



1 **Evaluating Vis–NIR spectroscopy for laboratory and in-situ prediction of forest soil**  
2 **organic carbon fractions**

3 Jiaxin Li <sup>a,b,d</sup>, Pierre Roudier <sup>c</sup>, Wenli Fei <sup>a</sup>, Lidu Shen <sup>a</sup>, Yage Liu <sup>a</sup>, Anzhi Wang <sup>a</sup>, Sam McNally <sup>d</sup>, Yuan  
4 Zhang <sup>a,\*</sup>, Jiabing Wu <sup>a,\*</sup>

5

6 <sup>a</sup> CAS Key Laboratory of Forest Ecology and Silviculture, Institute of Applied Ecology, Chinese  
7 Academy of Sciences, Shenyang, China

8 <sup>b</sup> University of Chinese Academy of Sciences, Beijing, China

9 <sup>c</sup> Bioeconomy Science Institute, Manaaki Whenua – Landcare Research, Palmerston North, New  
10 Zealand

11 <sup>d</sup> Bioeconomy Science Institute, Manaaki Whenua – Landcare Research group, Lincoln, New Zealand

12 \* *Corresponding author*: E-mail address: zhangyuan@iae.ac.cn; wujb@iae.ac.cn

13

14 **Abstract**

15 Forest soil organic carbon (SOC) stability is influenced by its relative composition of particulate organic  
16 carbon (POC) and mineral-associated organic carbon (MAOC) fractions. However, conventional SOC  
17 fractionation methods are labor-intensive and restrict large-scale monitoring of SOC dynamics. Visible–  
18 near infrared (Vis–NIR) spectroscopy offers a rapid alternative, yet its applicability for predicting SOC  
19 fractions in forest soils under field conditions remains poorly understood. This study developed an  
20 integrated framework to evaluate the feasibility of in-situ Vis–NIR spectroscopy for predicting SOC  
21 fractions by comparing four in-situ application workflows, including direct laboratory-to-field transfer,  
22 EPO-assisted transfer, direct in-situ modeling, and EPO-assisted in-situ modeling. Direct transfer of  
23 laboratory models to in-situ spectra resulted in substantial performance degradation due to moisture–  
24 driven spectral domain shifts (POC:  $R^2 = 0.80$ ; MAOC:  $R^2 = 0.59$ ). In contrast, direct in-situ modeling.



25 The highest accuracy for POC was achieved using EPO-corrected in-situ spectra ( $R^2 = 0.90$ ), whereas  
26 MAOC prediction performed best using uncorrected in-situ spectra ( $R^2 = 0.71$ ). Independent cross-year  
27 validation further demonstrated that environmental variability, particularly soil moisture, constrained  
28 model robustness. The analysis of the fitted models revealed distinct spectral mechanisms controlling  
29 SOC fraction predictions, linking POC to shortwave infrared organo–clay absorption features (~2200  
30 nm) and MAOC to visible wavelengths associated with iron oxides. These findings highlight the  
31 conditional feasibility of in-situ Vis–NIR spectroscopy for forest SOC fraction prediction and guide field-  
32 based soil carbon monitoring.

33

34 **Keywords:** Spectral transferability; External parameter orthogonalization; Particulate organic carbon;  
35 Mineral-associated organic carbon; In-situ spectroscopy; Model robustness

36

### 37 **Introduction**

38 Soil organic carbon (SOC) plays a central role in regulating the global carbon cycle, and forest soils  
39 constitute one of the largest terrestrial carbon reservoirs. However, the stability and persistence of SOC  
40 are governed not only by its total stock but also by the behavior of its constituent fractions. Increasing  
41 evidence indicates that the response of SOC to environmental change is strongly controlled by the relative  
42 proportions of particulate organic carbon (POC) and mineral-associated organic carbon (MAOC), which  
43 differ fundamentally in their formation pathways, residence time, persistence, and ecological functions  
44 (Lavalée et al., 2020; Zhou et al., 2024). POC represents a relatively labile carbon pool derived primarily



45 from partially decomposed plant residues and microbial products, whereas MAOC forms a more stable  
46 carbon pool through associations with mineral surfaces. Consequently, the distribution of these fractions  
47 provides key insights into soil carbon stability, vulnerability to loss, and long-term sequestration potential  
48 (Abramoff et al., 2022; Viscarra Rossel et al., 2024).

49 Despite their importance, quantifying SOC fractions remains challenging. Conventional approaches  
50 rely mainly on physical, chemical, or thermal fractionation techniques (Delahaie et al., 2024), with  
51 physical fractionation being the most widely used (Wu et al., 2025). Although conceptually robust, these  
52 methods are labor-intensive, time-consuming, and costly, limiting their applicability for large-scale  
53 monitoring or high-frequency assessments of soil carbon dynamics (De Vos et al., 2015). Developing  
54 rapid and scalable approaches to estimate SOC fractions is therefore essential to advancing soil carbon  
55 research and monitoring.

56 Diffuse reflectance spectroscopy in the visible–near infrared range (Vis–NIR, 350 to 2,500 nm)  
57 provides a promising alternative for rapid soil analysis. Because soil organic matter and minerals exhibit  
58 characteristic absorption features across the Vis–NIR region, this technique enables the simultaneous  
59 estimation of multiple soil properties (R. A. Viscarra Rossel et al., 2006; Rossel and Behrens, 2010).  
60 Over the past two decades, laboratory-based Vis–NIR spectroscopy has demonstrated strong predictive  
61 capability for bulk SOC and other soil properties (Jaconi et al., 2017). However, the application of Vis–  
62 NIR spectroscopy to SOC fractions remains comparatively limited. Existing studies have predominantly  
63 relied on mid-infrared spectroscopy to capture fundamental vibrational features of organic functional  
64 groups and minerals (Janik et al., 2007; Pacini et al., 2025), whereas relatively few have explored the  
65 potential of Vis–NIR spectroscopy to predict SOC fractions (Qi et al., 2024).



66 While most current spectral models for SOC fractions were developed using air-dried soils  
67 measured under controlled laboratory conditions, the recent emergence of portable Vis–NIR  
68 spectrometers offers a practical solution for acquiring soil spectra directly in the field (Conforti et al.,  
69 2018; Silvero et al., 2020). In contrast to laboratory spectroscopy, spectra recorded under field conditions  
70 are strongly influenced by environmental factors such as soil moisture, surface roughness, and sensor–  
71 soil contact conditions, which can substantially alter spectral signatures and reduce the transferability of  
72 laboratory-based models (Roudier et al., 2017; Koch et al., 2021). This challenge is even more  
73 pronounced in forest ecosystems when acquiring soil reflectance data via airborne or satellite remote  
74 sensing, as dense canopy cover and litter layers often prevent the spectral signal from reaching the soil  
75 surface. Consequently, the reliability of in-situ Vis–NIR spectroscopy for predicting SOC fractions in  
76 forest soils remains poorly understood.

77 Several methodological strategies have been proposed to address the discrepancy between  
78 laboratory and in-situ spectra. One approach involves directly transferring laboratory-calibrated models  
79 to in-situ spectra, while another involves calibrating models with in-situ measurements. In addition,  
80 spectral correction techniques such as external parameter orthogonalization (EPO) have been proposed  
81 to mitigate environmental interference by removing spectral variance attributable to external factors,  
82 such as soil moisture (Minasny et al., 2011; Roudier et al., 2017). However, the relative effectiveness of  
83 these strategies for predicting SOC fractions remains unclear, particularly in forest soils where  
84 environmental variability is high. Moreover, the spectral mechanisms underlying model predictions are  
85 rarely examined, limiting the ability to determine whether machine learning models capture physically  
86 meaningful soil signals.



87 To address these gaps, this study systematically evaluates the potential of Vis–NIR spectroscopy to  
88 predict forest SOC fractions using both laboratory and in-situ measurements. Specifically, we compare  
89 multiple modeling strategies, including laboratory calibration, laboratory-to-field model transfer,  
90 spectral-correction-assisted transfer, and direct in-situ modeling, under varying environmental conditions.  
91 The specific objectives of this study were to: (1) evaluate the feasibility of predicting forest SOC fractions  
92 (POC and MAOC) using in-situ Vis–NIR spectroscopy; (2) compare the effectiveness of different  
93 modeling strategies for field-based spectral prediction; (3) identify the dominant spectral regions  
94 contributing to model predictions and examine their consistency with known soil mineral–organic  
95 interaction mechanisms; and (4) assess the robustness of spectral models under environmental variability  
96 using independent cross-year validation.

## 97 **2. Materials and methods**

### 98 **2.1. Study area and soil samples**

99 The study was conducted on the northern slope of Changbai Mountain, within the Changbai Mountain  
100 National Nature Reserve (41°41'49"–42°51'18" N, 127°42'55"–128°16'48" E), Jilin Province, Northeast  
101 China (**Fig. 1a, b**). The area lies at an average elevation of 758 m and experiences a typical East Asian  
102 continental monsoon climate, with mean annual temperatures of 2.5–5.5°C, 700–900 mm of precipitation, and  
103 roughly 2300 hours of sunshine per year. Soils are predominantly albic dark brown soils developed on basalt,  
104 characterized by a humus-rich dark surface layer and a grayish albic subsurface horizon.

105 Soil sampling was conducted during two independent field campaigns in 2024 and 2025 across seven  
106 representative forest sites (**Fig. 1c**). The same sampling protocol and sample preparation procedures were  
107 applied in both years. The spatial arrangement of sampling points differed between the campaigns due to



108 field constraints, as described below. In 2024, three forest sites were sampled. At each site, nine  $2 \times 2$  m  
109 quadrats were established, and three sampling points were randomly selected within each quadrat. In 2025,  
110 four forest sites were sampled. At each site, five  $2 \times 2$  m quadrats were established, with four sampling points  
111 randomly selected per quadrat. At each sampling point, two adjacent soil cores were collected and combined  
112 to form a composite sample. Soil samples were taken from four depth intervals (0-10, 10-20, 20-30, and 30-  
113 40 cm), yielding  $\sim 1$  kg per composite sample. A total of 624 soil samples were collected across two sampling  
114 campaigns: 324 in 2024 and 300 in 2025. After removing visible roots and coarse debris, samples were placed  
115 in labeled plastic bags, air-dried, gently crushed, and sieved through a 2 mm mesh prior to laboratory analyses  
116 and spectral measurements.

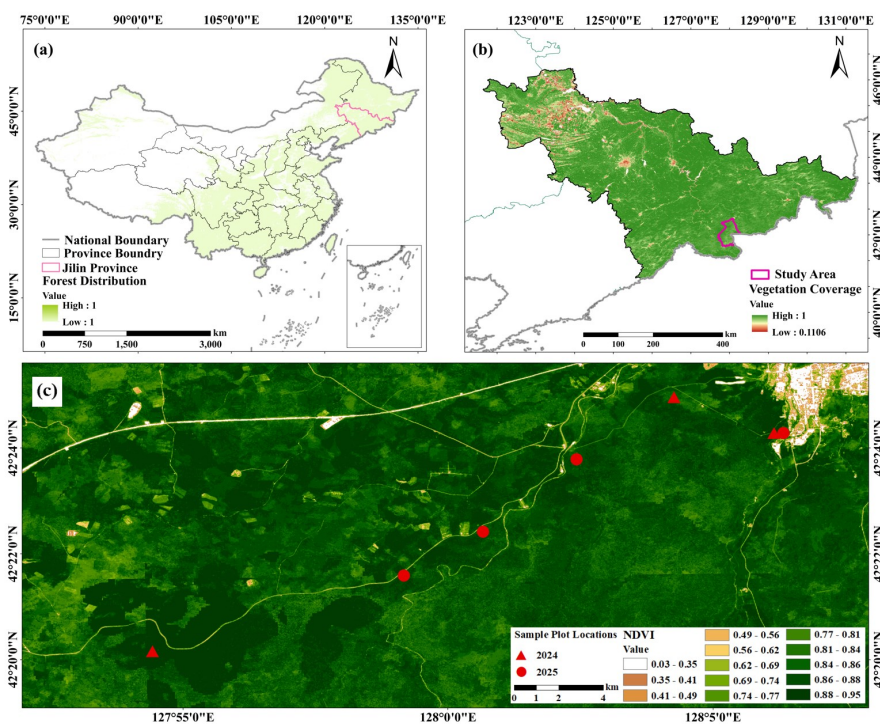


Figure 1. The location of the study area and soil sampling sites.

117 2.2 Laboratory measurement of SOC fractions



118 SOC was fractionated into particulate organic carbon (POC,  $\geq 53 \mu\text{m}$ ) and mineral-associated organic  
119 carbon (MAOC,  $< 53 \mu\text{m}$ ) using a size fractionation methodology similar to that outlined in Cotrufo et al.  
120 (2019). Briefly, 10 g of air-dried soil was mixed with 50 mL of  $5.0 \text{ g L}^{-1}$  sodium hexametaphosphate by  
121 manual shaking for 5 min, then oscillated at 90 rpm for 18 h. The suspension was then washed through a 53  
122  $\mu\text{m}$  sieve until the resulting filtrate was clear. The fraction retained on top of the sieve (e.g., POM,  $> 53 \mu\text{m}$ )  
123 and the fraction which passed through the sieve (e.g., MAOM,  $< 53 \mu\text{m}$ ) were collected in pre-weighed  
124 aluminum containers, dried at  $60^\circ\text{C}$  for at least 48 h. After drying, each fraction was gently ground with a  
125 mechanical grinder and sieved through a 0.15-mm mesh. Material retained on the sieve was re-ground until  
126 complete passage was achieved, ensuring sample homogeneity prior to elemental analysis. The carbon  
127 content for each isolated fraction was determined using an elemental analyzer (Vario MACRO cube).

## 128 **2.3 Spectral data acquisition and analysis**

### 129 **2.3.1 Laboratory and in-situ spectral measurements**

130 Both in-situ and laboratory spectra were acquired using an ASD FieldSpec4 spectrometer  
131 (Analytical Spectral Devices, Inc.) equipped with a high-intensity contact probe, covering the 350-2500  
132 nm range at 1 nm resolution. The spectrometer was calibrated before each measurement using a white  
133 Spectralon® panel.

134 For in-situ measurements, soil cores from the four depth intervals were measured immediately after  
135 extraction. Care was taken to avoid gravel, roots, litter, and other debris. Five replicate spectra per layer  
136 were averaged to reduce surface heterogeneity. This approach preserved vertical soil stratification and  
137 minimized the lag between sampling and measurement.



138 Laboratory spectra were obtained from air-dried, ground, and sieved samples placed in 8 cm  
139 diameter containers with leveled surfaces. For each sample, five replicate spectra were measured and  
140 averaged to produce a single representative spectrum. Consequently, 624 averaged spectra were obtained  
141 for in-situ measurements and 624 for laboratory measurements, resulting in 1248 spectra used for  
142 analysis.

### 143 **2.3.2 Spectral morphology analysis**

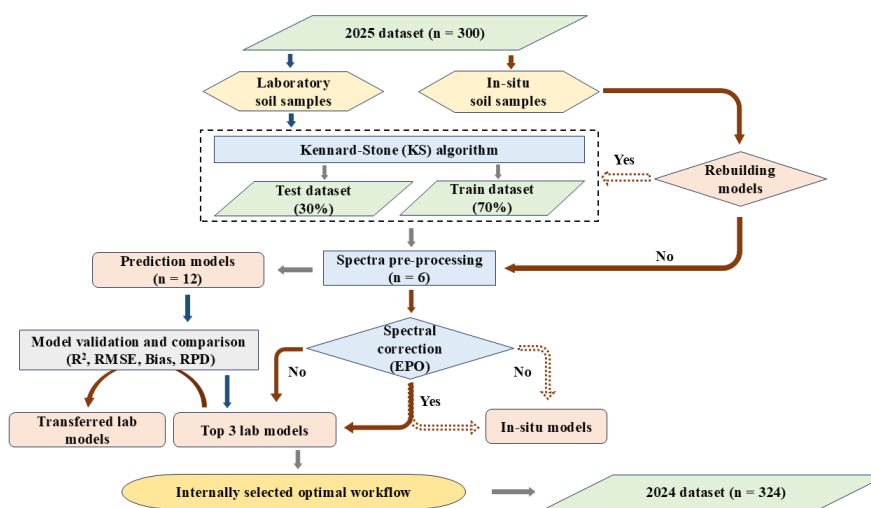
144 To provide a physically grounded interpretation of model outputs and to assess whether statistically  
145 significant wavelengths correspond to meaningful spectral features, spectral morphology was examined  
146 alongside machine-learning results. Spectra were grouped by soil horizon and averaged for each depth.  
147 The resulting depth-specific mean spectra from laboratory and in-situ measurements were compared,  
148 yielding eight representative spectral curves (**Fig. 4a**). The Morphological Interpretation of Reflectance  
149 Spectra (MIRS; Demattê et al., 2014) framework was applied to extract three key attributes: (1) absolute  
150 reflectance, (2) spectral shape (slope and curvature), and (3) absorption features (depth, width, and area).  
151 These morphological descriptors were used to interpret the influential spectral regions identified by the  
152 models and to link them to soil properties.

### 153 **2.4 Model development, transfer, and validation**

154 Model development and evaluation were conducted using a stepwise framework under laboratory  
155 and in-situ conditions. The workflow consisted of: (1) laboratory-based model development and  
156 evaluation, in which six preprocessing methods were combined with twelve machine learning algorithms  
157 (72 combinations) to identify robust approaches; (2) direct transfer of laboratory-calibrated models to in-  
158 situ spectra (Workflow A); (3) transfer of laboratory models to EPO-corrected in-situ spectra (Workflow



159 B); (4) direct in-situ model development using field-acquired spectra (Workflow C); and (5) EPO-  
 160 corrected in-situ model development (Workflow D) (Fig. 2). The laboratory-based modeling step served  
 161 as the baseline for evaluating the transferability and applicability of the four subsequent in-situ  
 162 application workflows. Based on laboratory performance and stability, the three best-performing models  
 163 were selected for subsequent in-situ analyses. Model robustness was further evaluated through  
 164 independent cross-year validation.



**Figure 2.** Flowchart of the modeling framework. Gray arrows: Shared Steps; Red arrows: In-situ-specific Steps (Dashed arrows distinguish laboratory-to-field model transfer from direct in-situ model development); Blue arrows: Lab-specific Steps.

165 **2.4.1 Dataset partitioning strategy**

166 The train and test datasets were partitioned using the Kennard-Stone (KS) algorithm (Kennard and  
 167 Stone, 1969) based on laboratory spectra to ensure a representative coverage of the spectral space. To  
 168 avoid potential data leakage arising from within-site correlations, the KS algorithm was implemented at  
 169 the sampling-location level rather than at the individual-sample level. Specifically, all depth layers



170 corresponding to the same sampling location were assigned collectively to either the calibration or test  
171 set. This grouping strategy ensured strict independence between datasets by preventing samples from the  
172 same site from being split across calibration and validation subsets. This design was adopted to maximize  
173 spectral diversity within each subset and to support the primary objective of this study, which focuses on  
174 assessing the transferability of laboratory-calibrated models to in-situ spectral measurements under  
175 varying environmental conditions, rather than on strict spatial or site-level extrapolation.

#### 176 **2.4.2 Spectral preprocessing**

177 Spectral preprocessing was applied to reduce noise, correct scattering effects, and highlight  
178 chemically relevant signals. Different preprocessing options were considered. The reflectance spectra  
179 (REF) were first truncated to 400-2500 nm to remove the noisier parts of the spectra. Reflectance data  
180 were then converted to absorbance ( $AB = \log_{10}(1/REF)$ ) to improve linearity between spectral response  
181 and soil constituent concentrations. Additional preprocessing included first derivative (FD) and standard  
182 normal variate (SNV), as well as their combinations with AB (AB\_FD, AB\_SNV). FD was calculated  
183 using a simple finite difference approach (i.e., first-order difference between adjacent wavelengths)  
184 without additional smoothing, whereas SNV corrected for scattering and intensity variations (Wang et  
185 al., 2022). The influence of different preprocessing approaches on model stability and transferability,  
186 assessed using internal and external validation sets, was evaluated during model development, and the  
187 strategy yielding the most robust performance was ultimately retained.

#### 188 **2.4.3 Laboratory model development**

189 To assess the predictive performance of different modeling approaches under controlled laboratory  
190 conditions and to select models suitable for in-situ application, twelve widely used spectral algorithms



191 were evaluated using laboratory spectra. These models cover a broad range of linear and nonlinear  
192 approaches commonly applied in soil spectroscopy studies. Linear models included Principal Component  
193 Regression (PCR), Partial Least Squares Regression (PLSR), and Lasso regression. Nonlinear  
194 approaches comprised memory-based learning (MBL), tree-based methods (Cubist, Random Forest,  
195 Gradient Boosting Decision Trees), kernel-based methods (SVMR, GPR), and neural networks (ANN,  
196 RNN, CNN). All models were trained and optimized using five-fold cross-validation on the 2025 training  
197 set. Model training and evaluation were implemented in R (v4.4.3) and Python (v3.9.0) using standard  
198 packages for modeling and machine learning. The detailed descriptions of the models and their  
199 hyperparameters are provided in Supplementary Material S1.

#### 200 **2.4.4 Spectral correction method**

201 In-situ soil spectral measurements are strongly influenced by environmental factors, particularly  
202 soil moisture, which may obscure spectral features associated with intrinsic soil properties and degrade  
203 predictive performance (Roudier et al., 2017). To harmonize measurement conditions between laboratory  
204 and in-situ spectra, the in-situ Vis–NIR spectra were corrected using the EPO method (Minasny et al.,  
205 2011).

206 Paired laboratory and in-situ spectra were constructed by matching samples collected from identical  
207 sampling locations and depth layers. For each paired observation, a difference spectrum (laboratory  
208 minus in-situ) was calculated to characterize systematic measurement-induced variability. The matrix of  
209 difference spectra was subjected to principal component analysis (PCA), and the resulting eigenvectors  
210 were used to define an orthogonal subspace representing external parameter effects. A projection matrix  
211 was subsequently constructed to remove this subspace from the in-situ spectra prior to model transfer or



212 in-situ model calibration.

213 The number of principal components retained ( $c$ ) was determined as the minimum number required  
214 to explain 95% of the variance in the difference spectra. This threshold was selected as a conservative  
215 compromise to remove dominant measurement-related variability while minimizing the risk of  
216 eliminating soil-intrinsic spectral information. The derived projection matrix was consistently applied in  
217 all subsequent analyses.

#### 218 **2.4.5 Transfer of laboratory models to in-situ spectra**

219 To evaluate the transferability of laboratory-calibrated models under field conditions, the three best-  
220 performing models identified in the laboratory evaluation were selected for subsequent in-situ analyses  
221 (**Fig. S4**). The laboratory-calibrated models served as the baseline reference for all transfer analyses. As  
222 Workflow A, laboratory-calibrated models were directly applied to in-situ spectra without recalibration  
223 or spectral correction. This direct transfer reflects the intrinsic ability of spectral–property relationships  
224 learned under controlled laboratory conditions to generalize to field-acquired spectra. As Workflow B,  
225 in-situ spectra were first corrected using EPO (Minasny et al., 2011) before model transfer. The same  
226 EPO projection matrix derived from the training dataset was applied consistently across workflows.

#### 227 **2.4.6 In-situ spectral model development with and without spectral correction**

228 Accordingly, direct in-situ spectral models were developed using field-acquired Vis–NIR spectra  
229 (Workflow C). The same three modeling algorithms identified from the laboratory evaluation were  
230 employed to ensure direct comparability with the transfer workflows. In addition, an EPO-assisted in-  
231 situ modeling workflow (Workflow D) was implemented, in which in-situ spectra were first corrected



232 using the same training-sample-derived EPO projection matrix before model calibration. For both  
233 Workflow C and Workflow D, model calibration, validation, and hyperparameter optimization strictly  
234 followed the procedures described above, ensuring methodological consistency and unbiased comparison  
235 across all workflows.

#### 236 **2.4.7 Cross-year environmental transfer evaluation**

237 The dataset collected in 2025 served as the model development dataset and was used for calibration–  
238 test partitioning, cross-validation, hyperparameter optimization, spectral preprocessing parameterization,  
239 and EPO transformation matrix construction.

240 In contrast, the dataset collected in 2024 was kept fully independent of all model development  
241 procedures and used exclusively as an external validation dataset to evaluate model transferability and  
242 robustness across different environmental conditions. All preprocessing parameters, EPO projection  
243 matrices, and model hyperparameters were determined solely based on the 2025 data and subsequently  
244 fixed prior to prediction on the 2024 dataset.

245 Differences in spectral distributions between the 2024 and 2025 datasets are illustrated in **Fig. S2**,  
246 reflecting variations in soil moisture conditions and measurement environments between sampling  
247 campaigns.

#### 248 **2.5 Evaluation metrics**

249 Model performance was assessed on both internal test sets and independent external validation sets  
250 using  $R^2$ , Root Mean Squared Error (RMSE), Bias, and the ratio of performance to deviation (RPD).  
251 During model development, a five-fold cross-validation (CV) procedure was applied to the training



252 dataset for model tuning and selection, and the optimal models were subsequently evaluated on the  
253 independent test and external validation datasets. The equations are as follows:

254 
$$R^2 = 1 - \frac{\sum_i^n (\hat{y}_i - y_i)^2}{\sum_i^n (y_i - \bar{y})^2} \quad (1)$$

255 
$$RMSE = \sqrt{\frac{\sum_i^n (\hat{y}_i - y_i)^2}{n}} \quad (2)$$

256 
$$Bias = \frac{\sum_i^n (y_i - \hat{y}_i)}{n} \quad (3)$$

257 
$$RPD = \frac{SD}{RMSE} \quad (4)$$

258 where  $n$  is the number of samples,  $y_i$  and  $\hat{y}_i$  are the measured and predicted values, respectively,  $\bar{y}$   
259 is the mean of measured values, and SD is the standard deviation of measured values.  $R^2$  indicates the  
260 proportion of variance explained by the model, RMSE reflects the average prediction error, Bias reveals  
261 systematic over- or underestimation, and RPD evaluates predictive capability. RPD was used as a relative  
262 indicator of predictive performance with respect to data variability and was interpreted in a comparative  
263 rather than categorical manner (R.A. Viscarra Rossel et al., 2006).

## 264 **2.6 Model interpretation analysis**

265 Model interpretation was conducted using SHapley Additive exPlanations (Lundberg and Lee, 2017)  
266 to quantify the contribution of individual spectral bands to model predictions. SHAP has been widely  
267 applied in soil spectroscopy to identify informative wavelengths and to improve the interpretability of  
268 machine learning models (Padarian et al., 2020; Haghi et al., 2021), particularly when combined with  
269 spectral morphology analysis (Viscarra Rossel et al., 2022; Ma et al., 2023). In this research, SHAP was  
270 conducted exclusively for the best-performing model identified for each target variable. Given the  
271 fundamental structural differences between Cubist (global, rule-based modeling) and MBL (locally



272 adaptive modeling), distinct SHAP strategies were implemented to ensure consistency between the  
273 prediction mechanisms and the attribution analysis.

#### 274 **2.6.1 SHAP interpretation of Cubist models**

275 Cubist operates under a single global prediction structure defined by rule-based regression functions.  
276 SHAP values were therefore computed directly from the fitted Cubist models to quantify wavelength-  
277 level contributions under a consistent global mapping between spectral inputs and target variables,  
278 enabling interpretation of global spectral importance patterns across the dataset.

#### 279 **2.6.2 Local SHAP interpretation for MBL models**

280 In contrast to Cubist, MBL generates predictions using locally calibrated models that vary among  
281 target samples. Because standard SHAP implementations assume a fixed global prediction function,  
282 applying global SHAP directly to MBL outputs may yield ambiguous interpretations.

283 To ensure consistency between prediction and attribution, SHAP analysis for MBL was conducted  
284 at the local model level. For each test sample, a local calibration set was constructed by retrieving its k  
285 nearest neighbors in spectral space using Euclidean distance, consistent with the neighborhood definition  
286 and optimal k used in MBL prediction (Wold et al., 1998; Daelemans, 2009) A local PLSR model was  
287 fitted within each neighborhood, and SHAP values were computed relative to this local prediction  
288 function using a model-agnostic framework (Sokol and Flach, 2020). Ten test samples were selected via  
289 stratified sampling across the full target range to balance representativeness and computational feasibility.  
290 SHAP values were calculated separately for each local model and aggregated by averaging absolute  
291 contributions to derive stable spectral importance patterns.



292 **3. Results**

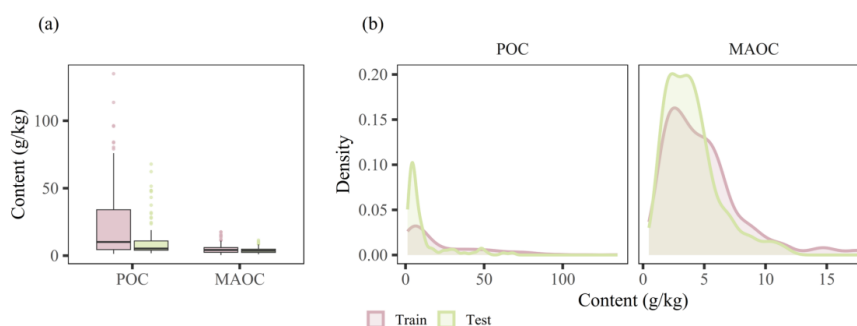
293 **3.1 Forest SOC fractions**

294 The mean contents of POC and MAOC were 19.04 and 4.56 g kg<sup>-1</sup> (**Fig.3**), respectively. Both POC  
295 and MAOC exhibited positively skewed distributions, with skewness values of 1.87 and 1.70,  
296 respectively, indicating the presence of high-value observations. The training and test sets showed  
297 comparable ranges (**Fig. 3**). For POC, values ranged from 1.28 to 134.88 g kg<sup>-1</sup> in the training set and  
298 from 1.76 to 67.89 g kg<sup>-1</sup> in the test set. For MAOC, the corresponding ranges were 0.48–17.62 g kg<sup>-1</sup>  
299 and 1.05–11.47 g kg<sup>-1</sup>, respectively (Fig. 3a). POC exhibited greater variability than MAOC across  
300 samples.

301 The relationships between SOC and its fractions were further examined using Pearson correlation  
302 analysis, showing a strong correlation between SOC and POC ( $r = 0.99$ ) and a moderate correlation  
303 between SOC and MAOC ( $r = 0.72$ ) (**Fig. S3**).

304 The distributions of POC and MAOC in the training and test sets were further compared using  
305 probability density functions (**Fig. 3b**). The density curves showed a high degree of overlap between the  
306 two subsets, indicating consistent distribution patterns and suggesting that the dataset partitioning is  
307 representative with minimal sampling bias.

308



**Figure 3.** (a) Distributions of particulate organic carbon (POC) and mineral-associated organic carbon (MAOC) for the training (210 samples) and test (90 samples) sets; (b) probability density functions showing the distributions of POC and MAOC in the training and test sets.

### 309 3.2 Spectral characteristics of laboratory and in-situ soil measurements

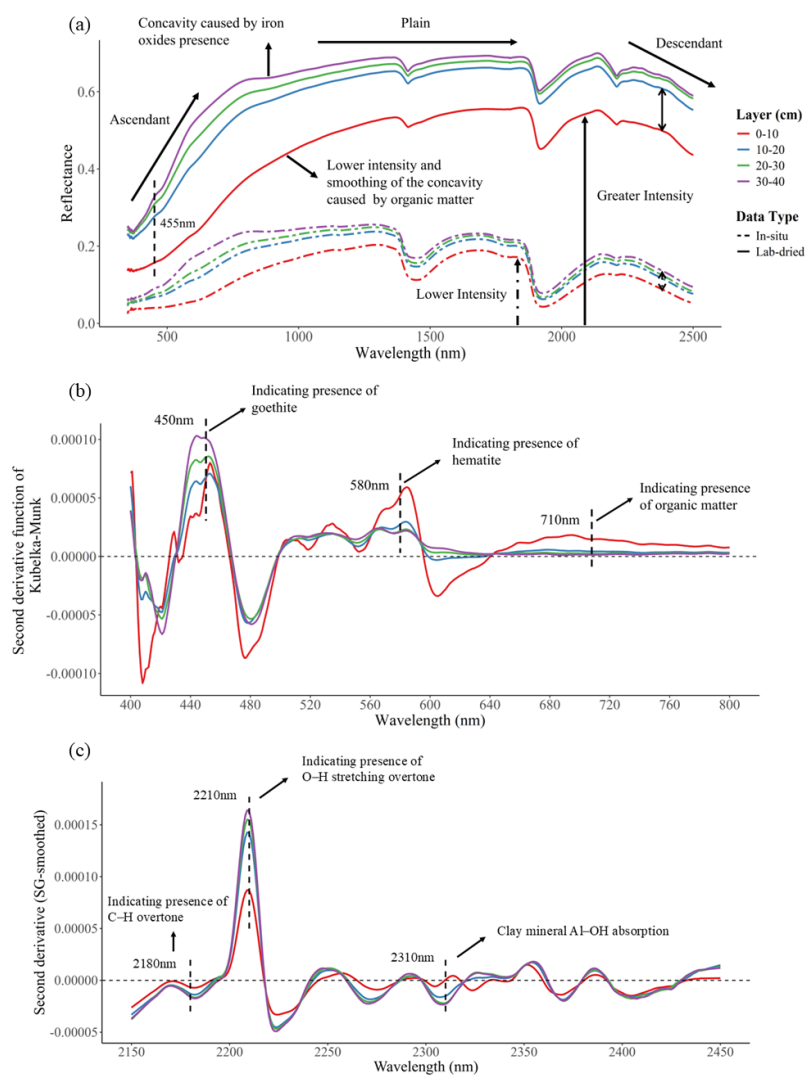
310 The mean reflectance of in-situ soil samples was, on average, 71.45% lower than that of laboratory  
311 air-dried samples across the 400–2500 nm spectral range (**Fig. 4**). Reflectance reduction was particularly  
312 evident near 1400 nm and 1900 nm. Principal component analysis further revealed a clear separation  
313 between laboratory and in-situ spectra (**Fig. S1**), indicating a distinct shift in the spectral domain.

314 Across all soil layers, the raw reflectance spectra (**Fig. 4a**) exhibited a consistent three-stage pattern:  
315 increasing reflectance in the visible region (400–700 nm), a relatively stable plateau in the near-infrared  
316 region (700–1300 nm), and a gradual decline in the shortwave infrared region (1300–2500 nm). Under  
317 laboratory conditions, subsurface layers (10–40 cm) generally showed higher reflectance than the surface  
318 layer (0–10 cm), while the surface layer exhibited lower reflectance and smoother spectral concavity.

319 The second derivative of the Kubelka–Munk function (**Fig. 4b**) revealed distinct spectral features  
320 in the visible region (Dematté et al., 2014). A feature around 580 nm showed higher intensity in the  
321 surface layer, whereas features near 450 nm were more pronounced in subsurface layers. The feature



322 around 710 nm also displayed clear vertical variation, with higher intensity in the surface layer. In the  
323 shortwave infrared region (**Fig. 4c**), several diagnostic features were observed. A prominent feature near  
324 2210 nm showed stronger intensity in subsurface layers, whereas a feature near 2180 nm was more  
325 pronounced in the surface layer. Additional features around 2310 nm also exhibited depth-dependent  
326 variation.



**Figure 4.** Spectral characteristics and feature analysis of soil samples. (a) Comparison of mean spectral reflectance (0–40 cm) between laboratory (solid lines) and in-situ (dashed lines) measurements; (b) Second derivative of Kubelka–Munk transformed laboratory spectra (400–800 nm); (c) Second derivative of Savitzky–Golay (SG) smoothed laboratory spectra (2150–2450 nm).

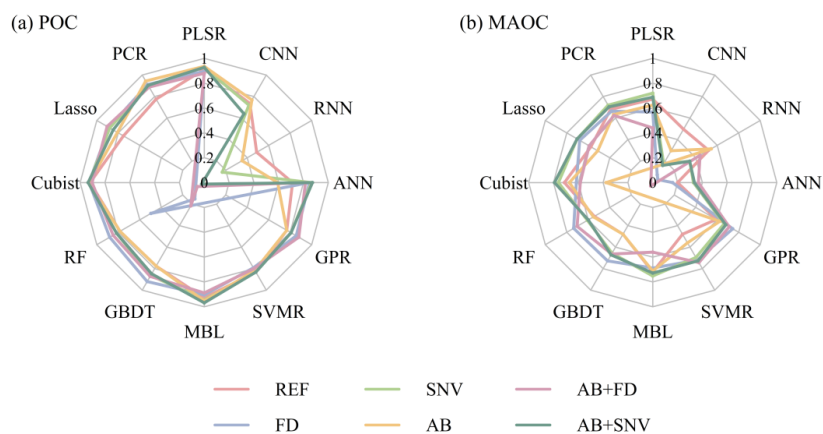
327 **3.3 Laboratory predictions of SOC fractions**

328 Among all evaluated models, PLSR, Cubist, and MBL consistently outperformed the others. For



329 POC prediction,  $R^2$  values ranged from 0.85 to 0.92 for PLSR, 0.88 to 0.92 for Cubist, and 0.86 to 0.96  
 330 for MBL (Fig. 5). In contrast, MAOC prediction exhibited substantially lower accuracy, with  $R^2$  values  
 331 ranging from 0.30 to 0.65 for PLSR, 0.48 to 0.74 for Cubist, and 0.45 to 0.69 for MBL.

332 For both of the carbon fractions, the AB\_SNV preprocessing consistently yielded the highest  
 333 prediction accuracy. Under this unified preprocessing condition, POC achieved the best performance  
 334 with MBL, with predictions yielding an  $R^2$  of 0.96, RMSE of 3.05, Bias of 0.48, and RPD of 4.67. By  
 335 comparison, under the same preprocessing conditions, MAOC prediction achieved its optimal  
 336 performance with Cubist, reaching an  $R^2$  of 0.74, an RMSE of 1.15, a Bias of -0.26, and an RPD of 1.90.



**Figure 5.** Comparison of predictive performance ( $R^2$ ) of twelve machine learning models for particulate organic carbon (POC) and mineral-associated organic carbon (MAOC) based on laboratory Vis-NIR spectra. Radar charts illustrate model performance under different spectral preprocessing methods. The radius along each axis indicates the  $R^2$  coefficient, with higher values indicating better predictive accuracy.

### 337 3.4 In-situ predictions of SOC fractions

338 Model transferability and in-situ applicability for POC and MAOC were assessed by comparing the  
 339 predictive performance of the three best-performing laboratory models (Fig. S4) across the four in-situ



340 application workflows (A–D), including two laboratory-to-field transfer workflows (A–B) and two direct  
341 in-situ modeling workflows (C–D), with the optimal model–preprocessing combinations identified under  
342 each workflow (**Fig. 6**).

343 The direct application of laboratory-calibrated models to in-situ spectra (Workflow A) resulted in  
344 pronounced declines in predictive performance across all carbon fractions. POC showed a substantial  
345 decline in predictive performance following transfer, with  $R^2$  decreasing from 0.85–0.96 under laboratory  
346 conditions to 0.04–0.70. In contrast, MAOC exhibited the weakest transfer performance, with  $R^2$   
347 declining from 0.30–0.74 to mostly below 0.43.

348 Workflow B partially mitigated the decline observed under Workflow A. For POC,  $R^2$  values  
349 increased to 0.23–0.85, respectively, with reductions in RMSE relative to Workflow A. However,  
350 improvements varied across model–preprocessing combinations; for example, under SNV + PLSR,  $R^2$   
351 for POC decreased compared with Workflow A. MAOC prediction showed only marginal gains in  $R^2$ .

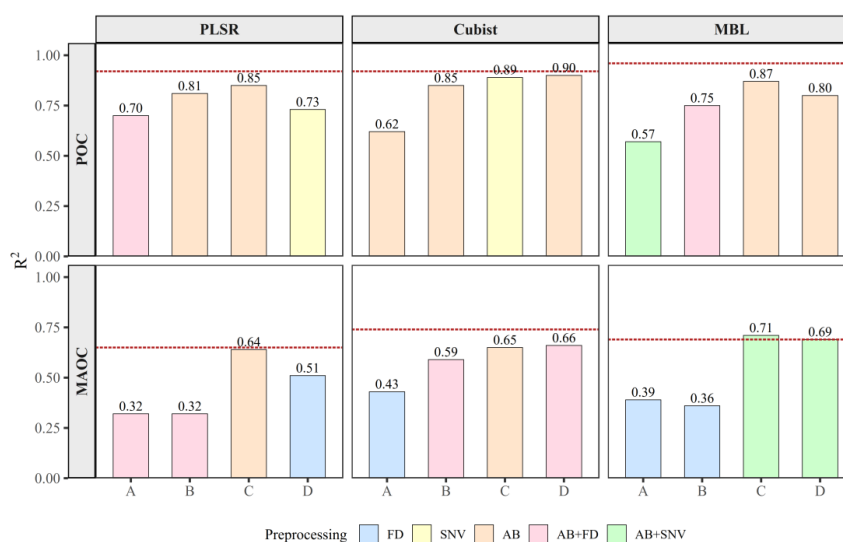
352 Workflow C consistently outperformed the transfer-based workflows. POC prediction improved  
353 substantially, with  $R^2$  values ranging from 0.57 to 0.89. MAOC prediction also improved, with  $R^2$   
354 increasing to 0.50–0.71.

355 Calibrating models using EPO-corrected spectra (Workflow D) achieved performance comparable  
356 to or slightly better than Workflow C, although improvements were primarily observed for the Cubist  
357 model (POC:  $R^2$  0.74–0.90; MAOC:  $R^2$  0.53–0.66). For PLSR and MBL, EPO correction led to minor  
358 reductions or negligible changes in predictive accuracy compared with Workflow C.

359 For in-situ applications, Workflow D provided the most robust and stable performance for POC,



360 with the Cubist model and AB preprocessing yielding the best balance of accuracy and robustness. In  
 361 contrast, MAOC prediction performed best under Workflow C with MBL and AB-SNV preprocessing,  
 362 whereas EPO correction did not provide additional improvement. The optimal in-situ modeling strategy  
 363 identified for each target variable was subsequently evaluated using the independent 2024 dataset.



**Figure 6.** Performance of PLSR, Cubist, and MBL models across the four in-situ application workflows: (A) direct transfer of laboratory-calibrated models to in-situ spectra; (B) transfer of laboratory-calibrated models to EPO-corrected in-situ spectra; (C) direct in-situ model development; and (D) EPO-corrected in-situ model development. Acronyms in the legend denote preprocessing methods: REF, raw reflectance spectra; FD, first derivative; SNV, standard normal variate; and AB, absorbance transformation. The red dashed line indicates the optimal R<sup>2</sup> obtained under laboratory conditions.

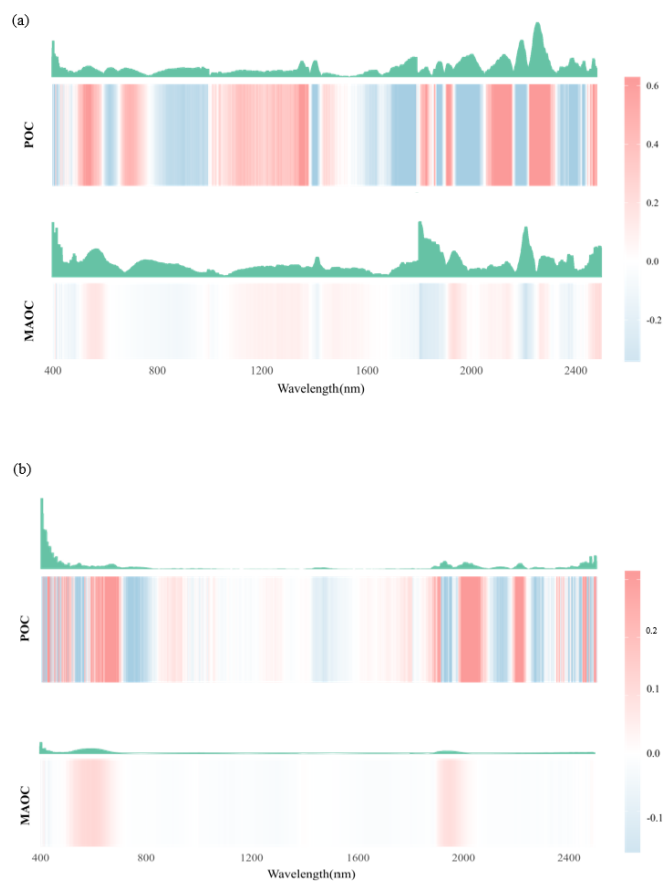
### 364 3.5 Interpretation of the optimal SOC fractions prediction models

365 SHAP analysis was first conducted to interpret the best-performing laboratory models (models using  
 366 spectra preprocessed with the AB\_SNV transformation (**Fig. 7a**). The patterns of relative contributions  
 367 for POC and MAOC were relatively stable across wavelengths, reflecting the consistent spectral  
 368 responses observed under controlled laboratory conditions. For POC, the most influential wavelengths



369 were primarily concentrated between 2,200 and 2,400 nm, which are commonly associated with aliphatic  
370 C–H absorption features and mineral–organic matter interactions. In contrast, MAOC contributions were  
371 mainly located in the 500–600 nm region, likely related to iron oxide absorption features and aromatic  
372 organic matter components.

373 For the in-situ models, the optimal POC model was developed using EPO-corrected spectra. SHAP  
374 analysis indicates that EPO enhances chemically relevant spectral bands, resulting in wavelength  
375 contribution patterns that more closely resemble those observed under laboratory conditions (**Fig. 7b**).  
376 This suggests that EPO effectively suppresses environmental noise while preserving meaningful spectral  
377 information for POC prediction, which is consistent with its underlying. In contrast, the best MAOC  
378 model was constructed directly from uncorrected in-situ spectra and exhibited weak, localized SHAP  
379 responses across the spectral range. This indicates that environmental noise in field-acquired spectra  
380 strongly masks the characteristic spectral signals of MAOC, and that EPO correction may further remove  
381 weak but informative features. These results highlight the inherent difficulty of retrieving MAOC from  
382 in-situ hyperspectral observations.



**Figure 7.** Feature importance for different kinds of organic carbon derived from the SHAP values of the optimal models. (a) Laboratory models; (b) In-situ models. The heatmap below each line graph shows SHAP-derived feature importance for stratified sampling. The bars above each heatmap represent absolute mean SHAP values across all samples.

### 383 3.6 Cross-year transferability of the prediction models

384 To evaluate model robustness across varying environmental states, an independent dataset collected  
385 in 2024 was used to externally validate models developed on the 2025 dataset. Because both datasets  
386 were acquired at the same site but under different field conditions, the cross-year comparison mainly  
387 reflects environmental effects, particularly differences in soil moisture and in-situ measurement states,

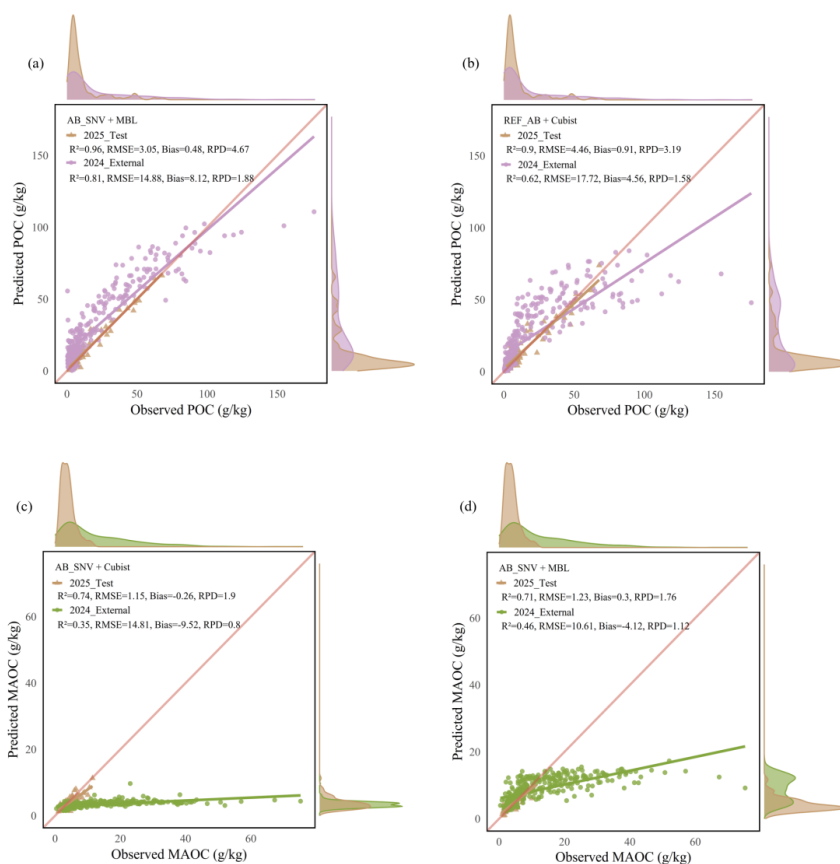


388 rather than true temporal dynamics. Therefore, this analysis focuses on environmental transferability  
389 driven by spectral domain shifts.

390 Overall, predictive performance declined on the independent 2024 dataset compared with the 2025  
391 test dataset (**Fig. 8**), indicating reduced robustness under altered environmental conditions. For  
392 laboratory-based models, POC predictions retained relatively high accuracy on the external dataset ( $R^2$   
393 = 0.81, RPD = 1.88), suggesting moderate environmental transferability. In contrast, MAOC predictions  
394 showed substantial degradation, with an  $R^2$  of 0.35 and RPD below 1.0, indicating very limited  
395 transferability for this fraction.

396 In-situ models exhibited consistently lower predictive performance on the external dataset. POC  
397 predictions declined to fair accuracy ( $R^2 = 0.62$ ), whereas MAOC predictions remained unreliable (RPD  
398 < 1.0). Comparison of marginal density distributions (**Fig. 8**) further revealed that laboratory-based POC  
399 predictions exhibited greater overlap between the calibration and external datasets, whereas in-situ  
400 models showed pronounced distributional shifts, indicative of a stronger spectral-domain mismatch.

401



**Figure 8.** Comparison of environmental transferability between optimal laboratory and in-situ models. The left column shows the prediction performance of laboratory-based models, while the right column shows the prediction performance of in-situ models. Marginal density distributions of the test (2025) and external (2024) datasets are shown along the top and right axes.

## 402 4. Discussion

### 403 4.1 Feasibility of Vis-NIR prediction for SOC fractions

404 This study systematically evaluated 12 machine learning algorithms to predict forest soil organic  
 405 carbon fractions from laboratory Vis-NIR spectra. Under controlled laboratory conditions, the models  
 406 achieved reliable predictions for POC and MAOC (Fig. 5), confirming the feasibility of Vis-NIR



407 spectroscopy for estimating SOC fractions in forest soils. Similar levels of accuracy have been reported  
408 in previous studies, with  $R^2$  values of 0.65 and 0.67 for POC and MAOC, respectively (Dai et al., 2024).  
409 Among the evaluated algorithms, PLSR, Cubist, and MBL consistently ranked among the best-  
410 performing models, indicating that both linear and nonlinear approaches can achieve high predictive  
411 accuracy when spectral variability is minimized and measurement conditions are well controlled.

412 However, previous studies have reported inconsistent conclusions regarding the optimal modeling  
413 strategy, with model performance varying across datasets, target fractions, and spectral configurations  
414 (Das et al. 2023; Munnaf and Mouazen, 2022). Together with our results, this variability suggests that no  
415 single algorithm can be universally recommended for Vis-NIR-based prediction of SOC fractions.  
416 Instead, systematic comparison of multiple modeling strategies remains necessary to identify models that  
417 best match the spectral characteristics and compositional properties of the target fractions.

418 When the best-performing laboratory models were directly applied to in-situ spectra, prediction  
419 accuracy declined markedly (**Fig. 6**). This reduction is consistent with the pronounced spectral  
420 discrepancies between laboratory and in-situ measurements (**Figs. 4a, S1**), where in-situ spectra  
421 exhibited lower reflectance due to soil moisture effects. Moisture-induced attenuation masked intrinsic  
422 absorption features, reduced the signal-to-noise ratio, and shifted the spectral feature space, thereby  
423 limiting the applicability of models trained under controlled laboratory conditions (Piccini et al., 2024).  
424 Under this transfer scenario, POC retained only moderate predictability, whereas MAOC prediction  
425 deteriorated substantially, reflecting its weaker and more indirect spectral expression.

426 By contrast, models developed directly from in-situ spectra achieved higher predictive accuracy  
427 than laboratory models applied to in-situ spectra. Because these models incorporate environmental



428 variability during calibration, they better capture both intrinsic spectral–chemical relationships and  
429 systematic field-induced spectral variation. These results highlight the importance of calibrating models  
430 directly under in-situ conditions when accurate field predictions are required, particularly for fractions  
431 such as MAOC that are sensitive to environmental noise.

#### 432 **4.2 Role of EPO correction across in-situ spectral application**

433 EPO played a differentiated role across the in-situ spectral application scenarios examined in this  
434 study, highlighting both its advantages and limitations depending on the modeling strategy and target  
435 variable. For laboratory-to-field model transfer, EPO substantially improved the applicability of  
436 laboratory-calibrated models to in-situ spectral data. Compared with the direct application of laboratory  
437 models to raw in-situ spectra, the use of EPO-corrected spectra consistently resulted in higher predictive  
438 performance and lower prediction errors (**Fig. 6**), particularly for the PLSR and Cubist models. This  
439 improvement demonstrates the effectiveness of EPO in reducing spectral discrepancies between  
440 laboratory and field measurements.

441 By removing spectral components orthogonal to the target soil properties—such as variability due  
442 to soil moisture, surface roughness, and particle-size heterogeneity—EPO reduces non-informative  
443 spectral variance and improves the statistical compatibility between laboratory and in-situ spectral  
444 domains. As a result, models calibrated under controlled conditions become more transferable to field-  
445 acquired spectra.

446 In contrast, the influence of EPO on models calibrated directly with in-situ spectra was strongly  
447 model-dependent. For PLSR and MBL, EPO correction resulted in reduced predictive performance,  
448 suggesting that the orthogonalization process may inadvertently remove spectral information relevant to



449 soil carbon fractions. Both PLSR and MBL rely strongly on covariance structures and spectral similarity  
450 patterns, which may be partially disrupted when spectral components are removed during EPO correction.  
451 Conversely, Cubist models exhibited improved performance after EPO correction. The rule-based  
452 structure and local linear regression framework of Cubist allow it to exploit the simplified feature space  
453 generated by EPO while retaining relevant predictive information, explaining its stronger robustness  
454 under EPO-corrected in-situ modeling conditions.

455 These results indicate that the effectiveness of EPO depends strongly on the modeling objective and  
456 algorithmic structure. Rather than acting as a universally optimal preprocessing method, EPO functions  
457 as a conditional tool whose benefits are most evident when reducing measurement-condition mismatch  
458 between laboratory and in-situ spectra.

#### 459 **4.3 Spectral band contributions and mechanistic interpretation of SOC fractions**

460 SHAP analysis of the optimal laboratory models showed that the spectral regions contributing to  
461 SOC fraction predictions align with known soil organic–mineral interactions (**Fig. 7a**). These  
462 contributions match the diagnostic spectral patterns identified in the visible and shortwave infrared  
463 regions (**Fig. 4**), indicating that the models capture physically meaningful relationships rather than purely  
464 statistical correlations.

465 For the POC, dominant contributions were concentrated in the shortwave infrared region,  
466 characterized by absorption features associated with organic functional groups. These signals are directly  
467 associated with plant-derived residues and microbial products, which constitute the main components of  
468 particulate organic matter (Demattê et al., 2004; Stenberg et al., 2010). As a result, POC exhibits a  
469 relatively strong and direct spectral expression, enabling robust and stable model predictions across



470 different modeling scenarios.

471 In contrast, MAOC showed stronger contributions in the visible region, where spectral variability  
472 is largely controlled by mineralogical features. This reflects the fact that MAOC is primarily associated  
473 with mineral surfaces rather than existing as a spectrally distinct organic phase. Consequently, its spectral  
474 representation is mediated through mineral-related signals, which are not uniquely specific to MAOC  
475 but are influenced by broader soil compositional characteristics (Schwertmann and Taylor, 1989; Chen  
476 et al., 2014; Kleber et al., 2015).

477 The comparatively lower prediction accuracy of MAOC is therefore not attributable to a single  
478 factor but rather reflects the combined effects of its spectral representation and the data characteristics.  
479 First, the mineral-mediated nature of MAOC results in a less specific, more context-dependent spectral  
480 signal than POC. Second, the narrower concentration range and lower variability of MAOC observed in  
481 this study reduce the effective signal-to-noise ratio available for model calibration. Third, because  
482 mineral-associated spectral signals are particularly sensitive to environmental conditions such as soil  
483 moisture and surface heterogeneity (Dai, 2025), which can substantially alter spectral responses in the  
484 field (Roudier et al., 2017), MAOC predictions are more susceptible to spectral-domain shifts under in  
485 situ and cross-year conditions.

486 These factors highlight two fundamentally different modes of spectral detectability: a direct,  
487 organic-dominated signal for POC in the shortwave infrared region, and an indirect, mineral-mediated  
488 signal for MAOC in the visible region (Breure et al., 2025). This contrast is consistent with the conceptual  
489 basis of soil organic matter fractionation, in which size-based separation distinguishes functionally  
490 distinct carbon pools. POC represents a relatively labile pool dominated by plant residues and microbial



491 products, characterized by rapid turnover and strong intrinsic spectral signatures. In contrast, MAOC  
492 represents a more stable pool in which organic compounds are associated with mineral surfaces and  
493 protected from decomposition, resulting in a weaker and less specific spectral expression mediated by  
494 mineral-related signals (Wang et al., 2024; Yang et al., 2023).

495 These differences indicate that the observed variation in spectral contributions and model  
496 performance reflects not only spectroscopic behavior but also fundamental differences in the form and  
497 function of SOC fractions. This mechanistic distinction explains the contrasting predictive accuracy and  
498 transferability between POC and MAOC. It highlights the importance of considering both compositional  
499 characteristics and data structure when interpreting and applying Vis–NIR-based soil carbon models.

#### 500 **4.4 Challenges in the robustness of soil carbon models under cross-year environmental** 501 **variability**

502 The cross-year validation results highlight the challenges of achieving robust Vis–NIR-based  
503 prediction of soil carbon fractions under changing environmental conditions. The observed performance  
504 degradation is primarily attributable to environmental state–induced spectral domain shifts rather than to  
505 model instability alone (**Fig. S2**). Variations in soil moisture and associated changes in surface conditions  
506 substantially alter the statistical structure of in-situ spectra (Roudier et al., 2017), thereby limiting model  
507 generalization beyond the calibration domain.

508 Differences in robustness between laboratory-based and in-situ models reflect contrasting levels of  
509 spectral control. Laboratory spectra are obtained under standardized conditions, resulting in more stable  
510 spectral–chemical relationships and relatively stronger transferability for POC. In contrast, in-situ spectra  
511 integrate multiple sources of uncontrolled variability, including moisture fluctuations and surface



512 heterogeneity, which increase spectral domain mismatch when models are applied across environmental  
513 states. These findings demonstrate that high predictive accuracy in a single sampling campaign does not  
514 necessarily guarantee stable performance across years or conditions.

515       The influence of spectral correction further illustrates the conditional nature of model robustness.  
516 EPO improved laboratory-to-field transferability for POC by suppressing moisture-dominated spectral  
517 variance that is largely orthogonal to the target variable. However, the benefits of EPO were limited or  
518 even negative for MAOC. This likely reflects a violation of the orthogonality assumption underlying  
519 EPO: unlike POC, MAOC is closely associated with mineral surfaces and surface-bound water, resulting  
520 in statistical coupling between moisture-related spectral features and MAOC content.

521       SHAP analysis further supports this mechanistic interpretation by demonstrating that dominant  
522 spectral contributions captured by the models correspond to known soil mineralogical and organic  
523 interaction processes. Together, these results indicate that robust application of Vis–NIR spectroscopy  
524 for predicting forest SOC fractions requires explicit consideration of environment–driven spectral  
525 domain shifts and independent validation under contrasting field conditions.

526       From a methodological perspective, the optimal modeling strategy depends strongly on the intended  
527 application context. Laboratory-based models remain valuable for controlled spectral characterization  
528 and for establishing standardized spectral–chemical relationships. However, when the objective is  
529 reliable field prediction under environmentally variable conditions, direct in-situ calibration offers a more  
530 robust strategy because it incorporates environmental variability into model development.

531       The role of EPO should likewise be interpreted conditionally. EPO is most effective when the  
532 objective is to improve laboratory-to-field model transfer by reducing systematic spectral discrepancies,



533 particularly for spectrally intrinsic fractions such as POC. In contrast, when models are calibrated directly  
534 using in-situ spectra, the benefits of EPO become model- and target-dependent and may even be negative  
535 for proxy-mediated variables such as MAOC. Overall, these findings demonstrate that robust field  
536 spectroscopy depends not only on algorithm selection but also on aligning preprocessing and calibration  
537 strategies with the spectral expression of the target soil property.

## 538 **5. Conclusion**

539 This study evaluated the capability of Vis–NIR spectroscopy to predict SOC fractions in forest soils  
540 using either laboratory or in-situ spectra, multiple modeling strategies, spectral pre-processing options,  
541 and independent cross-year validation. Under laboratory conditions, reliable predictions were obtained  
542 for both POC and MAOC. The AB\_SNV preprocessing strategy produced the most consistent results,  
543 with MBL performing best for POC and Cubist providing the most accurate predictions for MAOC.

544 Model performance declined substantially when laboratory-calibrated models were directly  
545 transferred to in-situ spectra, primarily due to environmental spectral discrepancies associated with soil  
546 moisture and surface conditions. This illustrates the importance of dedicated transfer strategies. Spectral  
547 correction using EPO improved laboratory-to-field model transfer for POC, but its effectiveness varied  
548 across algorithms and target fractions. In contrast, models calibrated directly using in-situ spectra  
549 achieved more stable predictions, indicating that incorporating environmental variability during  
550 calibration is essential for reliable field applications.

551 The interpretation of the calibrated models further demonstrated that they captured physically  
552 meaningful spectral information: POC predictions were mainly associated with organic absorption  
553 features in the shortwave infrared region, whereas MAOC predictions relied more strongly on mineral



554 proxy signals in the visible range.

555 Overall, this study demonstrates the feasibility of directly quantifying SOC fractions (POC and

556 MAOC) in-situ using Vis–NIR spectroscopy under forest field conditions. Providing new methodological

557 insights into soil spectroscopy under field conditions and establishing a foundation for advancing in situ

558 monitoring of forest soil carbon dynamics.

559



560 **Declaration of competing interest**

561 The authors declare that they have no known competing financial interests or personal relationships that  
562 could have appeared to influence the work reported in this paper.

563

564 **Author contributions**

565 JL performed the measurements and analyzed the data; JL and YZ conceived and executed the research  
566 and wrote the manuscript draft; WF, LS, and YL gave suggestions about the approach and contributed  
567 extensively to the paper. AW, JW, and PR reviewed and edited the manuscript. All co-authors revised the  
568 paper.

569

570 **Acknowledgements**

571 This work was supported by the National Key Research and Development Program of China  
572 (2022YFF1300501), the National Natural Science Foundation of China (32271873), the Youth Start-up  
573 Fund, Institute of Applied Ecology, Chinese Academy of Sciences, the China Postdoctoral Science  
574 Foundation (2023M743699), and the Program of China Scholarship Council (Grant No. 202504910271).

575 The authors thank the reviewers for their constructive comments, which have helped improve our  
576 manuscript.

577

578 **Data availability**

579 The datasets generated during the current study are available from the corresponding author upon  
580 reasonable request.

581



582 **References**

- 583 Abramoff, R.Z., Guenet, B., Zhang, H., Georgiou, K., Xu, X., Viscarra Rossel, R.A., Yuan, W., Ciais, P., 2022.  
584 Improved global-scale predictions of soil carbon stocks with millennial version 2. *Soil Biology and*  
585 *Biochemistry* 164, 108466. <https://doi.org/10.1016/j.soilbio.2021.108466>
- 586 Breure, T.S., De Rosa, D., Panagos, P., Cotrufo, M.F., Jones, A., Lugato, E., 2025. Revisiting the soil carbon  
587 saturation concept to inform a risk index in european agricultural soils. *Nat Commun* 16, 2538.  
588 <https://doi.org/10.1038/s41467-025-57355-y>
- 589 Chen, C., Dynes, J.J., Wang, J., Sparks, D.L., 2014. Properties of fe-organic matter associations via coprecipitation  
590 versus adsorption. *Environ. Sci. Technol.* 48, 13751–13759. <https://doi.org/10.1021/es503669u>
- 591 Conforti, M., Matteucci, G., Buttafuoco, G., 2018. Using laboratory vis-NIR spectroscopy for monitoring some  
592 forest soil properties. *J Soils Sediments* 18, 1009–1019. <https://doi.org/10.1007/s11368-017-1766-5>
- 593 Cotrufo, M.F., Ranalli, M.G., Haddix, M.L., Six, J., Lugato, E., 2019. Soil carbon storage informed by particulate  
594 and mineral-associated organic matter. *Nat. Geosci.* 12, 989–994. <https://doi.org/10.1038/s41561-019-0484-6>
- 595
- 596 Daelemans, W., 2009. Memory-based language processing [WWW Document]. URL  
597 [https://www.cambridge.org/core/books/memorybased-language-](https://www.cambridge.org/core/books/memorybased-language-processing/8B3E1D0E5A2EA90E9DF509B67E410033)  
598 [processing/8B3E1D0E5A2EA90E9DF509B67E410033](https://www.cambridge.org/core/books/memorybased-language-processing/8B3E1D0E5A2EA90E9DF509B67E410033) (accessed 4.27.25).
- 599 Dai, L., 2025. Prediction of soil organic carbon fractions in tropical cropland using a regional visible and near-  
600 infrared spectral library and machine learning.
- 601 Dai, L., Xue, J., Lu, R., Wang, Z., Chen, Z., Yu, Q., Shi, Z., Chen, S., 2024. In-situ prediction of soil organic carbon  
602 contents in wheat-rice rotation fields via visible near-infrared spectroscopy. *Soil & Environmental Health*  
603 *2*, 100113. <https://doi.org/10.1016/j.seh.2024.100113>
- 604 Das, B., Chakraborty, D., Singh, V.K., Das, D., Sahoo, R.N., Aggarwal, P., Murgaokar, D., Mondal, B.P., 2023.  
605 Partial least square regression based machine learning models for soil organic carbon prediction using  
606 visible–near infrared spectroscopy. *Geoderma Regional* 33, e00628.  
607 <https://doi.org/10.1016/j.geodrs.2023.e00628>
- 608 De Vos, B., Cools, N., Ilvesniemi, H., Vesterdal, L., Vanguelova, E., Carnicelli, S., 2015. Benchmark values for  
609 forest soil carbon stocks in europe: Results from a large scale forest soil survey. *Geoderma* 251–252, 33–  
610 46. <https://doi.org/10.1016/j.geoderma.2015.03.008>
- 611 Delahaie, A.A., Cécillon, L., Stojanova, M., Abiven, S., Arbelet, P., Arrouays, D., Baudin, F., Bispo, A., Boulonne,  
612 L., Chenu, C., Heinonsalo, J., Jolivet, C., Karhu, K., Martin, M., Pacini, L., Poeplau, C., Ratié, C., Roudier,  
613 P., Saby, N.P.A., Savignac, F., Barré, P., 2024. Investigating the complementarity of thermal and physical  
614 soil organic carbon fractions. *SOIL* 10, 795–812. <https://doi.org/10.5194/soil-10-795-2024>
- 615 Demattê, J.A., Bellinaso, H., Romero, D., Fongaro, C., 2014. Morphological interpretation of reflectance spectrum  
616 (MIRS) using libraries looking towards soil classification. *Scientia Agricola* 71, 509–520.  
617 <https://doi.org/10.1590/0103-9016-2013-0365>
- 618 Demattê, J.A.M., Campos, R.C., Alves, M.C., Fiorio, P.R., Nanni, M.R., 2004. Visible–NIR reflectance: A new  
619 approach on soil evaluation. *Geoderma* 121, 95–112. <https://doi.org/10.1016/j.geoderma.2003.09.012>
- 620 Haghi, K.R., Pérez-Fernández, E., Robertson, J., 2021. Prediction of various soil properties for a national spatial  
621 dataset of scottish soils based on four different chemometric approaches: A comparison of near infrared  
622 and mid-infrared spectroscopy. *Geoderma* 396, 115071. <https://doi.org/10.1016/j.geoderma.2021.115071>
- 623 Jaconi, A., Don, A., Freibauer, A., 2017. Prediction of soil organic carbon at the country scale: Stratification  
624 strategies for near-infrared data. *European Journal of Soil Science* 68, 919–929.



- 625 <https://doi.org/10.1111/ejss.12485>
- 626 Janik, L., Skjemstad, J., Shepherd, K., Spouncer, L., 2007. The prediction of soil carbon fractions using mid-infrared-  
627 partial least square analysis. *Australian Journal of Soil Research* 45, 73–81.  
628 <https://doi.org/10.1071/SR06083>
- 629 Kaiser, K., Guggenberger, G., 2003. Mineral surfaces and soil organic matter. *European Journal of Soil Science* 54,  
630 219–236. <https://doi.org/10.1046/j.1365-2389.2003.00544.x>
- 631 Kennard, R.W., Stone, L.A., 1969. Computer aided design of experiments. *Technometrics* 11, 137–148.  
632 <https://doi.org/10.1080/00401706.1969.10490666>
- 633 Kleber, M., Eusterhues, K., Keiluweit, M., Mikutta, C., Mikutta, R., Nico, P.S., 2015. Chapter one - mineral–organic  
634 associations: Formation, properties, and relevance in soil environments, in: Sparks, D.L. (Ed.), *Advances*  
635 *in Agronomy*. Academic Press, pp. 1–140. <https://doi.org/10.1016/bs.agron.2014.10.005>
- 636 Koch, M., Schodlok, M.C., Guggenberger, G., Stadler, S., 2021. Effects of water tension and surface roughness on  
637 soil hyperspectral reflectance. *Geoderma* 385, 114888. <https://doi.org/10.1016/j.geoderma.2020.114888>
- 638 Lavallee, J.M., Soong, J.L., Cotrufo, M.F., 2020. Conceptualizing soil organic matter into particulate and mineral-  
639 associated forms to address global change in the 21st century. *Global Change Biology* 26, 261–273.  
640 <https://doi.org/10.1111/gcb.14859>
- 641 Lundberg, S.M., Lee, S.-I., 2017. A unified approach to interpreting model predictions, in: *Proceedings of the 31st*  
642 *International Conference on Neural Information Processing Systems, NIPS'17*. Curran Associates Inc.,  
643 Red Hook, NY, USA, pp. 4768–4777.
- 644 Ma, Y., Minasny, B., Demattê, J.A.M., McBratney, A.B., 2023. Incorporating soil knowledge into machine-learning  
645 prediction of soil properties from soil spectra. *European Journal of Soil Science* 74, e13438.  
646 <https://doi.org/10.1111/ejss.13438>
- 647 Minasny, B., McBratney, A.B., Bellon-Maurel, V., Roger, J.-M., Gobrecht, A., Ferrand, L., Joalland, S., 2011.  
648 Removing the effect of soil moisture from NIR diffuse reflectance spectra for the prediction of soil organic  
649 carbon. *Geoderma* 167–168, 118–124. <https://doi.org/10.1016/j.geoderma.2011.09.008>
- 650 Munnaf, M.A., Mouazen, A.M., 2022. Removal of external influences from on-line vis-NIR spectra for predicting  
651 soil organic carbon using machine learning. *CATENA* 211, 106015.  
652 <https://doi.org/10.1016/j.catena.2022.106015>
- 653 Pacini, L., Schiedung, M., Stojanova, M., Roudier, P., Arbelet, P., Barré, P., Baudin, F., Cambou, A., Cécillon, L.,  
654 Heinonsalo, J., Karhu, K., McNally, S., Omondigabe, P., Poeplau, C., Saby, N.P.A., 2025. Predicting the  
655 proportion of centennially stable soil organic carbon using mid-infrared spectroscopy. *Geoderma* 462,  
656 117536. <https://doi.org/10.1016/j.geoderma.2025.117536>
- 657 Padarian, J., McBratney, A.B., Minasny, B., 2020. Game theory interpretation of digital soil mapping convolutional  
658 neural networks. *SOIL* 6, 389–397. <https://doi.org/10.5194/soil-6-389-2020>
- 659 Piccini, C., Metzger, K., Debaene, G., Stenberg, B., Götzinger, S., Borůvka, L., Sandén, T., Bragazza, L., Liebisch,  
660 F., 2024. In-field soil spectroscopy in vis–NIR range for fast and reliable soil analysis: A review. *European*  
661 *Journal of Soil Science* 75, e13481. <https://doi.org/10.1111/ejss.13481>
- 662 Qi, M., Chen, S., Wei, Y., Zhou, H., Zhang, S., Wang, M., Zheng, J., Viscarra Rossel, R.A., Chang, J., Shi, Z., Luo,  
663 Z., 2024. Using visible-near infrared spectroscopy to estimate whole-profile soil organic carbon and its  
664 fractions. *Soil & Environmental Health* 2, 100100. <https://doi.org/10.1016/j.seh.2024.100100>
- 665 Rossel, R.A.V., Behrens, T., 2010. Using data mining to model and interpret soil diffuse reflectance spectra.  
666 *Geoderma, Diffuse reflectance spectroscopy in soil science and land resource assessment* 158, 46–54.  
667 <https://doi.org/10.1016/j.geoderma.2009.12.025>



- 668 Roudier, P., Hedley, C. b., Ross, C. w., 2015. Prediction of volumetric soil organic carbon from field-moist intact  
669 soil cores. *European Journal of Soil Science* 66, 651–660. <https://doi.org/10.1111/ejss.12259>
- 670 Roudier, P., Hedley, C.B., Lobsey, C.R., Viscarra Rossel, R.A., Leroux, C., 2017. Evaluation of two methods to  
671 eliminate the effect of water from soil vis–NIR spectra for predictions of organic carbon. *Geoderma* 296,  
672 98–107. <https://doi.org/10.1016/j.geoderma.2017.02.014>
- 673 Schwertmann, U., Taylor, R. m., 1989. Iron oxides, in: *Minerals in Soil Environments*. John Wiley & Sons, Ltd, pp.  
674 379–438. <https://doi.org/10.2136/sssabookser1.2ed.e8>
- 675 Silvero, N.E.Q., Di Raimo, L.A.D.L., Pereira, G.S., Magalhães, L.P. de, Terra, F. da S., Dassan, M.A.A., Salazar,  
676 D.F.U., Demattê, J.A.M., 2020. Effects of water, organic matter, and iron forms in mid-IR spectra of soils:  
677 Assessments from laboratory to satellite-simulated data. *Geoderma* 375, 114480.  
678 <https://doi.org/10.1016/j.geoderma.2020.114480>
- 679 Sokol, K., Flach, P., 2020. LIMETree: Interactively customisable explanations based on local surrogate multi-output  
680 regression trees. <https://doi.org/10.48550/arXiv.2005.01427>
- 681 Stenberg, B., Viscarra Rossel, R.A., Mouazen, A.M., Wetterlind, J., 2010. Chapter five - visible and near infrared  
682 spectroscopy in soil science, in: Sparks, D.L. (Ed.), *Advances in Agronomy*. Academic Press, pp. 163–215.  
683 [https://doi.org/10.1016/S0065-2113\(10\)07005-7](https://doi.org/10.1016/S0065-2113(10)07005-7)
- 684 Viscarra Rossel, R.A., Behrens, T., Ben-Dor, E., Brown, D.J., Demattê, J.A.M., Shepherd, K.D., Shi, Z., Stenberg,  
685 B., Stevens, A., Adamchuk, V., Aichi, H., Barthès, B.G., Bartholomeus, H.M., Bayer, A.D., Bernoux, M.,  
686 Böttcher, K., Brodský, L., Du, C.W., Chappell, A., Fouad, Y., Genot, V., Gomez, C., Grunwald, S., Gubler,  
687 A., Guerrero, C., Hedley, C.B., Knadel, M., Morrás, H.J.M., Nocita, M., Ramirez-Lopez, L., Roudier, P.,  
688 Campos, E.M.R., Sanborn, P., Sellitto, V.M., Sudduth, K.A., Rawlins, B.G., Walter, C., Winowiecki, L.A.,  
689 Hong, S.Y., Ji, W., 2016. A global spectral library to characterize the world’s soil. *Earth Sci. Rev.* 155,  
690 198–230. <https://doi.org/10.1016/j.earscirev.2016.01.012>
- 691 Viscarra Rossel, R.A., Behrens, T., Ben-Dor, E., Chabrilat, S., Demattê, J.A.M., Ge, Y., Gomez, C., Guerrero, C.,  
692 Peng, Y., Ramirez-Lopez, L., Shi, Z., Stenberg, B., Webster, R., Winowiecki, L., Shen, Z., 2022. Diffuse  
693 reflectance spectroscopy for estimating soil properties: A technology for the 21st century. *European*  
694 *Journal of Soil Science* 73, e13271. <https://doi.org/10.1111/ejss.13271>
- 695 Viscarra Rossel, R.A., McGlynn, R.N., McBratney, A.B., 2006. Determining the composition of mineral-organic  
696 mixes using UV–vis–NIR diffuse reflectance spectroscopy. *Geoderma* 137, 70–82.  
697 <https://doi.org/10.1016/j.geoderma.2006.07.004>
- 698 Viscarra Rossel, R. A., Walvoort, D.J.J., McBratney, A.B., Janik, L.J., Skjemstad, J.O., 2006. Visible, near infrared,  
699 mid infrared or combined diffuse reflectance spectroscopy for simultaneous assessment of various soil  
700 properties. *Geoderma* 131, 59–75. <https://doi.org/10.1016/j.geoderma.2005.03.007>
- 701 Viscarra Rossel, R.A., Zhang, M., Behrens, T., Webster, R., 2024. A warming climate will make Australian soil a net  
702 emitter of atmospheric CO<sub>2</sub>. *npj Clim Atmos Sci* 7, 1–11. <https://doi.org/10.1038/s41612-024-00619-z>
- 703 Wang, H.-P., Chen, P., Dai, J.-W., Liu, D., Li, J.-Y., Xu, Y.-P., Chu, X.-L., 2022. Recent advances of chemometric  
704 calibration methods in modern spectroscopy: Algorithms, strategy, and related issues. *TrAC Trends in*  
705 *Analytical Chemistry* 153, 116648. <https://doi.org/10.1016/j.trac.2022.116648>
- 706 Wang, Z., Chen, S., Lu, R., Zhang, X., Ma, Y., Shi, Z., 2024. Non-linear memory-based learning for predicting soil  
707 properties using a regional vis-NIR spectral library. *Geoderma* 441, 116752.  
708 <https://doi.org/10.1016/j.geoderma.2023.116752>
- 709 Wold, S., Antti, H., Lindgren, F., Öhman, J., 1998. Orthogonal signal correction of near-infrared spectra.  
710 *Chemometrics and Intelligent Laboratory Systems* 44, 175–185. [38](https://doi.org/10.1016/S0169-</a></p></div><div data-bbox=)



711 7439(98)00109-9  
712 Wu, Y., Peng, Z., Wang, X., Huang, J., Yang, L., Liu, L., 2025. Warmer climate reduces the carbon storage, stability  
713 and saturation levels in forest soils. *Earth's Future* 13, e2024EF004988.  
714 <https://doi.org/10.1029/2024EF004988>  
715 Yang Y., Wu F., Wu Q., Zhu J., Ni X., 2023. Soil organic carbon associated with iron oxides in terrestrial ecosystems:  
716 Content, distribution and control. *Chin. Sci. Bull.* 68, 695–704. <https://doi.org/10.1360/TB-2022-0728>  
717 Zhou, Z., Ren, C., Wang, C., Delgado-Baquerizo, M., Luo, Y., Luo, Z., Du, Z., Zhu, B., Yang, Y., Jiao, S., Zhao, F.,  
718 Cai, A., Yang, G., Wei, G., 2024. Global turnover of soil mineral-associated and particulate organic carbon.  
719 *Nat Commun* 15, 5329. <https://doi.org/10.1038/s41467-024-49743-7>  
720