minasny, B., McBratney, A. B. and Searle,
980 R.: Baseline high-resolution maps of organic carbon content in Australian soils. *Sci Data*, 10,
981 181, <https://doi.org/10.1038/s41597-023-02056-8>, 2023.

982 Walden, L., Serrano, O., Zhang, M., Shen, Z., Sippo, J. Z., Bennett, L. T., Maher, D. T.,
983 Lovelock, C. E., Macreadie, P. I. and Gorham, C.: Multi-scale mapping of Australia's
984 terrestrial and blue carbon stocks and their continental and bioregional drivers. *Commun*
985 *Earth Environ*, 4, 189, <https://doi.org/10.1038/s43247-023-00838-x>, 2023.

986 Wang, B., Waters, C., Orgill, S., Gray, J., Cowie, A., Clark, A. and Li Liu, D.: High resolution
987 mapping of soil organic carbon stocks using remote sensing variables in the semi-arid
988 rangelands of eastern Australia. *Sci Total Environ*, 630, 367-378,
989 <https://doi.org/10.1016/j.scitotenv.2018.02.204>, 2018a.

990 Wang, B., Gray, J. M., Waters, C. M., Anwar, M. R., Orgill, S. E., Cowie, A. L., Feng, P. and Li
991 Liu, D.: Modelling and mapping soil organic carbon stocks under future climate change in
992 south-eastern Australia. *Geoderma*, 405, 115442,
993 <https://doi.org/10.1016/j.geoderma.2021.115442>, 2022.

994 Wang, B., Waters, C., Orgill, S., Cowie, A., Clark, A., Li Liu, D., Simpson, M., McGowen, I.
995 and Sides, T.: Estimating soil organic carbon stocks using different modelling techniques in
996 the semi-arid rangelands of eastern Australia. *Ecol Indic*, 88, 425-438,
997 <https://doi.org/10.1016/j.ecolind.2018.01.049>, 2018b.

998 Wang, Y. P., Zhang, H., Ciais, P., Goll, D., Huang, Y., Wood, J. D., Ollinger, S. V., Tang, X. and
999 Prescher, A. K.: Microbial activity and root carbon inputs are more important than soil carbon
1000 diffusion in simulating soil carbon profiles. *J Geophys Res Biogeosci*, 126, e2020JG006205,
1001 <https://doi.org/10.1029/2020JG006205>, 2021.

1002 Wieder, W., Grandy, A., Kallenbach, C., Taylor, P. and Bonan, G.: Representing life in the Earth
1003 system with soil microbial functional traits in the MIMICS model. *Geosci Model Dev*, 8,
1004 1789-1808, <https://doi.org/10.5194/gmd-8-1789-2015>, 2015.

1005 Wiesmeier, M., Barthold, F., Spörlein, P., Geuß, U., Hangen, E., Reischl, A., Schilling, B.,
1006 Angst, G., von Lützw, M. and Kögel-Knabner, I.: Estimation of total organic carbon storage
1007 and its driving factors in soils of Bavaria (southeast Germany). *Geoderma Regional*, 1, 67-
1008 78, <https://doi.org/10.1016/j.geodrs.2014.09.001>, 2014.

1009 Wiesmeier, M., Urbanski, L., Hobbey, E., Lang, B., von Lützw, M., Marin-Spiotta, E., van
1010 Wesemael, B., Rabot, E., Ließ, M. and Garcia-Franco, N.: Soil organic carbon storage as a
1011 key function of soils-A review of drivers and indicators at various scales. *Geoderma*, 333,
1012 149-162, <https://doi.org/10.1016/j.geoderma.2018.07.026>, 2019.

1013 Wynn, J. G., Bird, M. I., Vellen, L., Grand-Clement, E., Carter, J. and Berry, S. L.: Continental-
1014 scale measurement of the soil organic carbon pool with climatic, edaphic, and biotic controls.
1015 *Global Biogeochem Cycles*, 20, <https://doi.org/10.1029/2005GB002576>, 2006.

1016 Zhang, H., Goll, D. S., Wang, Y. P., Ciais, P., Wieder, W. R., Abramoff, R., Huang, Y., Guenet,
1017 B., Prescher, A. K. and Viscarra Rossel, R. A.: Microbial dynamics and soil physicochemical

1018 properties explain large-scale variations in soil organic carbon. *Glob Change Biol*, 26, 2668-
1019 2685, <https://doi.org/10.1111/gcb.14994>, 2020.

1020