

1 **Improvement of Soil Properties Maps using an Iterative Residual Correction Method**

2 Chengcheng Xu<sup>1</sup>, Elia Scudiero<sup>2,3</sup>, Ray Anderson<sup>3,2</sup>, Nathaniel Chaney<sup>1</sup>

3 <sup>1</sup> Department of Civil and Environmental Engineering, Duke University, Durham, NC 27705,  
4 USA

5 <sup>2</sup> Department of Environmental Sciences, University of California Riverside, Riverside, CA  
6 92521, USA

7 <sup>3</sup> United States Department of Agriculture – Agricultural Research Service, George E. Brown  
8 Jr. Salinity Laboratory, Agricultural Water Efficiency and Salinity Research Unit, Riverside,  
9 CA 92507, USA

10 *Correspondence to:* Chengcheng Xu (Chengcheng.xu@duke.edu)

11

12 **Short Summary**

13 ~~Accurate soil information is vital. This study developed a method to improve existing~~  
14 ~~probabilistic soil maps — spatially continuous maps providing prior estimates — by~~  
15 ~~correcting their probability distributions as new soil data arrives. By iteratively adjusting~~  
16 ~~previous predictions, the method increases both accuracy and certainty of soil maps. Its~~  
17 ~~application in California enhanced predictions for several soil properties. This method can~~  
18 ~~be further used for more soil properties and regions.~~

19

20 **Abstract**

21 Accurate mapping of soil properties is vital for many applications, ~~yet~~ existing models for  
22 digital soil maps ~~often~~ underestimate their spatial variability or prediction uncertainties,

**Deleted:** Accurate soil information is important. This study developed a new method that improves existing soil maps by correcting their probability distributions using newly collected soil measurements. By repeatedly adjusting previous predictions, the method makes soil maps more accurate and more certain. The application in California improved the performance of predictions for soil texture, organic matter, and bulk density. This method can be further used for more soil properties and regions....

**Deleted:** including precision irrigation and fertilization. However,...

35 which introduces risk for applications such as irrigation and drainage management. This  
 36 study introduces an approach — iterative residual correction (IRC) — to update existing  
 37 probabilistic soil maps when new soil observations become available. We demonstrated  
 38 its application for enhanced soil mapping performance using a Californian case study. To  
 39 implement this, we first generate prior probabilistic soil property maps using a pruned  
 40 hierarchical Random Forest (pHRF) method. These prior estimates are then refined by  
 41 integrating additional soil profile data and iteratively adjusting residuals of distribution of  
 42 soil properties (reducing differences between observations and prior predictions) pixel by  
 43 pixel. For this purpose, we employed Random Forest regressors to gradually adjusts the  
 44 soil property distributions and incrementally corrects prior bias. Updated soil maps were  
 45 evaluated over California and at 1-km resolution to test the methodology, using additional  
 46 soil observations from the World Soil Information Service, the Soil Characterization  
 47 Database, the University of California, Riverside, and the United States Department of  
 48 Agriculture Agricultural Research Service. Posterior soil texture predictions achieved an  
 49 RMSE below 10%, a 7% reduction over priors, RMSE and spatial representation for soil  
 50 organic matter and bulk density also improved. Furthermore, the method reduced  
 51 prediction uncertainties, (narrower prediction intervals compared to the priors) and  
 52 enforced physical constraints on soil property bounds. Looking forward, this IRC method  
 53 offers a scalable pathway to improve existing probabilistic soil maps, providing a strategy  
 54 for the evolution of digital soil products as new soil observations emerge.

- Deleted: including agricultural
- Deleted: fertilization.
- Deleted: a hybrid
- Deleted: that combines prior soil predictions with
- Deleted: improve
- Deleted: to demonstrate its application. We
- Deleted: observation
- Deleted: It
- Deleted: statistical shape of
- Deleted: bias of
- Deleted: knowledge with observed
- Deleted: information. We
- Deleted: soil mapping
- Deleted: . For residual correction, we compiled laboratory-measured soil profile data
- Deleted: three primary sources:
- Deleted: (WoSIS),
- Deleted: National
- Deleted: (SCD), and field measurements conducted by
- Deleted: ,
- Deleted: (UCR)
- Deleted: USDA-ARS
- Deleted: Salinity Laboratory. From the evaluations, the posterior
- Deleted: show
- Deleted: of less than
- Deleted: compared to the
- Deleted: (pHRF-derived soil maps). For
- Deleted: (SOM)
- Deleted: oven-dry
- Deleted: (BD), the RMSE
- Deleted: decreased, as the priors initially underestimated their spatial variation. Although posterior SOM and BD predictions were less accurate than other soil properties, this was expected since they are dynamic soil properties and their response to environment and anthropogenic activities is more difficult to simulate. The residual correction also showed
- Deleted: , as demonstrated by
- Deleted: . This method also applied optimization with
- Deleted: , such as ensuring the bounds of soil prop... [1]
- Deleted: ¶ ... [2]

## 1 Introduction

Soils play an important role in regulating Earth's water, energy, and nutrient cycles (Vereecken et al., 2016). Soil maps guide agricultural practices, ecosystem management, hydraulic modeling, and climate studies, such as crop modeling, flood risk assessment, groundwater management, and climate change (Vereecken et al., 2022). The importance of soil maps has increased with the advent of precision agriculture, including site-specific seeding, irrigation, and fertilization recommendations that intrinsically depend on high-resolution soil properties (Jiang et al., 2011; Li et al., 2019; Mueller et al., 2001; Ortuani et al., 2016). However, the accuracy and reliability of these management actions heavily depend on the quality of soil maps as a critical decision-making input. Traditional soil surveys involve field observations, laboratory analyses, and expert interpretation, but are labor-intensive and expensive (Grunwald et al., 2011; Rossiter et al., 2022; Soil Survey Staff et al., 2023). These limitations have driven the development of digital soil mapping (DSM) techniques. DSM leverages decades of soil data collection and sharing, establishing quantitative models to generate georeferenced soil maps (McBratney et al., 2003).

Digital soil maps are typically derived from existing soil surveys, geostatistical models, machine learning, or hybrid approaches. Soil survey-based soil mapping method, which use low, high, and representative values to describe soil property distributions for each soil component (Soil Survey Staff et al., 2023). The method typically approximates each soil component as a triangular distribution (Chaney et al., 2016; Soil survey staff, 2023), potentially oversimplifying multi-modal distributions of soil properties in some cases

134 [\(Haghverdi et al., 2020; Nussbaum et al., 2023\)](#). Additionally, estimating soil properties  
135 [from synthetic sampling within a map unit could create artificial spatial patterns, adding](#)  
136 [noises into the mapping results \(Chaney et al., 2019\)](#). Developments such as Latin-  
137 [hypercube sampling and landscape adaptive covariance functions have improved the](#)  
138 [representation of spatial patterns of soil properties \(Minasny and McBratney, 2006\)](#). Yet,  
139 [soil survey-based approaches remain valuable particularly in areas where soil profile data](#)  
140 [is limited \(Nauman et al., 2024\)](#). Geostatistical models often require presumed  
141 [parameterization and are constrained by stationarity assumptions, which is difficult to](#)  
142 [apply in areas with insufficient field knowledge \(Oliver and Webster, 2014\)](#). To address  
143 [these challenges, non-parametric models, such as Random Forest, trained with hybridized](#)  
144 [soil data that combine soil surveys with georeferenced soil profiles show potentials in](#)  
145 [improving soil mapping, particularly for large-scale maps \(Chaney et al., 2019; Nauman et](#)  
146 [al., 2024\)](#).

147  
148 [Map of soil properties have been observed with bias compared to field observations in](#)  
149 [certain areas due to many factors \(Hengl et al., 2017; Powers et al., 2011\)](#). At the  
150 [measurement level, sampling methods may favor certain landscape positions or soil](#)  
151 [conditions, causing a clustered representation \(Ramcharan et al., 2018\)](#). In areas with  
152 [coarse sampling density, models trained on unrepresentative data are likely to deviate](#)  
153 [from actual observations \(Sharififar et al., 2019\)](#). Commonly used DSM models can show  
154 [bias. For example, Random Forest classifier favors the majority class \(Chen et al., 2004\)](#),  
155 [and Random Forest regressors struggle to capture extreme values \(Nauman et al., 2024\)](#).

156 Furthermore, certain areas may not be fully captured by the DSM model and the selected  
157 feature space, such as areas with complex glacial pattern, parent material transitions, and  
158 alluvial processes (unaddressed problem in SOLUS; SoilGrids 2.0; (Nauman et al., 2024;  
159 Poggio et al., 2021)). Model-based solutions include using ensemble models to enhance  
160 accuracy compared to a single model (Sylvain et al., 2021). Post-processing methods,  
161 such as regression kriging and bias-corrected decision trees, can also be used (Hengl et  
162 al., 2004). Yet, kriging-based methods have limitations in areas with high spatial  
163 heterogeneity and abrupt transitions, where stationary assumptions do not meet. Non-  
164 parametric models can be used for bias correction that overcome the limitation of making  
165 presumed distributions.

166  
167 Quantifying uncertainties in DSM is important for its practical applications (Schmidinger  
168 and Heuvelink, 2023). DSM products represent soil properties as multi-dimensional  
169 matrices showing vertical and horizontal soil variation (Vereecken et al., 2022), with each  
170 pixel containing weighted possible values and their prediction uncertainties. These  
171 uncertainties can be represented either as continuous values through prediction intervals  
172 or as discrete classifications with associated class probabilities (Chaney et al., 2016,  
173 2019; Hengl et al., 2017; Ramcharan et al., 2018). Common quantification approaches  
174 include geostatistical techniques like kriging, where the nugget term accounts for  
175 measurement errors while kriging variance reflects spatial uncertainty patterns (Chilès and  
176 Delfiner, 2012; Takoutsing et al., 2022), and machine learning methods such as Quantile  
177 Random Forest (QRF) which generates probability distributions from decision tree outputs

178 [using values of soil properties \(Poggio et al., 2021; Shi et al., 2024\)](#). For discrete  
179 [classifications, uncertainty derives from soil raster probabilities during soil taxa](#)  
180 [classification \(Chaney et al., 2016; Odgers et al., 2015\)](#). Given the data-driven nature of  
181 [DSM and frequent limitations in soil profile availability, integrating multiple qualified data](#)  
182 [sources improves the amount of soil data and reduce prediction uncertainties \(Nauman et](#)  
183 [al., 2024\)](#), particularly in regions where predictions must rely more heavily on legacy soil  
184 [data](#).

185  
186 In this study, we present a hybrid DSM approach combining pruned Hierarchical Random

187 Forest (pHRF) [with iterative residual correction \(IRC\) method \(Xu et al., 2025\)](#). The pHRF  
188 [method leverages the National Cooperative Soil Survey \(NCSS\)](#) soil survey data and  
189 georeferenced soil taxa information to generate prior distributions, while additional soil

190 profiles correct biases in [prior](#) predictions. This method builds on development in previous

191 research while addressing specific limitations. Sylvain et al. (2021) applied XGBoost

192 (sequential decision trees) and ensemble models to correct deterministic soil property

193 maps, demonstrating reduced bias for many soil properties [\(Sylvain et al., 2021\)](#). Zhang et

194 al. (2010) introduced a bias-correction technique with Random Forest models to mitigate

195 their tendency to regress toward mean values, though not in DSM contexts [\(Zhang and Lu,](#)

196 [2012\)](#). Our approach extends these concepts by probabilistically updating posterior

197 distributions at each location through an iterative correction process that continues until

198 convergence across vertical intervals. Vertical correlations are maintained through layer-

199 by-layer residual correction, which preserves inter-layer correlations while dynamically

Deleted: predictions with residual correction. The pHRF method leverages NCSS...

Deleted: posterior

Deleted: (Sylvain et al., 2021)

Deleted: (Zhang and Lu, 2012)

205 optimizing the feature space at each correction step. Unlike methods requiring  
206 distributional assumptions, our non-parametric framework adapts to diverse landscapes  
207 and data scenarios. The models implement residual correction by minimizing the  
208 differences between priors and new observations to adjust posterior distributions, with the  
209 entire process continuing until property variations stabilize between different iterations.  
210 This method aims to improve the accuracy and reliability of soil property maps, supporting  
211 decision-making in relevant applications.

Deleted: observed

## 213 2 Methods

214 This study introduces a hybrid framework for digital soil mapping (DSM) that updates  
215 existing probabilistic soil property maps using newly collected soil observations. The  
216 framework combines prior soil property estimates with an iterative residual correction  
217 (IRC) method. The IRC method integrates additional georeferenced soil profiles (soil  
218 observations not used to train prior soil maps) and employs non-parametric models to  
219 adjust the distribution of prior estimates, thereby correcting biases in the prior soil maps.

Deleted: approach

Deleted: , combining

Deleted: derived from the pruned hierarchical Random Forest (pHRF) method followed

Deleted: ).

Deleted: soil profiles

219 adjust the distribution of prior estimates, thereby correcting biases in the prior soil maps.  
220

Deleted: and improving

Deleted: accuracy of

Deleted: property predictions.

221 The following sections first describe the general residual correction framework (Section  
222 2.1). To illustrate the method concretely, we then provide a worked example using one

Deleted: workflow of the pHRF

Deleted: , the steps

223 randomly selected soil column to demonstrate how the feature space is constructed and  
224 updated across two consecutive iterations (Section 2.1.1). Building on this example, we  
225 detail the key components of the IRC method: the iterative update of feature space  
226 (Section 2.1.2), the convergence criteria for residual correction (Section 2.1.3), and the

Deleted: ,

240 ~~process for updating posterior soil properties with physical constraints (Section 2.1.4).~~

Deleted: generation of updated

241 ~~Finally, we present the California case study (Section 2.2), describing the soil datasets~~

Deleted: property maps.

242 ~~used (Section 2.2.1) and the implementation details for applying the IRC method over~~

243 ~~California (Section 2.2.2).~~

244

## 245 ~~2.1 Iterative Residual Correction Framework for DSM~~

Deleted: 2.1 Soil Data  
To correct ...

246 ~~Residual correction is implemented to address underestimated soil property variation in~~

247 ~~prior maps (tendency to underestimate high values and overestimate low values,~~

248 ~~smoothing out soil variation across landscape). The overall workflow of the IRC method~~

249 ~~consists of three components: (1) prior map generation (Figure 1a), (2) residual preparation~~

250 ~~(Figure 1b), and (3) iterative correction (Figure 1c).~~

251

252 ~~First, probabilistic prior soil property maps are generated or retrieve probabilistic soil~~

253 ~~property maps from an existing DSM product as the prior soil maps (Figure 1a). These~~

254 ~~maps represent the initial estimates of soil properties and their associated uncertainties.~~

255 ~~Second, a residual preparation step is carried out to enable correction using new soil~~

256 ~~profile observations (Figure 1b). The preparation involves four key steps: (1) adding~~

257 ~~additional soil profiles from new field measurements or databases; (2) spatially aligning~~

258 ~~these profiles with the corresponding pixels in the prior soil maps using geographic~~

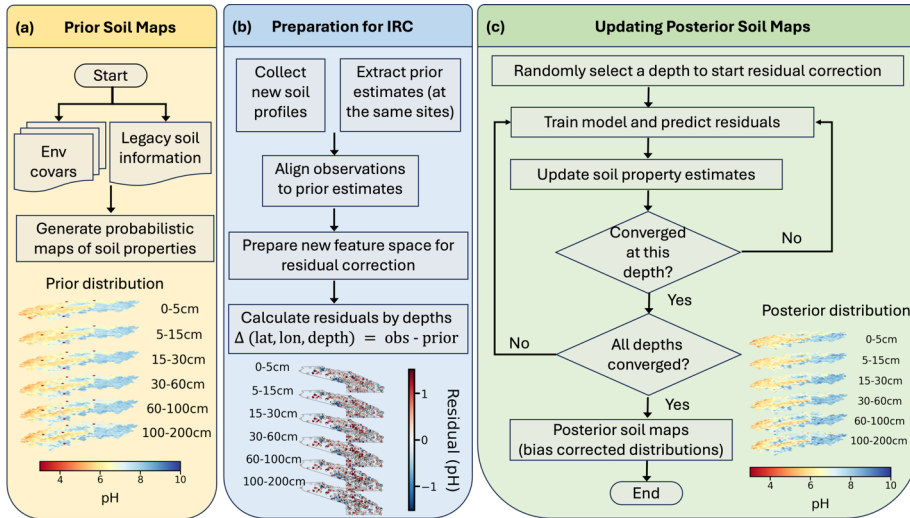
259 ~~coordinates; (3) vertically aligning observations with prior predictions at matching depth~~

260 ~~intervals; and (4) calculating residuals depth by depth as the difference between observed~~

261 ~~values and prior predictions. During this stage, the feature space for residual modeling is~~

266 also prepared, consisting of static environmental covariates (which remain fixed  
267 throughout iterations) and dynamic soil covariates (which are updated iteratively). Detailed  
268 construction of the feature space is described in Section 2.1.1.  
269  
270 Finally, iterative residual correction is performed to update soil property estimates across  
271 depths (Figure 1c). During each iteration, the model predicts residuals for one depth layer  
272 at a time, with the layer selected randomly. A Random Forest regressor is trained to learn  
273 the relationship between residuals and the feature space at sampled locations, then  
274 interpolates residual corrections across the study area. Predicted residuals are added to  
275 the prior (or previous iteration's) estimates to generate updated soil property values. After  
276 each update, convergence is evaluated for the modeling depth by comparing the median  
277 difference between the current residuals and those from the previous iteration. Once this  
278 change falls below a predefined threshold, that depth is considered converged and  
279 excluded from subsequent updates. The algorithm then focuses on the remaining  
280 “unconverged” depths, until convergence is achieved across all layers. After convergence  
281 is verified for all depths, the final corrected residuals are added to the prior estimates to  
282 update the posterior distributions of soil properties.

283



284  
 285 Figure 1: Workflow for updating posterior soil property maps. The process begins with  
 286 panel (a), the preparation of environmental covariates (env covars) to generate  
 287 probabilistic maps of soil properties (prior soil maps). As illustrated in panel (b), the  
 288 preparation for residual correction involves adding additional soil profiles, spatially and  
 289 vertically aligning prior soil map values with new profile observations, calculating residuals  
 290 depth by depth, and preparing environmental covariates and soil covariates (new feature  
 291 space) for residual correction. Finally, as shown in panel (c), the iterative residual  
 292 correction step applies bias corrections across different depths, focusing on layers where  
 293 residuals have not yet stabilized. During each iteration, the model predicts residuals for  
 294 one depth at a time, randomly selecting a layer. Once residuals for a given depth converge,  
 295 that layer is excluded from further updates, allowing the model to concentrate on  
 296 remaining depths until all achieve stability. After verifying convergence across all depths,

Formatted: Font: Not Bold

297 the algorithm updates the posterior distribution of soil properties and produces the final  
298 soil maps (posterior soil property maps).

299

300 In this IRC framework, "prior probabilistic soil property maps" refer to spatially continuous  
301 soil property maps that provide an initial (prior) estimate of soil properties with associated  
302 uncertainty across the study area. These prior maps provide, for each pixel and depth  
303 interval, a distribution of possible soil property values with associated probabilities or  
304 weights. The IRC method does not require prior and new soil observations to be co-located  
305 at the same pixels. Instead, the method requires that a prior estimate exists at locations  
306 where new soil observations are available. By learning the relationship between residuals  
307 (differences between new observations and prior estimates) and environmental and soil  
308 covariates at sampled locations, the trained model can interpolate residual corrections  
309 across the study area.

310

### 311 **2.1.1 Worked Example**

312 The iterative residual correction method is further illustrated in Figure 2 using an example  
313 with a randomly selected soil column. Figure 2a shows the location of the selected soil  
314 column, where additional soil profile observations are available. The right panel displays  
315 the top-3 probable pH values (from prior soil maps) at each depth intervals (0–5 cm, 5–15  
316 cm, 15–30 cm, 30–60 cm, 60–100 cm, 100–200 cm), while the left panel shows the three  
317 weights (probabilities) associated with these pH values. In this simplified example, we use  
318 3 bins to represent the soil property distribution; however, in actual implementation, more

319 bins are maintained (typically top-12 probable values) to better capture soil variability. For  
320 this demonstration, Depth 2 (D<sub>2</sub>; 5–15 cm) is randomly selected as the modeling layer to  
321 initiate the iterative correction process. Only one layer is modeled and updated for a given  
322 iteration. Note that in real model execution, model generally processes over 3,000 soil  
323 columns simultaneously in California, though only one column is shown here for clarity.

324

325 In Figure 2b, the table details features used to train the Random Forest regressor for  
326 residual prediction. The feature space consists of environmental covariates that remain  
327 fixed across iterations and soil covariates that are updated iteratively:

328 (1) Environmental covariates (21 dimensions): These capture spatial variations in  
329 soil-forming factors and remain unchanged throughout all iterations. The covariates  
330 include remote sensing data (Sentinel-1, Sentinel-2, GOES land surface  
331 temperature) and terrain attributes, identical to those used in the prior mapping  
332 method (Xu et al., 2025).

333 (2) Depth information (1 dimension): The centroid (median value) of the soil depth  
334 interval for the modeling layer (e.g., 10 cm for the 5–15 cm layer), describing the  
335 vertical position in the soil profile.

336 (3) Representative soil property values (1 dimension): The expected value (weighted  
337 mean) of the soil property at each pixel in the modeling layer, representing the  
338 current best estimate. This is computed as the weighted sum of top-probable  
339 values.

340 (4) Top-probable soil property values (1 dimension): The current predictions at each  
341 pixel (residuals plus previous prediction of soil property values), reflecting both  
342 intra-pixel and inter-pixel soil heterogeneity.

343 (5) Inter-layer differences (5 dimensions): Differences in top-probable predicted soil  
344 property values between the modeling layer and the other five depth layers. For  
345 instance, if modeling Depth 2, the inter-layer differences would be  $(D_2-D_1)$ ,  $(D_2-D_3)$ ,  
346  $(D_2-D_4)$ ,  $(D_2-D_5)$ , and  $(D_2-D_6)$ . These features capture vertical correlations in the soil  
347 profile and aid in estimating spatial patterns.

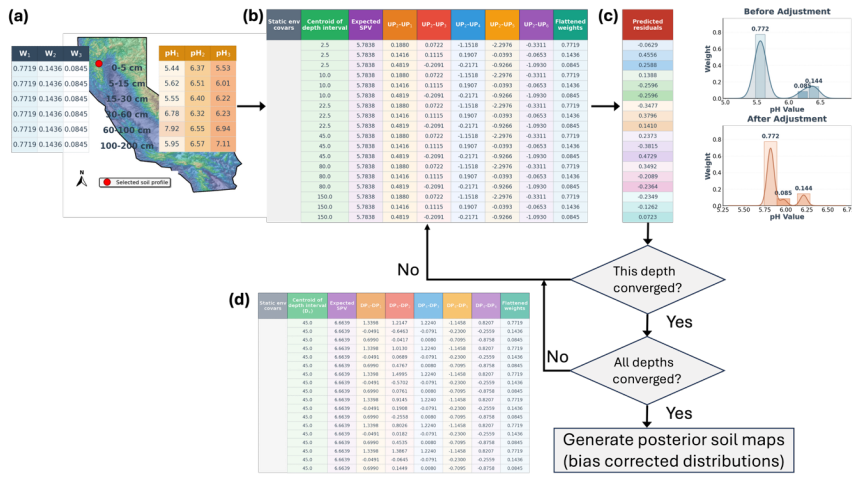
348 (6) Weights (1 dimension): Probabilities associated with each top-probable soil  
349 property value. These weights remain fixed throughout iterations.

350  
351 In summary, environmental covariates and weights remain static, while depth information,  
352 representative values, top-probable values, and inter-layer differences are updated across  
353 iterations based on the most recent soil property estimates.

354  
355 A Random Forest regressor is then trained using the feature space to predict residuals for  
356 the modeling layer ( $D_2$  in this example). The right panel in Figure 2c compares the  
357 distribution of pH values before and after residual adjustment in the current iteration. After  
358 applying the residual correction, convergence is checked for  $D_2$  by comparing the median  
359 difference between the current and previous residuals. If  $D_2$  has converged (difference  
360 below threshold), the algorithm proceeds to check whether all depth layers have  
361 converged. If all layers have converged, the iterative process terminates, and the final

362 posterior soil property maps are generated by adding the last predicted residuals to the  
363 prior values.  
364  
365 If either convergence check returns "No" (i.e.,  $D_2$  has not converged or other layers remain  
366 unconverged), the algorithm continues iterating. Here, the soil property values for  $D_2$  are  
367 updated by adding the predicted residuals to the previous pH values. These updated  
368 values are then used to reconstruct the feature space following the same structure  
369 described above, updating the representative values, top-probable values, and inter-layer  
370 differences. By updating soil covariates layer by layer and iteratively refining the feature  
371 space, the next prediction retains prior knowledge while integrating new information about  
372 soil heterogeneity and vertical relationships for soil profiles (Wu et al., 2025). A new  
373 iteration begins by randomly selecting another unconverged layer, and the process repeats  
374 until convergence is achieved across all depth layers.

375



376

377 **Figure 2: Schematic illustration of the iterative residual correction (IRC) method using a**  
 378 **worked example at a randomly selected soil column. (a) Prior distributions and**  
 379 **observation location: The map shows the location of the selected soil column within the**  
 380 **study area. The right panel displays the top-3 probable pH values at each of the six depth**  
 381 **intervals (0–5 cm, 5–15 cm, 15–30 cm, 30–60 cm, 60–100 cm, 100–200 cm), while the left**  
 382 **panel shows the three weights ( $w_1, w_2, w_3$ ) associated with these pH values. Depth 2 ( $D_2$ : 5–**  
 383 **15 cm) is randomly selected for this iteration. (b) Feature space components: The table**  
 384 **details the structure of the feature space used to train the Random Forest regressor for**  
 385 **residual prediction. The feature space comprises both static and dynamic components.**  
 386 **Static components include environmental covariates (satellite imagery, terrain attributes)**  
 387 **that remain unchanged throughout iterations, and weights ( $w_1, w_2, w_3$ ) associated with top-**  
 388 **probable values. Dynamic soil covariates that are updated in each iteration include: the**  
 389 **centroid of the depth interval (e.g., 10 cm for  $D_2$ ), the expected (representative) soil**

390 property value computed as the weighted mean, the top-probable soil property values  
391 reflecting intra-pixel heterogeneity, and inter-layer differences capturing vertical  
392 correlations (e.g.,  $D_2-D_1$ ,  $D_2-D_3$ ). (c) Residual correction and convergence workflow: A  
393 Random Forest model trained on the feature space predicts residuals for the modeling  
394 layer  $D_2$ . The right panel compares the pH distribution before and after residual  
395 adjustment. The flowchart below describes the convergence logic: after predicting and  
396 applying residuals to  $D_2$ , the algorithm evaluates whether  $D_2$  has converged. If  $D_2$  has  
397 converged, the algorithm checks whether all depth layers have achieved convergence. If  
398 both checks pass, the final posterior soil property maps are generated by adding the last  
399 converged residuals to the prior values. (d) If either check fails, the algorithm updates the  
400 soil property values for  $D_2$  by adding predicted residuals, reconstructs the feature space  
401 with the updated values, randomly selects another unconverged layer, and repeats the  
402 process. This iterative cycle continues until convergence is achieved across all six depth  
403 layers.

### 405 2.1.2 Convergence of Residual Correction

406 The residual correction process continues until the median difference between updated  
407 residuals and previous residuals falls below a predefined threshold. Convergence is  
408 achieved when the residuals stabilize across multiple iterations, indicating that further  
409 adjustments do not largely change the predictions. This stability ensures that the final  
410 posterior soil properties are reliable and consistent. The stopping criterion is a  
411 customizable parameter. In this work, it was set to the 5th percentile of the distribution of

Moved (insertion) [2]

Moved (insertion) [3]

412 value changes. To avoid over-correcting bias, only the last converged residuals are added  
413 to the prior prediction to generate the final posterior results.

414

### 415 **2.1.3 Update with Constraints**

416 During residual correction, a common issue arises where the addition of residuals to prior  
417 soil property values results in values that exceed physical bounds (such as sand content  $\geq$   
418 100%). To address this, a residual update process with constraints is implemented.

Moved (insertion) [4]

Deleted: as post processing, we only

419

420 As illustrated in Figure 2c to 2d, after the Random Forest regressor predicts residuals for  
421 the layer ( $D_2$ ), these residuals are added to the previous soil property values to generate  
422 updated predictions. Immediately after this addition step, the updated values are  
423 examined to check whether they fall within predefined physical bounds (e.g., 0% to 100%  
424 for particle size fractions, positive values for bulk density). This constraint check occurs  
425 before the convergence evaluation and before the updated values are used to reconstruct  
426 the feature space for the next iteration.

427

428 If any updated value exceeds the physical bounds, it is adjusted to the nearest valid bound  
429 (minimum or maximum). For example, if adding a residual of +15% to a prior sand content  
430 of 90% yields 105%, this value is capped at 100%. The "excess" residual (+5% in this case)  
431 is then redistributed proportionally (based on their weights) among the other top-probable  
432 values at the same pixel, ensuring that the total correction remains consistent with the  
433 model's prediction while maintaining physical plausibility. For particle size fractions (sand,

435 silt, clay), an additional compositional constraint ensures that the three fractions sum to  
436 100% at each pixel after residual correction.

437

## 438 **2.2 California Case Study: Soil Data and Model Implementation**

### 439 **2.2.1 Soil Data**

440 To demonstrate the IRC method, we apply it to soil property mapping in California. We use  
441 georeferenced soil profiles with laboratory measurements of soil properties. We compiled  
442 soil profile data from three primary sources: the World Soil Information Service (WoSIS),  
443 the National Soil Characterization Database (SCD), and field measurements conducted in  
444 California (Batjes et al., 2024; National Cooperative Soil Survey, 2018; Scudiero et al.,  
445 2024).

446

447 To ensure consistency across different data sources, we applied several quality control  
448 steps. First, we checked the physical plausibility of all soil property values by defining a  
449 valid range with specific minimum and maximum thresholds for each property. Any data  
450 point falling outside these ranges was considered an error and removed. For soil texture,  
451 we required the sum of sand, silt, and clay fractions to equal 100%. If a profile did not meet  
452 this compositional constraint, it was excluded. After quality check, the datasets are  
453 compatible because the WoSIS records for California are largely derived from the NCSS  
454 database, and both the SCD and WoSIS datasets follow standardized laboratory protocols,  
455 such as those from the Kellogg Soil Survey Laboratory (Soil, 1996; Soil Survey Staff, 2014).  
456 For our own field measurements, we used the Integral Suspension Pressure (ISP+) method

Deleted: .

458 to maintain precision for particle size analysis (Corwin and Scudiero, 2020; Scudiero et al.,  
459 2024).

460

461 During preprocessing, we harmonized all soil data, which was originally reported at  
462 different soil horizons, into six standard depth intervals: 0–5 cm, 5–15 cm, 15–30 cm, 30–  
463 60 cm, 60–100 cm, and 100–200 cm (Arrouays et al., 2014). The harmonization was  
464 performed using equal-area spline functions to interpolate soil property values from the  
465 original horizon depths to these standard intervals (Hartemink et al., 2010, p.201). The  
466 spline function fits a smooth curve through observed values at their measured depths,  
467 then calculates the area under this curve within each standardized depth interval and  
468 divides by the interval width to obtain the value. Location of soil profiles and their  
469 distribution of soil property values are presented in Figure 3. Six soil properties are studied;

470 sand content, silt content, clay content, pH, soil organic matter (log-scaled), and oven-dry  
471 bulk density. These samples were not co-located with the training samples used to  
472 generate the prior maps (samples at the same locations were already removed). The  
473 number of observations varies by soil property: pH has the most samples, followed by  
474 oven-dry bulk density and soil organic matter. The sample sizes across properties can also  
475 be inferred from the frequency histograms shown in the lower-left corner of each panel in  
476 Figure 3. Across all depths combined, each soil property has more than 11000  
477 observations in California. The number of observations generally decreases with depth,  
478 with depths below 1 m having notably fewer samples compared to shallower layers.

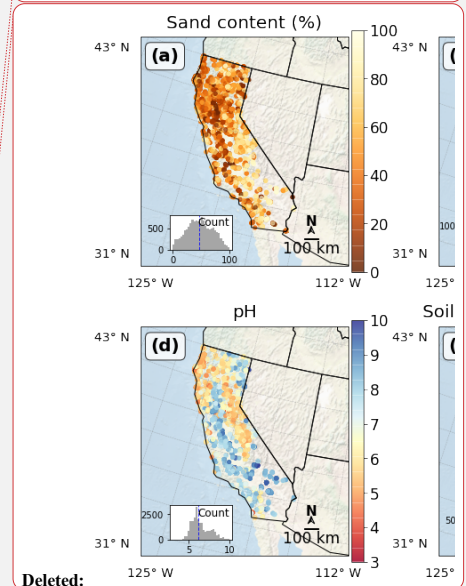
479

Deleted: standardized depth intervals.

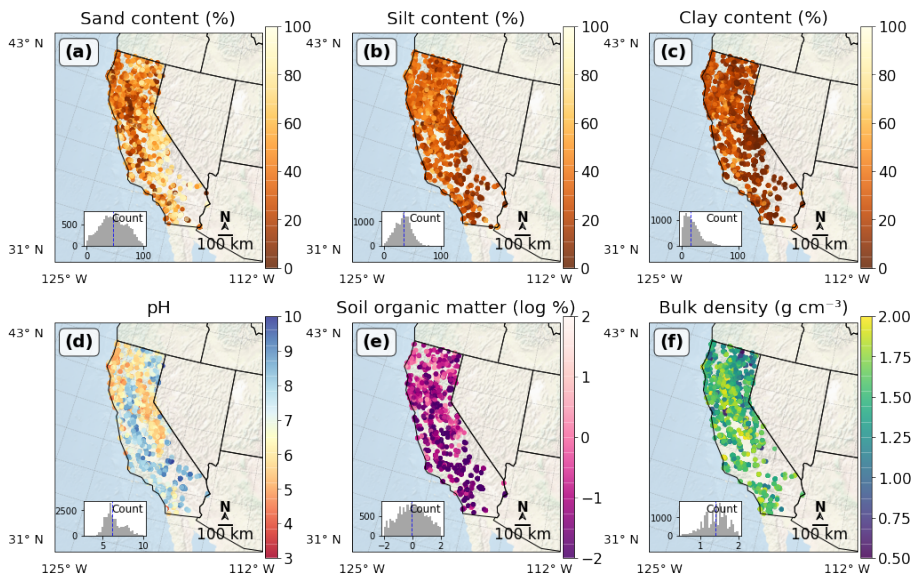
Deleted: 1

Deleted: —sand content, silt content, clay content, pH, soil organic matter (log-scaled), and bulk density

Formatted: Font: SimSun



Deleted:



485

486 **Figure 3:** Spatial distribution and statistical characteristics of soil properties observations

487 across California. The figure presents six soil parameters mapped using an Albers Equal

488 Area projection: (a) sand content (mass %), (b) silt content (mass %), (c) clay content

489 (mass %), (d) pH, (e) soil organic matter (log-scaled mass %), and (f) bulk density (g/cm<sup>3</sup>).

490 Each subplot displays sample locations as colored points, with field-collected samples

491 shown as triangles to distinguish them from WoSIS (circles) and SCD (squares) samples.

492 Distribution histograms in the lower left corner of each subplot show the frequency

493 distribution of values, with blue dashed lines indicating median values. Distance scale bar

494 and compass rose are provided in the right corner. Note that the total number of soil

495 measurements varies by property and generally decreases with depth beyond the surface

496 layer, with the surface layers and depths below 1 m generally having fewer observations.

Deleted: 1

Formatted: Font: Not Bold

Formatted: Font: Not Bold

Deleted: .

Formatted: Font: Not Bold

Formatted: Font: Not Bold

499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520

**2.2.1.1 World Soil Information Service (WoSIS)**  
The World Soil Information Service (WoSIS), managed by the International Soil Reference and Information Centre (ISRIC), aggregates global soil data from diverse sources, including national soil institutes, research organizations, and collaborative initiatives like the Global Soil Partnership (GSP) and the International Network of Soil Information Institutions (INSII). The database provides soil properties for different soil horizons, georeferenced in decimal degrees, and undergoes quality controls (Batjes et al., 2024). In California, WoSIS typically offers 2,000 to over 5,000 soil observations for the modeling soil property. Samples below 1-m depth are fewer than those from shallower layers.

**2.2.1.2 Soil Characterization Database (SCD)**

The Soil Characterization Database (SCD) is a subset of the National Cooperative Soil Survey (NCSS) database (National Cooperative Soil Survey, 2018). It records soil properties for each soil horizon within a soil profile (pedon), including soil texture, bulk density, and water retention. In California, SCD provides between 500 and over 1,000 soil samples per layer for the studied soil property. Each soil profile is georeferenced and includes metadata such as site location, land use, and sampling methods.

**2.2.1.3 Ground Truth Soil Sampling and Measurements**

Additional soil sampling was conducted to complement georeferenced soil profiles in California for model training and evaluation. These data are reported in (Scudiero et al.,

Formatted: Font color: Accent 2

Deleted: 2.1.1 World Soil Information Service (WoSIS)  
The World Soil Information Service (WoSIS), managed by the International Soil Reference and Information Centre (ISRIC), aggregates global soil data from diverse sources, including national soil institutes, research organizations, and collaborative initiatives like the Global Soil Partnership (GSP) and the International Network of Soil Information Institutions (INSII). The database provides soil properties for multiple depth intervals, georeferenced in decimal degrees (WGS84), and undergoes quality controls (Batjes et al., 2024). In California, WoSIS typically offers 2,000 to over 5,000 soil observations for soil property of interest. Samples below 1-meter depth are fewer than those from shallower layers.

Formatted: Heading 4

Deleted: (National Cooperative Soil Survey, 2018).

Deleted: The data are collected using standardized laboratory methods...

Deleted: properties

Deleted: (WGS84)

Formatted: Font color: Accent 2

Deleted: truth soil sampling

Deleted: measurements

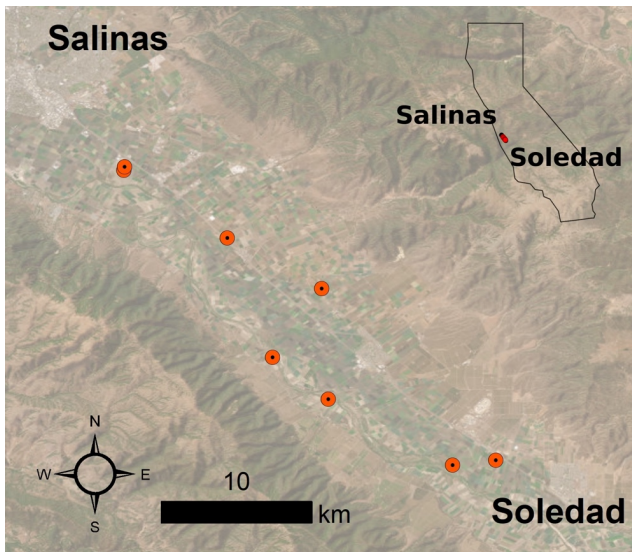
Formatted: Font color: Text 1

Formatted: Heading 4

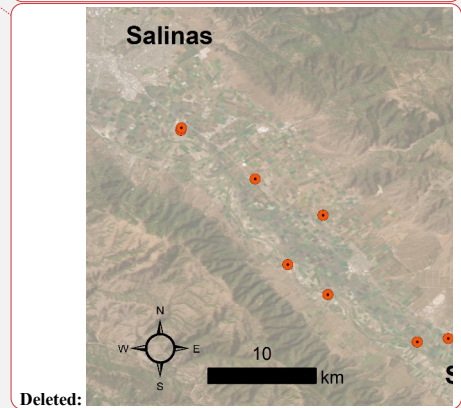
Formatted: Font color: Text 1

Deleted: Scudiero et al. (2024)

544 2024) and are briefly discussed here. Multiple fields located between Salinas and Soledad  
 545 in California's Salinas Valley were selected to collect soil particle size fraction data (Figure  
 546 4). These fields, presented as red dots in Figure 4, were chosen because they were  
 547 accessible, unfarmed during the sampling period, and spread across different parts of the  
 548 valley.



Deleted:  
 Formatted: Font color: Text 1  
 Formatted: Font color: Text 1  
 Deleted: California's  
 Formatted: Font color: Text 1  
 Deleted: 2  
 Deleted: 2  
 Formatted: Font color: Text 1  
 Formatted: Font color: Text 1  
 Deleted:  
 Formatted: Font color: Text 1



550  
 551 Figure 4: Map of sampling fields in the Salinas Valley in California. Each red dot represents  
 552 a sampling field between Salinas and Soledad. An inset map (top right) shows the location  
 553 of the sampling area within California. Scale bar and direction indicator are provided in the  
 554 left corner. Basemap: Esri World Imagery. Source: Esri, Maxar, Earthstar Geographics, and  
 555 the GIS User Community.

Deleted: 2  
 Formatted: Font: Not Bold, Font color: Text 1  
 Formatted: Font: Not Bold, Font color: Text 1  
 Formatted: Font: Not Bold, Font color: Text 1  
 Deleted: Environmental Systems Research Institute (  
 Deleted: )  
 Deleted: , with imagery and data provided by  
 Formatted: Font: Not Bold, Font color: Text 1  
 Formatted: Font: Not Bold, Font color: Text 1  
 Formatted: Font: Not Bold, Font color: Text 1  
 Deleted: ¶

577 [Soil apparent electrical conductivity \(ECa\) was measured across fields using an](#)  
578 [electromagnetic induction \(EMI\) sensor connected to a GPS receiver. Following the ECa-](#)  
579 [directed soil sampling protocols of Corwin and Scudiero \(Corwin and Scudiero, 2020\), the](#)  
580 [most representative soil samples were identified with ESAP software package and the](#)  
581 [Response Surface Sampling Design algorithm \(Lesch et al., 2000; Lesch, 2005\). 0-0.8 and](#)  
582 [0-1.6 m soil profiles were further analyzed and followed with the expectation that ECa was](#)  
583 [a regional proxy for the field-scale variability of particle size fraction.](#)

584

585 [To measure particle size fraction, soil samples were then collected from multiple depths](#)  
586 [\(0-0.1, 0.1-0.4, and 0.4-1.2 m\) across fields. After collection, the samples were air-dried,](#)  
587 [ground, and sieved to remove particles larger than 2 mm; and then measured using the](#)  
588 [Integral Suspension Pressure method \(The improved integral suspension pressure method](#)  
589 [\(ISP+\) for precise particle size analysis of soil and sedimentary materials; Wolfgang Durner,](#)  
590 [Sascha C. Iden\) using PARIO™ system \(METER Group AG, Munich, Germany\).](#)

591

592 **2.2.2 Model Implementation for the California Case Study**

593 [For the California case study, prior soil property maps were generated using the pruned](#)  
594 [hierarchical Random Forest \(pHRF\) method \(Xu et al., 2025\). The pHRF-derived soil maps](#)  
595 [were developed with soil pedons from the National Soil Information System \(NASIS\) and](#)  
596 [part of SCD \(the remaining data not used in IRC method\). After gaining prior estimate of soil](#)  
597 [properties, the IRC method was then applied using the additional soil observations from](#)  
598 [WoSIS, SCD, and field measurements, which were not used in generating the prior maps.](#)

Deleted: . ¶

### 2.2.1 Prior Soil Properties ¶

Prior soil properties maps were generated using a pruned hierarchical Random Forest (pHRF) methodology (Xu et al., 2025). Figure 3A illustrates the workflow of the pHRF method. The DSM method begins by integrating soil covariates and soil pedons with taxonomic names to generate probabilistic maps of soil classes. These maps are then linked to a harmonized soil properties database (Chaney et al., 2019), which estimates the distribution of soil properties linked to each soil component. By combining these inputs, the pHRF method produces probabilistic maps of soil properties, serving as the prior distributions for subsequent residual correction. The pHRF method implements several key features: (1) it efficiently incorporates soil covariates such as Sentinel-1 and Sentinel-2 satellite data, GOES land surface temperature, to capture detailed land heterogeneity; (2) it uses a "moving polygon" algorithm to preserve natural landscape boundaries, ensuring spatial consistency; (3) it integrates soil pedons and soil surveys with soil properties estimates (harmonized soil properties database) to increase the availability of soil information; and (4) it employs a hierarchical structure in soil classification and pruned less plausible prediction that sharpens the prediction interval of prediction and reduces uncertainties in soil property estimates (Xu et al., 2025). The method addresses data imbalances, such as unevenly distributed soil observations and underrepresented soil classes. While the pHRF-derived soil property maps have demonstrated effectiveness in reducing uncertainties, certain properties, such as bulk density, still exhibit bias (Xu et al., 2025). This shows the need for further calibration. The pHRF method produces probabilistic maps of soil properties. Each pixel contains prior distribution of soil property values and their weights. ¶

### 2.2.2 Updating Posterior Soil Properties Maps ¶

Updating the posterior soil properties maps involves correcting prior soil property estimates by incorporating additional soil profiles and correcting the residuals (the differences between observed values and prior predictions). The process begins with the preparation for residual correction (Figure 3B) — calculating residuals between additional soil profiles and co-located prior soil data depth by depth. By adding these residuals to the prior distributions, the statistical shape of the probability distribution is adjusted (updated property; UP). Non-parametric model Random Forest regressors are selected for the adjustments, as they can flexibly adapt to changes in the distribution shape without relying on predefined assumptions. Additionally, soil covariates are prepared for residual correction at each depth as feature space for predictive models. ¶

The iterative residual correction method is further ... [4]

Formatted: Font: Not Bold

The convergence threshold for each soil property was set to the 5th percentile of the distribution of value changes between iterations.

Model training and evaluation were performed using out-of-bag (OOB) sampling, with OOB samples (samples withheld from the training process and not used to fit the models) that shared the same geolocation as training samples removed to prevent data leakage and reduce spatial autocorrelation effects. In each iteration, a new Random Forest model is trained to update residuals for one specific depth interval, and the same set of OOB samples remains excluded throughout to ensure independent validation.

### 3 Results

The iterative residual correction (IRC) method is applied to adjust pHRF-derived prior soil properties, including particle size fractions (sand, silt, clay), pH, oven-dry bulk density (BD), and soil organic matter (SOM) over California. This correction addresses biases in the prior soil property maps and updates the posterior distributions of these properties. These soil properties are important for land management and serve as essential inputs for pedotransfer functions. The residual correction is performed across California, covering six depth intervals: 0-5 cm, 5-15 cm, 15-30 cm, 30-60 cm, 60-100 cm, and 100-200 cm.

#### 3.1 Performance Evaluation of Posterior Soil Properties

Table 1 presents the performance metrics for the posterior predictions of six key soil properties: sand, silt, clay, pH, oven-dry bulk density (BD), and soil organic matter (SOM).

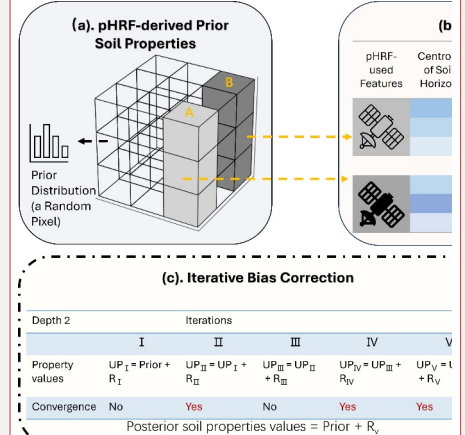
**Moved up [2]:** The residual correction process continues until the median difference between updated residuals and previous residuals falls below a predefined threshold. Convergence is achieved when the residuals stabilize across multiple iterations, indicating that further adjustments do not

**Moved up [3]:** This stability ensures that the final posterior soil properties are reliable and consistent.

**Moved up [4]:** During residual correction, a common issue arises where the addition of residuals to prior soil property values results in values that exceed physical bounds (such as sand content > 100%).

**Deleted:** significantly change the predictions.

**Deleted:** The stopping criteria is a customizable parameter. In this work, it was set to be 5<sup>th</sup> percentile of distribution of value changes. To avoid over-correcting bias (overfitting), only the last converged residuals are added to the prior prediction to generate the final posterior results. This method also addresses evaluation bias by achieving convergence across multiple iterations.



**Figure 4: Schematic workflow of iterative residual correction (IRC) for Soil Properties.** The workflow have three main components: (a) prior soil properties derived from the pruned hierarchical Random Forest (pHRF) method, (b) iterative optimization of the feature space, where  $W_1, W_2, W_3$  represent weights assigned to soil properties at each pixel;  $R_1, R_2, R_3$ ...

**Deleted:** To address this, an optimization process with constraints is implemented. During the first iteration, residuals are first computed without constraints. The updated soil property values are then examined to ensure they fall within predefined bounds (such as 0% to 100% for particle size fractions). If a value exceeds the bounds, it is adjusted to the nearest bound (minimum... [6])

**Deleted:** ).

832 The metrics include the root mean square error (RMSE), coefficient of determination ( $R^2$ ),  
833 and correlation coefficient ( $\rho$ ). For example, sand prediction shows an RMSE of 9.322, an  
834  $R^2$  of 0.841, and a correlation coefficient of 0.918. pH prediction shows an RMSE of 0.270,  
835 an  $R^2$  of 0.945, and a correlation coefficient of 0.972. These metrics are computed using  
836 out-of-bag (OOB) samples from random forest regressors. OOB samples are data points  
837 not included in the bootstrap samples used to train each tree in the random forest.  
838 Additionally, these metrics are evaluated by comparing the expected values of posterior  
839 predictions with co-located soil properties values; not computed on residuals.

840

841 Table 1 also shows variations in performance across different soil properties. SOM and  
842 bulk density show slightly worse metrics compared to particle size fractions and pH. For  
843 instance, SOM predictions have an RMSE of 1.961, an  $R^2$  of 0.608, and a correlation  
844 coefficient of 0.801, and bulk density predictions have an RMSE of 0.164, an  $R^2$  of 0.704,  
845 and a correlation coefficient of 0.843. Two main reasons can result in their lower  
846 performance. First, these properties are more dynamic in nature compared to particle size  
847 fractions and pH. SOM and bulk density can change over time due to factors such as land  
848 use practices. The prior predictions are trained using soil survey data that are older, while  
849 the posterior soil profiles used for evaluation may come from a different period. Second,  
850 SOM and bulk density are more challenging to model accurately. SOM is influenced by  
851 complex biological and soil-forming processes, such as decomposition rates and organic  
852 matter inputs. Similarly, bulk density is affected by soil compaction, organic matter  
853 content, and soil structure. All of them can vary spatially and temporally. [Depth-wise](#)

854 [analysis of model performance is provided in the Supplementary Information \(Table S1 and](#)  
855 [S2\).](#)

856

857 **Table 1: Performance metrics (RMSE,  $R^2$ , and correlation coefficient  $\rho$ ) for posterior**  
858 **predictions of soil properties, including sand, silt, clay, pH, oven-dry bulk density**  
859 **(BD), and soil organic matter (SOM). The table summarizes the range (minimum and**  
860 **maximum values) and accuracy metrics for each property averaged across all depth**  
861 **intervals.**

Property	Unit	Min	Max	RMSE	$R^2$	$\rho$
Sand	% mass	0.0	100.0	9.322	0.841	0.918
Silt	% mass	0.0	100.0	6.556	0.788	0.889
Clay	% mass	0.0	100.0	5.891	0.841	0.918
pH	$\log_{10}([H^+])$	3.0	10.0	0.270	0.945	0.972
BD (oven-dry)	$g/cm^3$	0.5	2.0	0.164	0.704	0.843
SOM	% mass	0.0	100.0	1.961	0.608	0.801

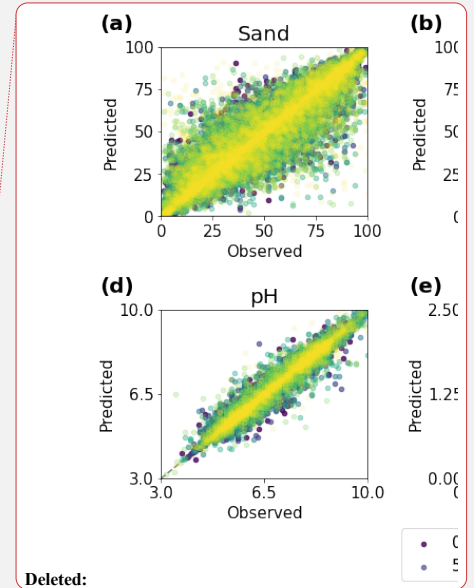
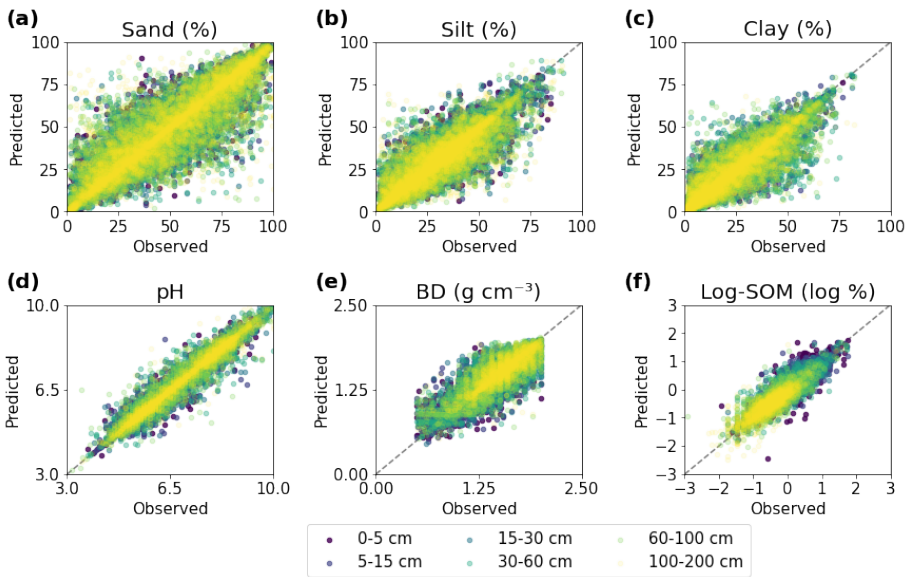
862

863 The posterior predictions of soil properties all align with the co-located observations and  
864 can capture the general trend of observations (Figure 5). Predictions of pH show the most  
865 concentrated clustering to the dashed line, indicating good agreement with observations  
866 across all depths. SOM and bulk density show relatively weaker performance compared to

867 other predicted soil properties. And this pattern of reduced accuracy persists throughout  
868 all depths.

869

870 As Figure 5 shows, the performance of the model tends to decline with increasing soil  
871 depth, except for SOM. This decline is primarily due to several reasons. First, the  
872 availability of soil data is often greater for shallower layers compared to deeper layers  
873 (such as > 1m), which limits the model's ability to learn patterns in deep layers. Second,  
874 remote sensing-derived soil covariates can only observe surface properties. Predictions  
875 for deeper layers rely on soil horizon information, soil profiles, geology, and parent  
876 material-related features. The certainty and quantity of them are less than easily  
877 measurable surface covariates. However, SOM shows better performance in deeper layers  
878 compared to surface layers. This is likely because surface SOM is highly variable due to  
879 factors like residue, land use, and management practices, while deeper SOM tends to be  
880 more stable.



881

882 **Figure 5: Evaluating posterior predictions with observations for six soil properties: (a)**  
 883 **sand, (b) silt, (c) clay, (d) pH, (e) bulk density (BD), and (f) log-scaled soil organic**  
 884 **matter (SOM). The left side shows scatter plots of posterior predictions versus**  
 885 **observations across six depth intervals, with each depth represented by a distinct**  
 886 **color. The dashed black line represents perfect prediction.**

887

888 **3.2 Comparison of Prior and Posterior Soil Predictions**

889 Prior and posterior predictions of soil properties are compared against co-located  
 890 observations to assess the added value of residual correction. The radar plots in Figure 6  
 891 illustrate the improvements achieved through the residual correction method using three  
 892 normalized metrics: 1-normalized absolute bias ( $1-|Bias|$ ), coefficient of determination

894 ( $R^2$ ), and 1-normalized RMSE by ranges of soil variability (1-nRMSE). These metrics are  
895 computed with values of soil properties, instead of on their residuals. Values in Figure 6  
896 closer to the outer edge of each plot indicate better model performance. Overall, all soil  
897 properties maintain reasonable normalized bias, and nRMSE (with nRMSE values  
898 consistently less than 0.02 for both prior and posterior predictions). However, the prior  
899 predictions tend to underestimate the variability of soil properties. As a result, the  
900 normalized metrics for prior and posterior predictions are similar, while the  $R^2$  values show  
901 some differences.

902

903 For all soil properties, posterior predictions consistently outperform prior predictions  
904 across all metrics. For particle size fractions,  $R^2$  values show the largest improvements:  
905 sand increases from 0.35 to 0.84, silt from 0.19 to 0.79, and clay from 0.25 to 0.84. The  
906 nRMSE metric also shows improvements. Sand decreases from 0.19 to 0.09, silt from 0.14  
907 to 0.07, and clay from 0.16 to 0.07, showing reductions in prediction errors using the  
908 residual correction.

909

910 Aggregating data from all depths. Figure 6 shows the degree of improvement across  
911 different soil properties. Prior pH predictions already demonstrate reasonable accuracy,  
912 with an  $R^2$  of 0.54 and nRMSE of 0.11. After the residual correction, these metrics improve  
913 to 0.94 for  $R^2$  and 0.04 for nRMSE. Bulk density and SOM show the biggest gains. For bulk  
914 density, the  $R^2$  increasing from 0.16 to 0.70 and nRMSE reducing from 0.18 to 0.11. Prior

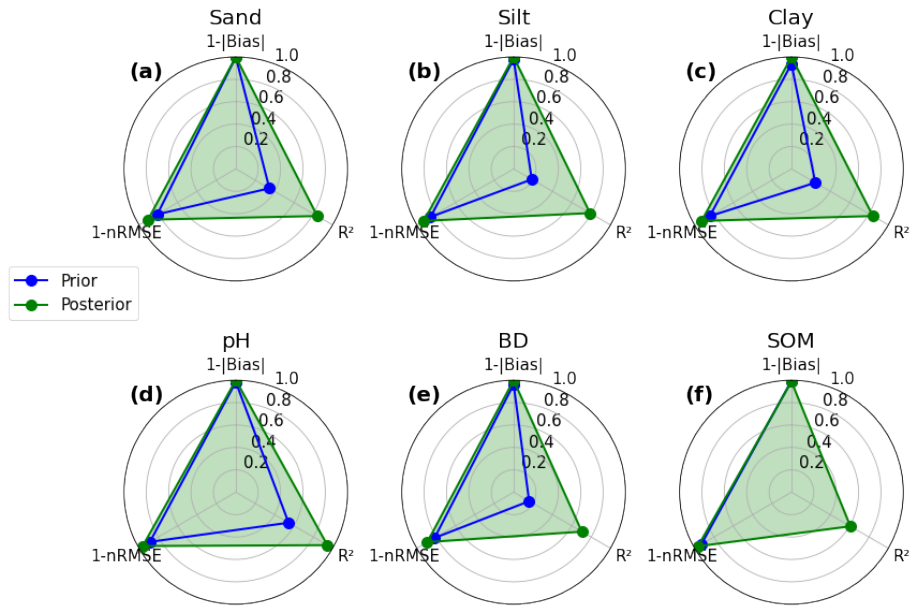
Deleted: ,

Deleted: .

Deleted: also

Deleted: different degrees

919 SOM are underfitted with a low  $R^2$  value. With the residual correction, the posterior SOM  
 920 show a positive  $R^2$  of 0.61. The nRMSE for SOM also improves from 0.07 to 0.04.



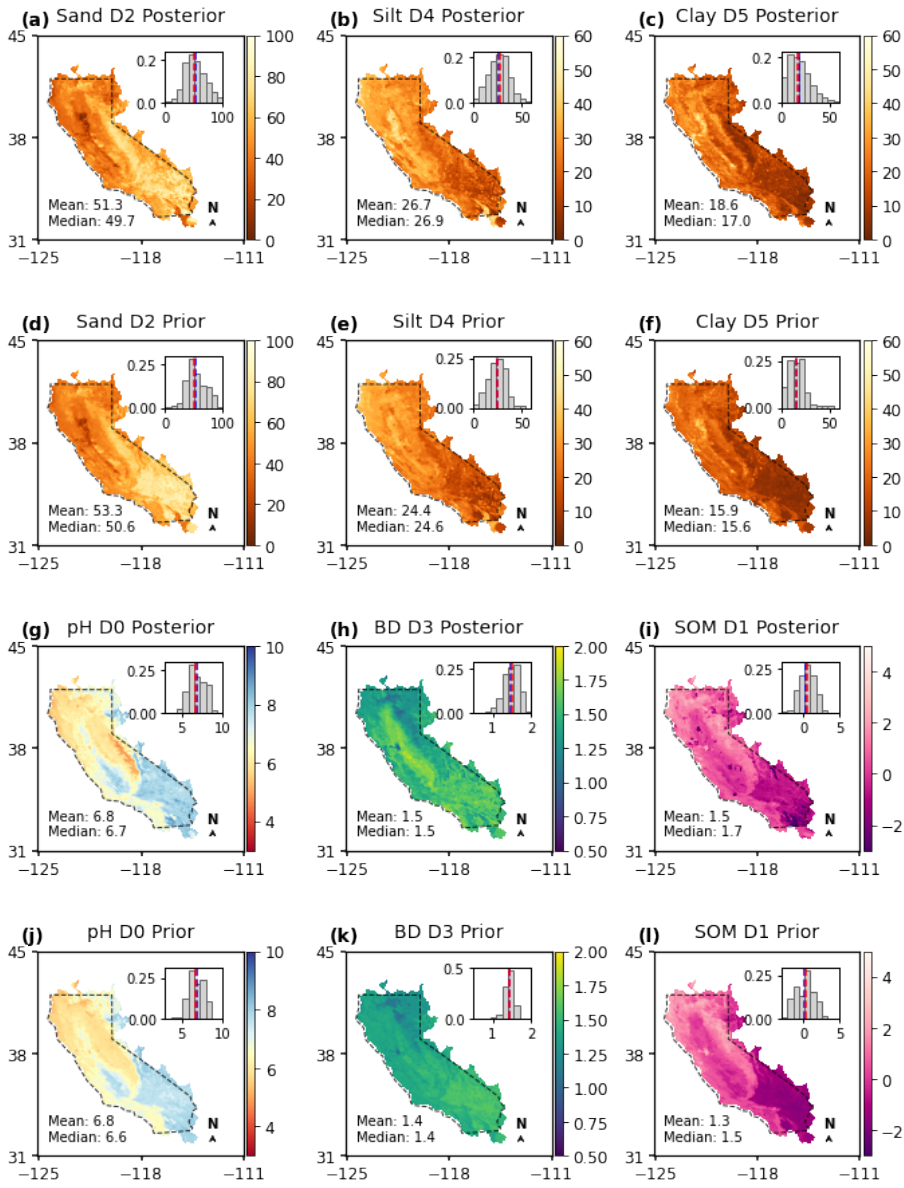
921  
 922 **Figure 6: Radar plots comparing the performance metrics of prior and posterior**  
 923 **predictions for six soil properties: (a) sand, (b) silt, (c) clay, (d) pH, (e) oven-dry bulk**  
 924 **density (BD), and (f) soil organic matter (SOM). Each plot presents three metrics: 1-**  
 925 **normalized absolute bias (1-|Bias|), coefficient of determination ( $R^2$ ), and 1-**  
 926 **normalized RMSE by ranges of soil variability (1-nRMSE). Prior predictions are shown**  
 927 **in blue, and posterior predictions in green. All metrics are scaled from 0 to 1, where**  
 928 **values closer to the outer edge of the plot indicate better model performance. The**  
 929 **green shaded area highlights the improvement achieved by the posterior predictions**  
 930 **over prior estimates.**

931

932 Horizontal spatial patterns of the six soil properties are presented in Figure 7. In the  
933 Central Valley California, soils are mostly medium textured with about 30% silt and lower  
934 sand content compared to surrounding areas. In the Mojave and Colorado Deserts, high  
935 sand contents (> 60%) with low clay contents are observed. SOM contents are also low in  
936 these areas. The histograms show how residual correction adjusts the distribution of soil  
937 properties.

938

939 For SOM and bulk density, the prior predictions often underestimate the observed  
940 variation. Figure 7 shows that the residual correction processes add noticeable spatial  
941 variations between prior and posterior soil maps. Prior bulk density values are often  
942 clustered around  $1.5 \text{ g/cm}^3$ , whereas the posterior histogram presents a broader range,  
943 spanning from  $1.25 \text{ g/cm}^3$  to  $1.6 \text{ g/cm}^3$ , capturing more heterogeneity of bulk density.  
944 Similarly, the residual correction adds soil heterogeneity to SOM. The posterior SOM can  
945 delineate water bodies, where SOM content is abruptly lower than the surrounding areas.  
946 Additionally, the posterior SOM maps present hill features in the desert areas.



947

948 **Figure 7: Spatial distribution of six soil properties (sand, silt, clay content, pH, bulk**

949 **density, and soil organic matter) across California. Maps of prior and posterior soil**  
950 **properties are compared. The corresponding frequency distributions of these soil**  
951 **properties are displayed in the right corner. Dashed polygons represent the**  
952 **continental part of California. In the histograms, the blue and red dashed lines**  
953 **represent the mean and median values, respectively. The maps labeled D0 to D5**  
954 **correspond to the first vertical layer down to the deepest layer. Note the map and**  
955 **distribution of soil organic matter (SOM) is log-scaled. Mean and median values are**  
956 **computed from the original SOM data.**

957

958 Soil profiles used for evaluating residual correction are grouped according to their  
959 corresponding pixel's land use classification from the National Land Cover Database  
960 (NLCD). Figure 8 presents selected vertical soil profiles of sand content, oven-dry bulk  
961 density, and SOM across three land use categories: forest, cultivated crops, and wetland.  
962 The number of samples varies by land use, with forests having the most, cultivated crops  
963 approximately half as many, and wetlands the fewest across California. To ensure a  
964 balanced visualization, a similar number of profiles are selected from each category. Sand  
965 content is chosen due to its broader range of variation (0-100%) compared to silt and clay  
966 (< 60% range). SOM and bulk density, which show relatively lower performance metrics,  
967 are included to assess the model's 'lower-bound performance'. These vertical profiles  
968 were not used during model training.

969

970 In Figure 8, solid lines represent the mean soil profiles for sand content, oven-dry bulk  
971 density, and SOM across forest, cultivated crops, and wetland land use categories. Blue  
972 lines, red lines, and green lines indicate prior, observation, and posterior predictions.  
973 Comparing the solid lines, the posterior predictions align more closely with the observed  
974 data compared to the prior estimates. However, the degree of alignment varies by soil  
975 property. For sand content and SOM, the posterior predictions show better agreement with  
976 observations, while bulk density predictions exhibit greater discrepancies, particularly in  
977 cultivated areas.

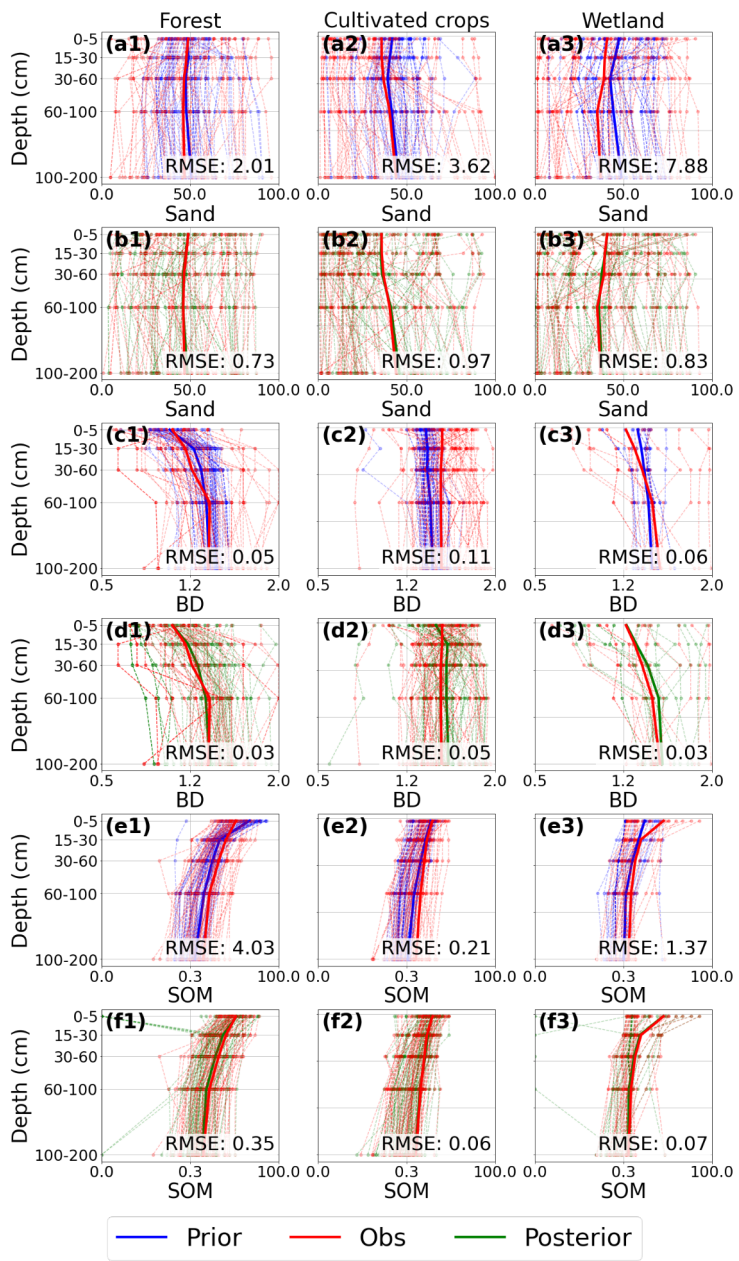
978

979 For sand content, the residual correction process improves estimates, especially in  
980 wetlands, with RMSE decreasing from 7.68 to 0.77 (%). Bulk density predictions perform  
981 better in forested and wetland areas. In cultivated crops, the posterior predictions show  
982 larger discrepancies. This suggests that bulk density is more challenging to predict in  
983 agricultural lands, particularly in shallow layers, likely due to agricultural activities. For  
984 SOM, the residual correction effectively improves estimates, especially in the surface  
985 layers of wetlands.

986

987 Dashed lines in Figure 8 represent individual soil profiles. Prior predictions often  
988 underestimated the variability in soil properties, struggling to capture extreme values. After  
989 the residual correction, the posterior predictions are better able to approximate these  
990 extremes. However, the correction process sometimes introduces additional noise. For  
991 example, some low SOM values (such as  $0.001 \text{ g/cm}^3$ ) were generated during residual

992 correction, even though such values are not presented in the observed data. It is likely due  
993 to that we used the van Bemmelen factor (1.724) to convert the prior soil organic matter to  
994 soil organic carbon.



996 **Figure 8: Vertical distribution of soil properties (sand content, oven-dry bulk density,**  
997 **and soil organic matter SOM) across three land use categories: forest, cultivated**  
998 **crops, and wetland. Prior estimates (blue), posterior estimates (green), and**  
999 **observations (red) are shown as depth profiles. Dashed lines represent individual**  
1000 **measurements, and solid lines show mean values. RMSE is computed elementwise to**  
1001 **evaluate model performance across all depths. X-axis and Y-axis represent value**  
1002 **ranges of a soil property and vertical depth intervals, respectively.**

1003

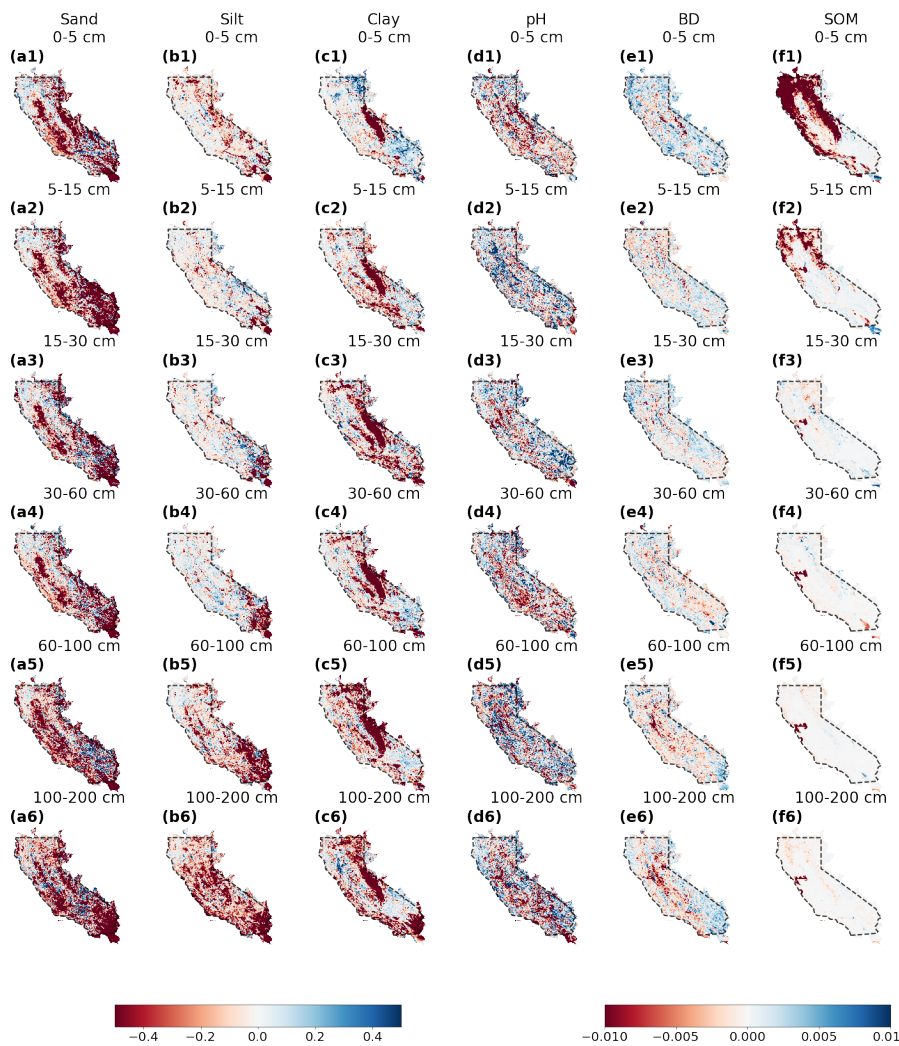
### 1004 **3.3 Uncertainty Analysis**

1005 Figure 9 shows the differences between 5% — 95% posterior and prior prediction interval  
1006 widths (PIWs) for six soil properties—sand, silt, clay, pH, bulk density, and SOM—from  
1007 surface to 2-m deep. The differences are calculated by subtracting the prior PIWs from the  
1008 posteriors. Red areas present a reduction in posterior PIW, indicating the residual  
1009 correction has reduced uncertainties of soil properties predictions. Blue pixels suggest the  
1010 opposite. White areas represent regions where the prior and posterior uncertainties are  
1011 similar.

1012

1013 In Figure 9, most pixels show reduced uncertainty for sand content after residual  
1014 correction, particularly in agricultural and desert regions. This improvement is attributed to  
1015 the inclusion of additional soil profile data from these areas. For clay content, the posterior  
1016 predictions consistently show reduced uncertainty across the Sierra Nevada Mountain  
1017 ranges. For SOM, the posterior PIWs improved in shallower layers (0-15 cm) over both the

1018 Coastal Ranges and the Sierra Nevada Mountains, with the coastal line showing notably  
1019 narrower PIWs. For pH, the results present a mixed pattern of PIWs after residual  
1020 correction, with some areas showing reduced uncertainty and others showing the  
1021 opposite. Similarly, bulk density exhibits a mixed pattern, though deeper layers (60 cm to 2  
1022 m) generally show reduced uncertainty in the Central Valley, California.



1023

1024 **Figure 9: Differences of 5% — 95% posterior and prior prediction interval widths (PIWs)**

1025 **for soil properties across different depths. Each column represents a specific soil**

1026 **property and rows show different depths. Black polygons represent the continental**

1027 **part of California. Differences between posterior and prior PIWs are in a red-to-blue**

1028 **color scale. Red pixels indicate a decrease in posterior PIW, indicating residual**  
1029 **correction reduces uncertainties. Vice versa for blue pixels. White areas indicate**  
1030 **similar extent of uncertainties. The left colorbar corresponds to sand, silt, clay with**  
1031 **wider ranges of PIW differences. The right colorbar represents other properties with**  
1032 **smaller PIW changes.**

1033

#### 1034 **4 Discussion**

##### 1035 **4.1 Limitations in Soil Profile Data**

1036 The effectiveness of residual correction depends on the spatial and vertical distribution of  
1037 soil profiles used to calculate residuals. In regions with sparse sampling, such as  
1038 California's desert areas (Figure 1), the limited number of profiles leads to interpolating the  
1039 entire area using limited observations. If soil heterogeneity is not captured by these limited  
1040 samples, the residual correction would overlook it. For soil texture, most data collected by  
1041 staff working on multiple projects under the National Institute of Food and Agriculture  
1042 (NIFA) and the Sustainable Agricultural Systems (SAS) programs range from the surface to  
1043 1.1 meters deep (additional field measurements used in this work). We use spline  
1044 interpolation to predict soil texture data beyond 1.1-m depths. It assumes vertical  
1045 continuity in soil properties, which may not reflect abrupt changes in subsurface layers.

1046

1047 Uncertainty also arises from converting some soil organic carbon (SOC) data to soil  
1048 organic matter (SOM). We used the van Bemmelen factor (1.724) to convert SOC to SOM  
1049 profiles. This factor does not hold true in scenarios such as organic-rich soils. Adding data

1050 quality controls—such as filtering profiles based on metadata (such as soil type, land  
1051 use)—could filter out samples that are not suitable for this conversion. However, this  
1052 conversion still has uncertainties, since even for mineral soils, this factor still has a certain  
1053 extent of variation depending on the organic matter composition (lower for soils with more  
1054 decomposed organic matter), soil types (forest soils or wetland soils with anaerobic  
1055 decomposition), and environmental influences (such as microbial activity).

1056

#### 1057 **4.2 Computational Challenges**

1058 The iterative residual correction process on distributions requires computational  
1059 resources, particularly when applied to large-extent or high-resolution datasets. This  
1060 process involves adjusting multiple values for each pixel, as each pixel represents a  
1061 distribution of soil properties. This process can be approached in two ways. The first  
1062 method involves correcting the residual values for each pixel, adding these residuals to  
1063 update the posterior values of soil properties, and then converting these updated values to  
1064 generate a posterior distribution of soil properties. The second method first converts all  
1065 pixel values into the same histogram bins and then corrects the shape of these histogram  
1066 bins for each pixel. Thus, the number of values retained per pixel affects computational  
1067 expense. Based on our experience, using method two, especially for soil texture, requires  
1068 100-bin histograms. Using method one with 20 most probable prior property values for  
1069 residual correction can achieve comparable results while reducing memory usage.

1070

1071 The iterative process of updating features and correcting residuals also plays a role. In our  
1072 simulations, we observed that subsequent residual corrections generally align with  
1073 previous ones. To ensure consistency, we require the corrections to converge more than  
1074 three times across different depths. For example, residual correction for a 1-km soil  
1075 property map over California takes approximately two hours after preprocessing the input  
1076 data. However, processing higher-resolution datasets, such as those at a 10-meter scale,  
1077 can demand significantly more computational resources. This highlights the trade-off  
1078 between resolution and computational efficiency in DSM projects.

1079

#### 1080 **4.3 Temporal and Spatial Constraints**

1081 The current method does not account for temporal changes in soil properties, limiting its  
1082 applicability to dynamic properties like soil organic matter or bulk density. Incorporating  
1083 temporal covariates (such as seasonal land surface temperature, recent land-use  
1084 changes) or stratifying soil profiles by collection date could address this. However, such  
1085 improvements rely on the availability of temporally resolved soil data, which are often  
1086 limited in quantities and sampling frequency.

1087

1088 Spatial clustering of soil samples poses another challenge. While duplicate profiles were  
1089 removed during data preprocessing, nearby samples may still share a certain level of  
1090 similarity due to spatial autocorrelation. This could lead to overly optimistic evaluation of  
1091 residual correction performance. Two methods can help address this issue:

1092 (1) Cross-validation with spatial considerations: Implement a cross-validation  
1093 method for splitting training and validation sets with attention to sample locations.  
1094 Ensure a minimum distance between training samples and evaluation data.

1095  
1096 (2) Independent dataset evaluation: Use independent datasets to evaluate the  
1097 model. CONUS-wide instrumental network, such as the U.S. Climate Reference  
1098 Network and the National Ecological Observatory Network, provide independent  
1099 soil data. However, these datasets have limitations as they were collected with  
1100 clustering to certain landscapes, potentially introducing bias in the evaluation.

1101

#### 1102 4.4 Similar Studies

1103 Several continental-scale DSM products (or methods) are compared, including the Soil  
1104 Survey Geographic Database (SSURGO), the Gridded National Soil Survey Geographic  
1105 Database (gNATSGO), the Probabilistic Layers for the Assessment of Soils (POLARIS), Soil-  
1106 Landscape Unified Synthesis (SOLUS), and the pruned Hierarchical Random Forest with  
1107 iterative bias correction (pHRF with IRC) soil properties. SSURGO is a traditional, polygon-  
1108 based product derived from expert field surveys and remains widely used in agricultural  
1109 applications (Soil Survey Staff et al., 2023). gNATSGO mainly builds on SSURGO by  
1110 rasterizing its map units to improve spatial coverage. And its estimation of soil properties  
1111 still rely on utilizing metadata of legacy soil data (Soil survey staff, 2023). These two still  
1112 inherit legacy data's limitations, such as scale inconsistency between soil map units and  
1113 derived soil maps, inconsistencies with field observations, and report distribution of soil

Deleted: (Soil Survey Staff et al., 2023)

Deleted: (Soil survey staff, 2023)

1116 properties with only three values (low end value, representative value, and high end value)  
1117 (Rossiter et al., 2022; Soil Survey Staff, 2025; Xu et al., 2025).

Field Code Changed

1118  
1119 Development of the following DSM products incorporates quantitative models in their  
1120 methodology. POLARIS produces probabilistic soil property maps using machine learning  
1121 and the DSMART algorithm (Chaney et al., 2016, 2019; Odgers et al., 2015), while the  
1122 uncertainties in the DSMART algorithm can propagate into POLARIS. SOLUS integrates  
1123 legacy soil data with georeferenced field observations and employs linear adjusted  
1124 Random Forest to predict soil properties (Nauman et al., 2024). SOLUS hierarchizes soil  
1125 data with different qualities into its training dataset, giving more attention to georeferenced  
1126 observations. However, since it also uses resampled soil data derived from polygon-based  
1127 soil map units, this process may introduce additional uncertainties into the final product.  
1128 The pHRF with IRC follows a different approach. Unlike most DSM methods that directly  
1129 predict soil properties from input data, this approach works in two steps: first, it generates  
1130 a prior estimate of soil taxa and property values, then iteratively adjusts these estimates to  
1131 improve model performance. In future work, the pHRF with IRC method will be applied on  
1132 large scale and assessed with more soil properties to evaluate its generalizability.

Deleted: ¶

... [7]

1133

## 1134 **5 Conclusion**

1135 The study introduces an iterative residual correction method for post processing used in a  
1136 Digital Soil Mapping (DSM) framework. The method integrates additional soil profile data  
1137 and iteratively optimizes the feature space to refine the distribution of soil properties until

1144 the residual correction model converges. Convergence is achieved when the median  
1145 difference between updated and previous predictions falls below a predefined threshold,  
1146 ensuring consistent predictions. The proposed DSM method operates through two primary  
1147 steps: (1) generating prior soil property maps using the pruned hierarchical Random Forest  
1148 (pHRF) approach, and (2) performing iterative residual correction on the priors. Residuals  
1149 (differences between observed values and prior predictions) are calculated and added to  
1150 the prior values of soil property to adjust the statistical shape of the probability distribution  
1151 pixel-by-pixel. The feature space, which includes soil covariates, depth information, and  
1152 vertical correlations, is iteratively optimized to capture incremental adjustments to  
1153 subsequent predictions.

1154  
1155 Using this method, we updated posterior distribution of soil properties for sand, silt, clay  
1156 content, soil pH, oven-dry bulk density, and soil organic matter over California. The results  
1157 show improvements in the accuracy of soil properties predictions, as shown by multiple  
1158 metrics including RMSE,  $R^2$ , and correlation coefficients. Furthermore, the iterative  
1159 residual correction model reduced prediction uncertainties, presenting narrower  
1160 prediction intervals compared to the priors.

1161  
1162 Several innovations contribute to the method's improvements. First, the integration of  
1163 additional soil profiles allows the model to further learn from georeferenced soil  
1164 information, complementing prior soil property estimates derived from traditional  
1165 surveys. Second, the iterative update of feature space captures both spatial and vertical

Deleted: optimization

1167 soil heterogeneity through a carefully selected combination of soil covariates and vertical  
1168 correlations among soil profile observations. Third, the convergence-based approach to  
1169 residual correction ensures stable output of posterior predictions while avoiding overfitting  
1170 since only converged residuals are added to the priors. Fourth, the implementation of  
1171 physical constraints and compositional data handling maintains the realism of predicted  
1172 soil properties. Future research could explore the application of this framework to other  
1173 soil properties and environmental contexts, such as soil hydraulic properties and CONUS-  
1174 wide simulation, to test the framework's generalization, supporting informed decision-  
1175 making in soil-related applications.

1176

#### 1177 **Data Availability**

1178 Data will be made available on request. [Code is available on](#)  
1179 [https://github.com/emmaxu43/IRC\\_CA/tree/main](https://github.com/emmaxu43/IRC_CA/tree/main).

1180

#### 1181 **Author Contributions**

1182 Chengcheng Xu and Nathaniel Chaney designed the study and developed the  
1183 methodology. Chengcheng Xu wrote the original draft and wrote the codes to produce the  
1184 methodology and analyses. Nathaniel Chaney supervised the work, provided resources  
1185 and funding, and helped guide the research direction. [Elia Scudiero provided funding,](#)  
1186 [project management, co-supervision.](#) Elia Scudiero and Ray Anderson provided soil  
1187 property samples from California that were used as part of the input dataset. Chengcheng

Deleted: Elia Scudiero

1189 Xu, Nathaniel Chaney, Elia Scudiero, and Ray Anderson discussed the results and  
1190 contributed to revising and editing the manuscript.

Deleted: Scudiero

1191

### 1192 **Competing Interests**

1193 The authors declare that they have no conflict of interest.

1194

### 1195 **Acknowledgements**

1196 This research was supported by the Agriculture and Food Research Initiative Competitive

Deleted: study

1197 Grant no. 2020-69012-31914 from the USDA National Institute of Food and Agriculture. The

Deleted: -NIFA-AFRI-006739 grant for sustainable agricultural systems

1198 authors want to thank Dr. Todd Skaggs for his and his teams' support for gathering input  
1199 data for this work. His and Dr. Ray Anderson's efforts are supported by USDA-ARS, Office  
1200 of National Programs (projects 2036-61000-019-000-D and 2036-61000-019-006-R). The  
1201 U.S. Department of Agriculture prohibits discrimination in all its programs and activities on  
1202 the basis of race, color, national origin, age, disability, and where applicable, sex, marital  
1203 status, familial status, parental status, religion, sexual orientation, genetic information,  
1204 political beliefs, reprisal, or because all or part of an individual's income is derived from  
1205 any public assistance program (not all prohibited bases apply to all programs). Persons  
1206 with disabilities who require alternative means for communication of program information  
1207 (braille, large print, audiotape, etc.) should contact USDA's TARGET Center at (202) 720-  
1208 2600 (voice and TDD). To file a complaint of discrimination, write to USDA, Director, Office  
1209 of Civil Rights, 1400 Independence Avenue, S.W., Washington, D.C. 20250-9410, or call

1214 (800) 795-3272 (voice) or (202) 720-6382 (TDD). USDA is an equal opportunity provider and  
1215 employer.

1216

#### 1217 **Financial Support**

1218 The study was supported by USDA-NIFA-AFRI-006739 grant for sustainable agricultural  
1219 systems.

1220

#### 1221 **References**

1222 [Arrouays, D., McKenzie, N., Hempel, J., Forges, A. R. de, and McBratney, A. B.:](#)  
1223 [GlobalSoilMap: Basis of the global spatial soil information system. CRC Press, 496 pp.,](#)  
1224 [2014.](#)

1225 Batjes, N. H., Calisto, L., and de Sousa, L. M.: Providing quality-assessed and standardised  
1226 soil data to support global mapping and modelling (WoSIS snapshot 2023), Earth System  
1227 Science Data, 16, 4735–4765, <https://doi.org/10.5194/essd-16-4735-2024>, 2024.

1228 Chaney, N. W., Wood, E. F., McBratney, A. B., Hempel, J. W., Nauman, T. W., Brungard, C.  
1229 W., and Odgers, N. P.: POLARIS: A 30-meter probabilistic soil series map of the contiguous  
1230 United States, Geoderma, <https://doi.org/10.1016/j.geoderma.2016.03.025>, 2016.

1231 Chaney, N. W., Minasny, B., Herman, J. D., Nauman, T. W., Brungard, C. W., Morgan, C. L.  
1232 S., McBratney, A. B., Wood, E. F., and Yimam, Y.: POLARIS Soil Properties: 30-m  
1233 Probabilistic Maps of Soil Properties Over the Contiguous United States, Water Resources  
1234 Research, <https://doi.org/10.1029/2018WR022797>, 2019.

1235 Chen, C., Liaw, A., and Breiman, L.: Using random forest to learn imbalanced data,  
1236 University of California, Berkeley, 110, 24, 2004.

1237 Chilès, J.-P. and Delfiner, P.: Geostatistics: modeling spatial uncertainty, in: Geostatistics:  
1238 modeling spatial uncertainty, John Wiley & Sons, Ltd, 147–237,  
1239 <https://doi.org/10.1002/9781118136188.ch3>, 2012.

1240 Corwin, D. L. and Scudiero, E.: Field-scale apparent soil electrical conductivity, Soil  
1241 Science Society of America Journal, 84, 1405–1441, <https://doi.org/10.1002/saj2.20153>,  
1242 2020.

- 1243 Grunwald, S., Thompson, J. A., and Boettinger, J. L.: Digital Soil Mapping and Modeling at  
1244 Continental Scales: Finding Solutions for Global Issues, *Soil Science Society of America*  
1245 *Journal*, 75, 1201–1213, <https://doi.org/10.2136/SSSAJ2011.0025>, 2011.
- 1246 Haghverdi, A., Najarchi, M., öztürk, H. S., and Durner, W.: Studying unimodal, bimodal, PDI  
1247 and bimodal-PDI variants of multiple soil water retention models: I. Direct model fit using  
1248 the extended evaporation and dewpoint methods, *Water (Switzerland)*, 12,  
1249 <https://doi.org/10.3390/w12030900>, 2020.
- 1250 [Hartemink, A. E., Hempel, J., Lagacherie, P., McBratney, A., McKenzie, N., MacMillan, R. A.,](#)  
1251 [Minasny, B., Montanarella, L., de Mendonça Santos, M. L., Sanchez, P., Walsh, M., and](#)  
1252 [Zhang, G.-L.: GlobalSoilMap.net – A New Digital Soil Map of the World, in: Digital Soil](#)  
1253 [Mapping: Bridging Research, Environmental Application, and Operation, edited by:](#)  
1254 [Boettinger, J. L., Howell, D. W., Moore, A. C., Hartemink, A. E., and Kienast-Brown, S.,](#)  
1255 [Springer Netherlands, Dordrecht, 423–428, https://doi.org/10.1007/978-90-481-8863-](#)  
1256 [5\\_33, 2010.](#)
- 1257 Hengl, T., Heuvelink, G. B., and Stein, A.: A generic framework for spatial prediction of soil  
1258 variables based on regression-kriging, *Geoderma*, 120, 75–93, 2004.
- 1259 Hengl, T., De Jesus, J. M., Heuvelink, G. B. M., Gonzalez, M. R., Kilibarda, M., Blagotić, A.,  
1260 Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas,  
1261 R., MacMillan, R. A., Batjes, N. H., Leenaars, J. G. B., Ribeiro, E., Wheeler, I., Mantel, S.,  
1262 and Kempen, B.: SoilGrids250m: Global gridded soil information based on machine  
1263 learning, *PLoS ONE*, <https://doi.org/10.1371/journal.pone.0169748>, 2017.
- 1264 Jiang, Q., Fu, Q., and Wang, Z.: Delineating site-specific irrigation management zones,  
1265 *Irrigation and Drainage*, 60, 464–472, <https://doi.org/10.1002/ird.588>, 2011.
- 1266 Lesch, S., Rhoades, J., and Corwin, D.: ESAP-95 version 2.01 R: User manual and tutorial  
1267 guide, *Research Rpt*, 146, 17, 2000.
- 1268 Lesch, S. M.: Sensor-directed response surface sampling designs for characterizing spatial  
1269 variation in soil properties, *Computers and Electronics in Agriculture*, 46, 153–179,  
1270 <https://doi.org/10.1016/j.compag.2004.11.004>, 2005.
- 1271 Li, N., Zhao, X., Wang, J., Sefton, M., and Triantafyllis, J.: Digital soil mapping based site-  
1272 specific nutrient management in a sugarcane field in Burdekin, *Geoderma*, 340, 38–48,  
1273 <https://doi.org/10.1016/j.geoderma.2018.12.033>, 2019.
- 1274 [McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping,](#)  
1275 [Geoderma](#), 117, 3–52, [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.
- 1276 Minasny, B. and McBratney, A. B.: A conditioned Latin hypercube method for sampling in  
1277 the presence of ancillary information, *Computers & geosciences*, 32, 1378–1388, 2006.

- 1278 Mueller, T. G., Pierce, F. J., Schabenberger, O., and Warncke, D. D.: Map Quality for Site-  
1279 Specific Fertility Management, *Soil Science Society of America Journal*, 65, 1547–1558,  
1280 <https://doi.org/10.2136/sssaj2001.6551547x>, 2001.
- 1281 National Cooperative Soil Survey: NCSS Soil Characterization Database (Lab Data Mart),  
1282 2018.
- 1283 Nauman, T. W., Kienast-Brown, S., Roecker, S. M., Brungard, C., White, D., Philippe, J., and  
1284 Thompson, J. A.: Soil landscapes of the United States (SOLUS): Developing predictive soil  
1285 property maps of the conterminous United States using hybrid training sets, *Soil Science  
1286 Society of America Journal*, 88, 2046–2065, <https://doi.org/10.1002/saj2.20769>, 2024.
- 1287 Nussbaum, M., Zimmermann, S., Walthert, L., and Baltensweiler, A.: Benefits of  
1288 hierarchical predictions for digital soil mapping—An approach to map bimodal soil pH,  
1289 *Geoderma*, 437, 116579, <https://doi.org/10.1016/j.geoderma.2023.116579>, 2023.
- 1290 Odgers, N. P., McBratney, A. B., and Minasny, B.: Digital soil property mapping and  
1291 uncertainty estimation using soil class probability rasters, *Geoderma*, 237,  
1292 <https://doi.org/10.1016/j.geoderma.2014.09.009>, 2015.
- 1293 Oliver, M. A. and Webster, R.: A tutorial guide to geostatistics: Computing and modelling  
1294 variograms and kriging, *CATENA*, 113, 56–69,  
1295 <https://doi.org/10.1016/j.catena.2013.09.006>, 2014.
- 1296 Ortuani, B., Chiaradia, E. A., Priori, S., L'Abate, G., Canone, D., Comunian, A., Giudici, M.,  
1297 Mele, M., and Facchi, A.: Mapping Soil Water Capacity Through EMI Survey to Delineate  
1298 Site-Specific Management Units Within an Irrigated Field, *Soil Science*, 181, 252,  
1299 <https://doi.org/10.1097/SS.000000000000159>, 2016.
- 1300 Poggio, L., De Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E., and  
1301 Rossiter, D.: SoilGrids 2.0: Producing soil information for the globe with quantified spatial  
1302 uncertainty, *SOIL*, 7, 217–240, <https://doi.org/10.5194/SOIL-7-217-2021>, 2021.
- 1303 Powers, J. S., Corre, M. D., Twine, T. E., and Veldkamp, E.: Geographic bias of field  
1304 observations of soil carbon stocks with tropical land-use changes precludes spatial  
1305 extrapolation, *Proceedings of the National Academy of Sciences*, 108, 6318–6322,  
1306 <https://doi.org/10.1073/pnas.1016774108>, 2011.
- 1307 Ramcharan, A., Hengl, T., Nauman, T., Brungard, C., Waltman, S., Wills, S., and Thompson,  
1308 J.: Soil Property and Class Maps of the Conterminous United States at 100-Meter Spatial  
1309 Resolution, *Soil Science Society of America Journal*, 82, 186–201,  
1310 <https://doi.org/10.2136/sssaj2017.04.0122>, 2018.
- 1311 Rossiter, D. G., Poggio, L., Beaudette, D., and Libohova, Z.: How well does digital soil  
1312 mapping represent soil geography? An investigation from the USA, *SOIL*, 8, 559–586,  
1313 <https://doi.org/10.5194/soil-8-559-2022>, 2022.

1314 Schmidinger, J. and Heuvelink, G. B. M.: Validation of uncertainty predictions in digital soil  
1315 mapping, *Geoderma*, 437, 116585, <https://doi.org/10.1016/j.geoderma.2023.116585>,  
1316 2023.

1317 [Scudiero, E., Corwin, D. L., Markley, P. T., Pourreza, A., Rounsaville, T., Bughici, T., and](#)  
1318 [Skaggs, T. H.: A system for concurrent on-the-go soil apparent electrical conductivity and](#)  
1319 [gamma-ray sensing in micro-irrigated orchards, \*Soil and Tillage Research\*, 235, 105899,](#)  
1320 [2024.](#)

1321 Sharififar, A., Sarmadian, F., Malone, B. P., and Minasny, B.: Addressing the issue of digital  
1322 mapping of soil classes with imbalanced class observations, *Geoderma*, 350, 84–92,  
1323 <https://doi.org/10.1016/j.geoderma.2019.05.016>, 2019.

1324 Shi, G., Sun, W., Shanguan, W., Wei, Z., Yuan, H., Zhang, Y., Liang, H., Li, L., Sun, X., Li, D.,  
1325 Huang, F., Li, Q., and Dai, Y.: A China dataset of soil properties for land surface modeling  
1326 (version 2), <https://doi.org/10.5194/essd-2024-299>, 29 August 2024.

1327 [Soil, K.: Survey laboratory methods manual, \*Soil Survey Investigations Report\*, 1996.](#)

1328 [Soil Survey Staff: Kellogg Soil Survey Laboratory methods manual, U.S. Department of](#)  
1329 [Agriculture, Natural Resources Conservation Service, Lincoln, Nebraska, 2014.](#)

1330 Soil survey staff: Gridded National Soil Survey Geographic (gNATSGO) Database for the  
1331 Conterminous United States, 2023. Natural Resources Conservation Service, United  
1332 States Department of Agriculture.

1333 Soil Survey Staff: Gridded Soil Survey Geographic (gSSURGO) Database for the  
1334 Conterminous United States, 2025. Natural Resources Conservation Service, United  
1335 States Department of Agriculture.

1336 Soil Survey Staff, Natural Resources Conservation Service, and United States Department  
1337 of Agriculture: Soil Survey Geographic (SSURGO) Database for the CONUS, 2023. Natural  
1338 Resources Conservation Service, United States Department of Agriculture.

1339 Sylvain, J.-D., Anctil, F., and Thiffault, É.: Using bias correction and ensemble modelling for  
1340 predictive mapping and related uncertainty: A case study in digital soil mapping,  
1341 *Geoderma*, 403, 115153, <https://doi.org/10.1016/j.geoderma.2021.115153>, 2021.

1342 Takoutsing, B., Heuvelink, G. B. M., Stoorvogel, J. J., Shepherd, K. D., and Aynekulu, E.:  
1343 Accounting for analytical and proximal soil sensing errors in digital soil mapping, *European*  
1344 *Journal of Soil Science*, 73, e13226, <https://doi.org/10.1111/ejss.13226>, 2022.

1345 Vereecken, H., Schnepf, A., Hopmans, J. W., Javaux, M., Or, D., Roose, T., Vanderborght, J.,  
1346 Young, M. H., Amelung, W., Aitkenhead, M., Allison, S. D., Assouline, S., Baveye, P., Bertl,  
1347 M., Brüggemann, N., Finke, P., Flury, M., Gaiser, T., Govers, G., Ghezzehei, T., Hallett, P.,  
1348 Hendricks Franssen, H. J., Heppell, J., Horn, R., Huisman, J. A., Jacques, D., Jonard, F.,

1349 Kollet, S., Lafolie, F., Lamorski, K., Leitner, D., McBratney, A., Minasny, B., Montzka, C.,  
1350 Nowak, W., Pachepsky, Y., Padarian, J., Romano, N., Roth, K., Rothfuss, Y., Rowe, E. C.,  
1351 Schwen, A., Šimůnek, J., Tiktak, A., Van Dam, J., van der Zee, S. E. A. T. M., Vogel, H. J.,  
1352 Vrugt, J. A., Wöhling, T., and Young, I. M.: Modeling Soil Processes: Review, Key  
1353 Challenges, and New Perspectives, *Vadose Zone Journal*, 15, vzi2015.09.0131,  
1354 <https://doi.org/10.2136/vzj2015.09.0131>, 2016.

1355 Vereecken, H., Amelung, W., Bauke, S. L., Bogaen, H., Brüggemann, N., Montzka, C.,  
1356 Vanderborght, J., Bechtold, M., Blöschl, G., Carminati, A., Javaux, M., Konings, A. G.,  
1357 Kusche, J., Neuweiler, I., Or, D., Steele-Dunne, S., Verhoef, A., Young, M., and Zhang, Y.:  
1358 Soil hydrology in the Earth system, *Nat Rev Earth Environ*, 3, 573–587,  
1359 <https://doi.org/10.1038/s43017-022-00324-6>, 2022.

1360 Wu, Y., Huang, Y., Chen, Z., Yao, Z., Fu, Y., Liu, K., Luo, X., and Wang, D.: Iterative Feature  
1361 Space Optimization through Incremental Adaptive Evaluation,  
1362 <https://doi.org/10.48550/arXiv.2501.14889>, 24 January 2025.

1363 Xu, C., Huang, J., Hartemink, A. E., and Chaney, N. W.: Pruned hierarchical Random Forest  
1364 framework for digital soil mapping: Evaluation using NEON soil properties, *Geoderma*, 459,  
1365 117392, <https://doi.org/10.1016/j.geoderma.2025.117392>, 2025.

1366 Zhang, G. and Lu, Y.: Bias-corrected random forests in regression, *Journal of Applied*  
1367 *Statistics*, 39, 151–160, <https://doi.org/10.1080/02664763.2011.578621>, 2012.

1368

▼ .....  
▲ .....  
**Page 2: [1] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**

▼ .....  
▲ .....  
**Page 2: [2] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**

▼ .....  
▲ .....  
**Page 22: [3] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**

▼ .....  
▲ .....  
**Page 23: [4] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**

▼ .....  
▲ .....  
**Page 24: [5] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**

▼ .....  
▲ .....  
**Page 24: [6] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**

▼ .....  
▲ .....  
**Page 44: [7] Deleted** **Chengcheng Xu** **3/1/26 4:00:00 PM**