

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

2 Chengcheng Xu¹, Elia Scudiero²³, Ray Anderson³², Nathaniel Chaney¹

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

5 ² Department of Environmental Sciences, University of California Riverside, Riverside, CA
6 92521, USA

7 ³ 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 emerges. 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,

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

43

44 **1 Introduction**

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

59
60 Digital soil maps are typically derived from existing soil surveys, geostatistical models,
61 machine learning, or hybrid approaches. Soil survey-based soil mapping method, which
62 use low, high, and representative values to describe soil property distributions for each soil
63 component (Soil Survey Staff et al., 2023). The method typically approximates each soil
64 component as a triangular distribution (Chaney et al., 2016; Soil survey staff, 2023),
65 potentially oversimplifying multi-modal distributions of soil properties in some cases

66 (Haghverdi et al., 2020; Nussbaum et al., 2023). Additionally, estimating soil properties
67 from synthetic sampling within a map unit could create artificial spatial patterns, adding
68 noises into the mapping results (Chaney et al., 2019). Developments such as Latin-
69 hypercube sampling and landscape adaptive covariance functions have improved the
70 representation of spatial patterns of soil properties (Minasny and McBratney, 2006). Yet,
71 soil survey-based approaches remain valuable particularly in areas where soil profile data
72 is limited (Nauman et al., 2024). Geostatistical models often require presumed
73 parameterization and are constrained by stationarity assumptions, which is difficult to
74 apply in areas with insufficient field knowledge (Oliver and Webster, 2014). To address
75 these challenges, non-parametric models, such as Random Forest, trained with hybridized
76 soil data that combine soil surveys with georeferenced soil profiles show potentials in
77 improving soil mapping, particularly for large-scale maps (Chaney et al., 2019; Nauman et
78 al., 2024).

79

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

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

98
99 Quantifying uncertainties in DSM is important for its practical applications (Schmidinger
100 and Heuvelink, 2023). DSM products represent soil properties as multi-dimensional
101 matrices showing vertical and horizontal soil variation (Vereecken et al., 2022), with each
102 pixel containing weighted possible values and their prediction uncertainties. These
103 uncertainties can be represented either as continuous values through prediction intervals
104 or as discrete classifications with associated class probabilities (Chaney et al., 2016,
105 2019; Hengl et al., 2017; Ramcharan et al., 2018). Common quantification approaches
106 include geostatistical techniques like kriging, where the nugget term accounts for
107 measurement errors while kriging variance reflects spatial uncertainty patterns (Chilès and
108 Delfiner, 2012; Takoutsing et al., 2022), and machine learning methods such as Quantile
109 Random Forest (QRF) which generates probability distributions from decision tree outputs

110 using values of soil properties (Poggio et al., 2021; Shi et al., 2024). For discrete
111 classifications, uncertainty derives from soil raster probabilities during soil taxa
112 classification (Chaney et al., 2016; Odgers et al., 2015). Given the data-driven nature of
113 DSM and frequent limitations in soil profile availability, integrating multiple qualified data
114 sources improves the amount of soil data and reduce prediction uncertainties (Nauman et
115 al., 2024), particularly in regions where predictions must rely more heavily on legacy soil
116 data.

117
118 In this study, we present a hybrid DSM approach combining pruned Hierarchical Random
119 Forest (pHRF) with iterative residual correction (IRC) method (Xu et al., 2025). The pHRF
120 method leverages the National Cooperative Soil Survey (NCSS) soil survey data and
121 georeferenced soil taxa information to generate prior distributions, while additional soil
122 profiles correct biases in prior predictions. This method builds on development in previous
123 research while addressing specific limitations. Sylvain et al. (2021) applied XGBoost
124 (sequential decision trees) and ensemble models to correct deterministic soil property
125 maps, demonstrating reduced bias for many soil properties (Sylvain et al., 2021). Zhang et
126 al. (2010) introduced a bias-correction technique with Random Forest models to mitigate
127 their tendency to regress toward mean values, though not in DSM contexts (Zhang and Lu,
128 2012). Our approach extends these concepts by probabilistically updating posterior
129 distributions at each location through an iterative correction process that continues until
130 convergence across vertical intervals. Vertical correlations are maintained through layer-
131 by-layer residual correction, which preserves inter-layer correlations while dynamically

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

139

140 **2 Methods**

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

147

148 The following sections first describe the general residual correction framework (Section
149 2.1). To illustrate the method concretely, we then provide a worked example using one
150 randomly selected soil column to demonstrate how the feature space is constructed and
151 updated across two consecutive iterations (Section 2.1.1). Building on this example, we
152 detail the key components of the IRC method: the iterative update of feature space
153 (Section 2.1.2), the convergence criteria for residual correction (Section 2.1.3), and the

154 process for updating posterior soil properties with physical constraints (Section 2.1.4).
155 Finally, we present the California case study (Section 2.2), describing the soil datasets
156 used (Section 2.2.1) and the implementation details for applying the IRC method over
157 California (Section 2.2.2).

158

159 **2.1 Iterative Residual Correction Framework for DSM**

160 Residual correction is implemented to address underestimated soil property variation in
161 prior maps (tendency to underestimate high values and overestimate low values,
162 smoothing out soil variation across landscape). The overall workflow of the IRC method
163 consists of three components: (1) prior map generation (Figure 1a), (2) residual preparation
164 (Figure 1b), and (3) iterative correction (Figure 1c).

165

166 First, probabilistic prior soil property maps are generated or retrieve probabilistic soil
167 property maps from an existing DSM product as the prior soil maps (Figure 1a). These
168 maps represent the initial estimates of soil properties and their associated uncertainties.

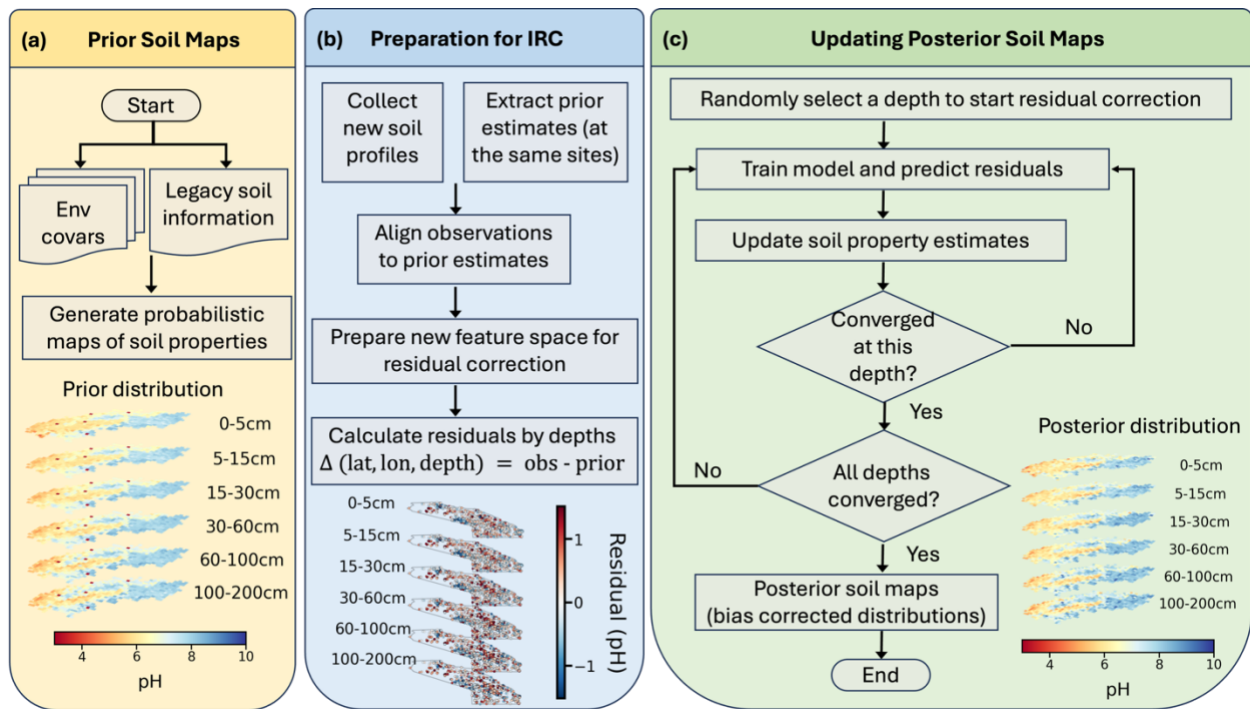
169 Second, a residual preparation step is carried out to enable correction using new soil
170 profile observations (Figure 1b). The preparation involves four key steps: (1) adding
171 additional soil profiles from new field measurements or databases; (2) spatially aligning
172 these profiles with the corresponding pixels in the prior soil maps using geographic
173 coordinates; (3) vertically aligning observations with prior predictions at matching depth
174 intervals; and (4) calculating residuals depth by depth as the difference between observed
175 values and prior predictions. During this stage, the feature space for residual modeling is

176 also prepared, consisting of static environmental covariates (which remain fixed
177 throughout iterations) and dynamic soil covariates (which are updated iteratively). Detailed
178 construction of the feature space is described in Section 2.1.1.

179

180 Finally, iterative residual correction is performed to update soil property estimates across
181 depths (Figure 1c). During each iteration, the model predicts residuals for one depth layer
182 at a time, with the layer selected randomly. A Random Forest regressor is trained to learn
183 the relationship between residuals and the feature space at sampled locations, then
184 interpolates residual corrections across the study area. Predicted residuals are added to
185 the prior (or previous iteration's) estimates to generate updated soil property values. After
186 each update, convergence is evaluated for the modeling depth by comparing the median
187 difference between the current residuals and those from the previous iteration. Once this
188 change falls below a predefined threshold, that depth is considered converged and
189 excluded from subsequent updates. The algorithm then focuses on the remaining
190 “unconverged” depths, until convergence is achieved across all layers. After convergence
191 is verified for all depths, the final corrected residuals are added to the prior estimates to
192 update the posterior distributions of soil properties.

193



194

195 Figure 1: Workflow for updating posterior soil property maps. The process begins with

196 panel (a), the preparation of environmental covariates (env covars) to generate

197 probabilistic maps of soil properties (prior soil maps). As illustrated in panel (b), the

198 preparation for residual correction involves adding additional soil profiles, spatially and

199 vertically aligning prior soil map values with new profile observations, calculating residuals

200 depth by depth, and preparing environmental covariates and soil covariates (new feature

201 space) for residual correction. Finally, as shown in panel (c), the iterative residual

202 correction step applies bias corrections across different depths, focusing on layers where

203 residuals have not yet stabilized. During each iteration, the model predicts residuals for

204 one depth at a time, randomly selecting a layer. Once residuals for a given depth converge,

205 that layer is excluded from further updates, allowing the model to concentrate on

206 remaining depths until all achieve stability. After verifying convergence across all depths,

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

209

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

220

221 **2.1.1 Worked Example**

222 The iterative residual correction method is further illustrated in Figure 2 using an example
223 with a randomly selected soil column. Figure 2a shows the location of the selected soil
224 column, where additional soil profile observations are available. The right panel displays
225 the top-3 probable pH values (from prior soil maps) at each depth intervals (0–5 cm, 5–15
226 cm, 15–30 cm, 30–60 cm, 60–100 cm, 100–200 cm), while the left panel shows the three
227 weights (probabilities) associated with these pH values. In this simplified example, we use
228 3 bins to represent the soil property distribution; however, in actual implementation, more

229 bins are maintained (typically top-12 probable values) to better capture soil variability. For
230 this demonstration, Depth 2 (D_2 ; 5–15 cm) is randomly selected as the modeling layer to
231 initiate the iterative correction process. Only one layer is modeled and updated for a given
232 iteration. Note that in real model execution, model generally processes over 3,000 soil
233 columns simultaneously in California, though only one column is shown here for clarity.

234

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

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

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

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

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

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

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

260

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

264

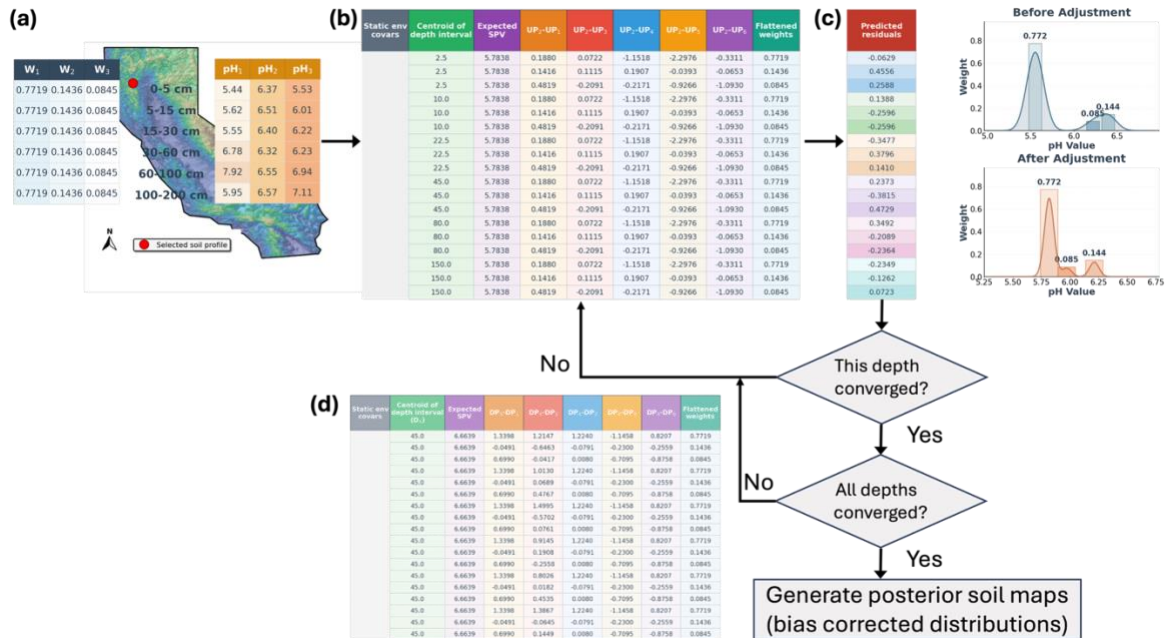
265 A Random Forest regressor is then trained using the feature space to predict residuals for
266 the modeling layer (D_2 in this example). The right panel in Figure 2c compares the
267 distribution of pH values before and after residual adjustment in the current iteration. After
268 applying the residual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

272 posterior soil property maps are generated by adding the last predicted residuals to the
273 prior values.

274

275 If either convergence check returns "No" (i.e., D_2 has not converged or other layers remain
276 unconverged), the algorithm continues iterating. Here, the soil property values for D_2 are
277 updated by adding the predicted residuals to the previous pH values. These updated
278 values are then used to reconstruct the feature space following the same structure
279 described above, updating the representative values, top-probable values, and inter-layer
280 differences. By updating soil covariates layer by layer and iteratively refining the feature
281 space, the next prediction retains prior knowledge while integrating new information about
282 soil heterogeneity and vertical relationships for soil profiles (Wu et al., 2025). A new
283 iteration begins by randomly selecting another unconverged layer, and the process repeats
284 until convergence is achieved across all depth layers.

285



286

287 Figure 2: Schematic illustration of the iterative residual correction (IRC) method using a

288 worked example at a randomly selected soil column. (a) Prior distributions and

289 observation location: The map shows the location of the selected soil column within the

290 study area. The right panel displays the top-3 probable pH values at each of the six depth

291 intervals (0–5 cm, 5–15 cm, 15–30 cm, 30–60 cm, 60–100 cm, 100–200 cm), while the left

292 panel shows the three weights (w_1 , w_2 , w_3) associated with these pH values. Depth 2 (D_2 ; 5–

293 15 cm) is randomly selected for this iteration. (b) Feature space components: The table

294 details the structure of the feature space used to train the Random Forest regressor for
295 residual prediction. The feature space comprises both static and dynamic components.

296 Static components include environmental covariates (satellite imagery, terrain attributes)

297 that remain unchanged throughout iterations, and weights (w_1 , w_2 , w_3) associated with top-

298 probable values. Dynamic soil covariates that are updated in each iteration include: the

299 centroid of the depth interval (e.g., 10 cm for D_2), the expected (representative) soil

300 property value computed as the weighted mean, the top-probable soil property values
301 reflecting intra-pixel heterogeneity, and inter-layer differences capturing vertical
302 correlations (e.g., D_2-D_1 , D_2-D_3). (c) Residual correction and convergence workflow: A
303 Random Forest model trained on the feature space predicts residuals for the modeling
304 layer D_2 . The right panel compares the pH distribution before and after residual
305 adjustment. The flowchart below describes the convergence logic: after predicting and
306 applying residuals to D_2 , the algorithm evaluates whether D_2 has converged. If D_2 has
307 converged, the algorithm checks whether all depth layers have achieved convergence. If
308 both checks pass, the final posterior soil property maps are generated by adding the last
309 converged residuals to the prior values. (d) If either check fails, the algorithm updates the
310 soil property values for D_2 by adding predicted residuals, reconstructs the feature space
311 with the updated values, randomly selects another unconverged layer, and repeats the
312 process. This iterative cycle continues until convergence is achieved across all six depth
313 layers.

314

315 **2.1.2 Convergence of Residual Correction**

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

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

324

325 **2.1.3 Update with Constraints**

326 During residual correction, a common issue arises where the addition of residuals to prior
327 soil property values results in values that exceed physical bounds (such as sand content >
328 100%). To address this, a residual update process with constraints is implemented.

329

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

337

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

344 silt, clay), an additional compositional constraint ensures that the three fractions sum to
345 100% at each pixel after residual correction.

346

347 **2.2 California Case Study: Soil Data and Model Implementation**

348 **2.2.1 Soil Data**

349 To demonstrate the IRC method, we apply it to soil property mapping in California. We use
350 georeferenced soil profiles with laboratory measurements of soil properties. We compiled
351 soil profile data from three primary sources: the World Soil Information Service (WoSIS),
352 the National Soil Characterization Database (SCD), and field measurements conducted in
353 California (Batjes et al., 2024; National Cooperative Soil Survey, 2018; Scudiero et al.,
354 2024).

355

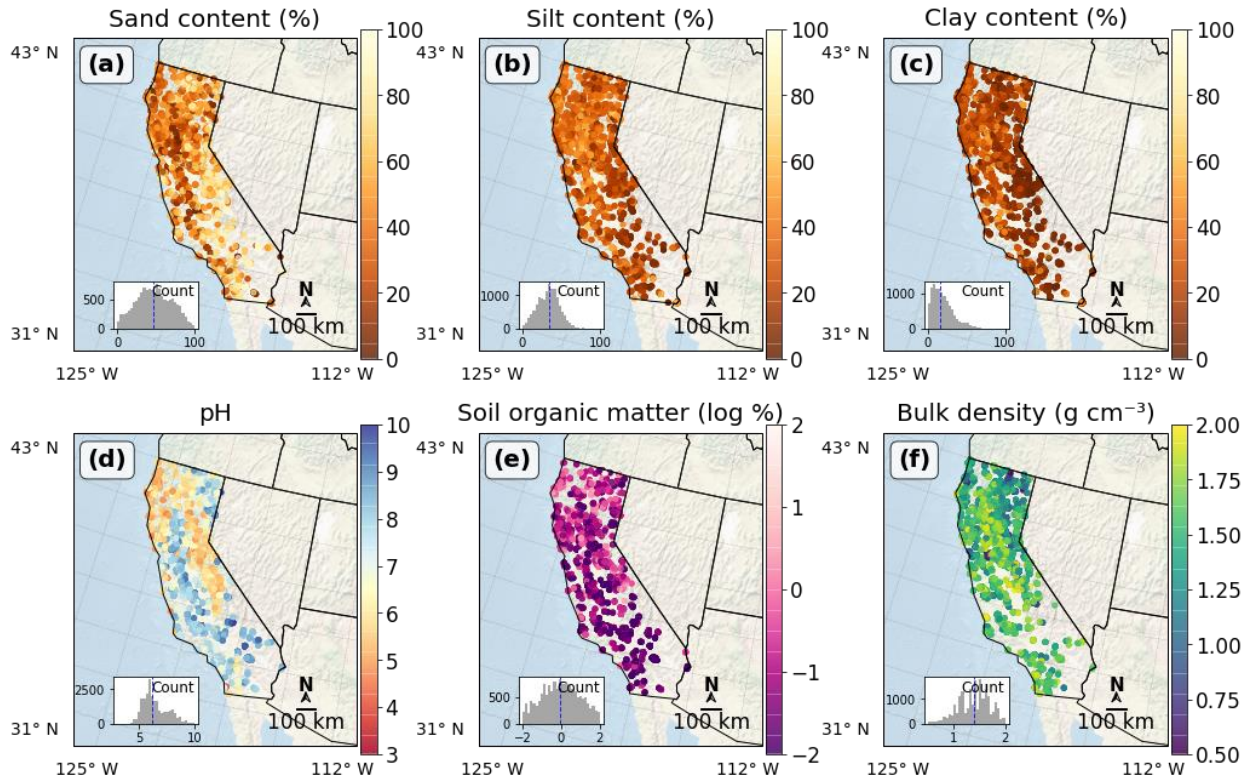
356 To ensure consistency across different data sources, we applied several quality control
357 steps. First, we checked the physical plausibility of all soil property values by defining a
358 valid range with specific minimum and maximum thresholds for each property. Any data
359 point falling outside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

366 to maintain precision for particle size analysis (Corwin and Scudiero, 2020; Scudiero et al.,
367 2024).

368

369 During preprocessing, we harmonized all soil data, which was originally reported at
370 different soil horizons, into six standard depth intervals: 0–5 cm, 5–15 cm, 15–30 cm, 30–
371 60 cm, 60–100 cm, and 100–200 cm (Arrouays et al., 2014). The harmonization was
372 performed using equal-area spline functions to interpolate soil property values from the
373 original horizon depths to these standard intervals (Hartemink et al., 2010, p.201). The
374 spline function fits a smooth curve through observed values at their measured depths,
375 then calculates the area under this curve within each standardized depth interval and
376 divides by the interval width to obtain the value. Location of soil profiles and their
377 distribution of soil property values are presented in Figure 3. Six soil properties are studied:
378 sand content, silt content, clay content, pH, soil organic matter (log-scaled), and oven-dry
379 bulk density. These samples were not co-located with the training samples used to
380 generate the prior maps (samples at the same locations were already removed). The
381 number of observations varies by soil property: pH has the most samples, followed by
382 oven-dry bulk density and soil organic matter. The sample sizes across properties can also
383 be inferred from the frequency histograms shown in the lower-left corner of each panel in
384 Figure 3. Across all depths combined, each soil property has more than 11000
385 observations in California. The number of observations generally decreases with depth,
386 with depths below 1 m having notably fewer samples compared to shallower layers.

387



388

389

390

391

392

393

394

395

396

397

398

399

400

Figure 3: Spatial distribution and statistical characteristics of soil properties observations across California. The figure presents six soil parameters mapped using an Albers Equal Area projection: (a) sand content (mass %), (b) silt content (mass %), (c) clay content (mass %), (d) pH, (e) soil organic matter (log-scaled mass %), and (f) bulk density (g/cm³). Each subplot displays sample locations as colored points, with field-collected samples shown as triangles to distinguish them from WoSIS (circles) and SCD (squares) samples. Distribution histograms in the lower left corner of each subplot show the frequency distribution of values, with blue dashed lines indicating median values. Distance scale bar and compass rose are provided in the right corner. Note that the total number of soil measurements varies by property and generally decreases with depth beyond the surface layer, with the surface layers and depths below 1 m generally having fewer observations.

401 **2.2.1.1 World Soil Information Service (WoSIS)**

402 The World Soil Information Service (WoSIS), managed by the International Soil Reference
403 and Information Centre (ISRIC), aggregates global soil data from diverse sources, including
404 national soil institutes, research organizations, and collaborative initiatives like the Global
405 Soil Partnership (GSP) and the International Network of Soil Information Institutions (INSII).
406 The database provides soil properties for different soil horizons, georeferenced in decimal
407 degrees, and undergoes quality controls (Batjes et al., 2024). In California, WoSIS typically
408 offers 2,000 to over 5,000 soil observations for the modeling soil property. Samples below
409 1-m depth are fewer than those from shallower layers.

410

411 **2.2.1.2 Soil Characterization Database (SCD)**

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

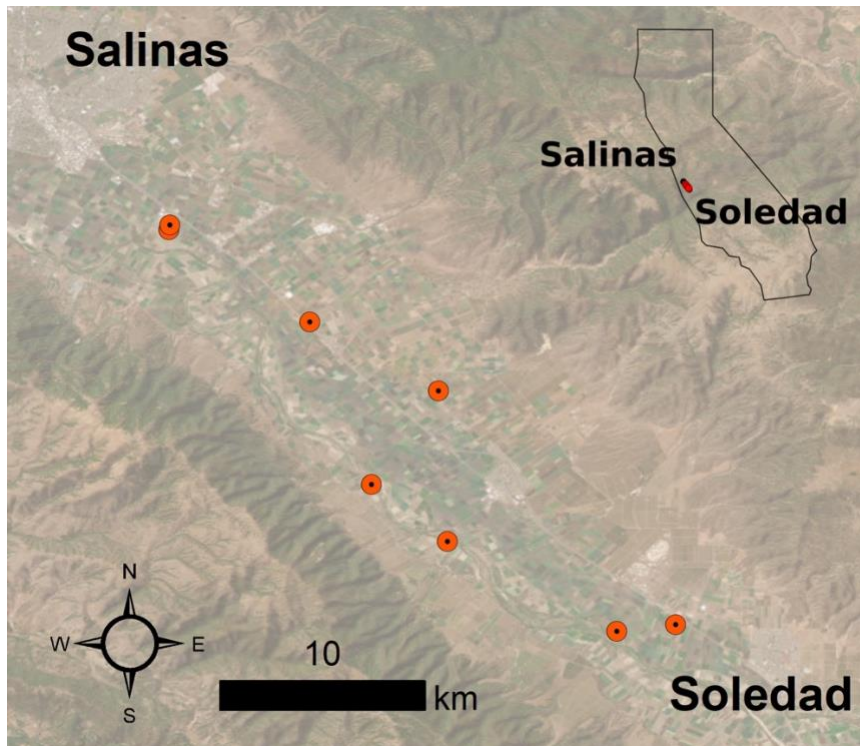
418

419 **2.2.1.3 Ground Truth Soil Sampling and Measurements**

420 Additional soil sampling was conducted to complement georeferenced soil profiles in
421 California for model training and evaluation. These data are reported in (Scudiero et al.,
422 2024) and are briefly discussed here. Multiple fields located between Salinas and Soledad

423 in California's Salinas Valley were selected to collect soil particle size fraction data (Figure
424 4). These fields, presented as red dots in Figure 4, were chosen because they were
425 accessible, unfarmed during the sampling period, and spread across different parts of the
426 valley.

427



428

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

434

435 Soil apparent electrical conductivity (ECa) was measured across fields using an
436 electromagnetic induction (EMI) sensor connected to a GPS receiver. Following the ECa-
437 directed soil sampling protocols of Corwin and Scudiero (Corwin and Scudiero, 2020), the
438 most representative soil samples were identified with ESAP software package and the
439 Response Surface Sampling Design algorithm (Lesch et al., 2000; Lesch, 2005). 0-0.8 and
440 0-1.6 m soil profiles were further analyzed and followed with the expectation that ECa was
441 a regional proxy for the field-scale variability of particle size fraction.

442
443 To measure particle size fraction, soil samples were then collected from multiple depths
444 (0–0.1, 0.1–0.4, and 0.4–1.2 m) across fields. After collection, the samples were air-dried,
445 ground, and sieved to remove particles larger than 2 mm; and then measured using the
446 Integral Suspension Pressure method (The improved integral suspension pressure method
447 (ISP+) for precise particle size analysis of soil and sedimentary materials; Wolfgang Durner,
448 Sascha C. Iden) using PARIO™ system (METER Group AG, Munich, Germany).

449

450 **2.2.2 Model Implementation for the California Case Study**

451 For the California case study, prior soil property maps were generated using the pruned
452 hierarchical Random Forest (pHRF) method (Xu et al., 2025). The pHRF-derived soil maps
453 were developed with soil pedons from the National Soil Information System (NASIS) and
454 part of SCD (the remaining data not used in IRC method). After gaining prior estimate of soil
455 properties, the IRC method was then applied using the additional soil observations from
456 WoSIS, SCD, and field measurements, which were not used in generating the prior maps.

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

459

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

466

467 **3 Results**

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

475

476 **3.1 Performance Evaluation of Posterior Soil Properties**

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

479 The metrics include the root mean square error (RMSE), coefficient of determination (R^2),
480 and correlation coefficient (ρ). For example, sand prediction shows an RMSE of 9.322, an
481 R^2 of 0.841, and a correlation coefficient of 0.918. pH prediction shows an RMSE of 0.270,
482 an R^2 of 0.945, and a correlation coefficient of 0.972. These metrics are computed using
483 out-of-bag (OOB) samples from random forest regressors. OOB samples are data points
484 not included in the bootstrap samples used to train each tree in the random forest.
485 Additionally, these metrics are evaluated by comparing the expected values of posterior
486 predictions with co-located soil properties values; not computed on residuals.

487

488 Table 1 also shows variations in performance across different soil properties. SOM and
489 bulk density show slightly worse metrics compared to particle size fractions and pH. For
490 instance, SOM predictions have an RMSE of 1.961, an R^2 of 0.608, and a correlation
491 coefficient of 0.801, and bulk density predictions have an RMSE of 0.164, an R^2 of 0.704,
492 and a correlation coefficient of 0.843. Two main reasons can result in their lower
493 performance. First, these properties are more dynamic in nature compared to particle size
494 fractions and pH. SOM and bulk density can change over time due to factors such as land
495 use practices. The prior predictions are trained using soil survey data that are older, while
496 the posterior soil profiles used for evaluation may come from a different period. Second,
497 SOM and bulk density are more challenging to model accurately. SOM is influenced by
498 complex biological and soil-forming processes, such as decomposition rates and organic
499 matter inputs. Similarly, bulk density is affected by soil compaction, organic matter
500 content, and soil structure. All of them can vary spatially and temporally. Depth-wise

501 analysis of model performance is provided in the Supplementary Information (Table S1 and
 502 S2).

503

504 **Table 1: Performance metrics (RMSE, R^2 , and correlation coefficient ρ) for posterior**
 505 **predictions of soil properties, including sand, silt, clay, pH, oven-dry bulk density**
 506 **(BD), and soil organic matter (SOM). The table summarizes the range (minimum and**
 507 **maximum values) and accuracy metrics for each property averaged across all depth**
 508 **intervals.**

Property	Unit	Min	Max	RMSE	R^2	ρ
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

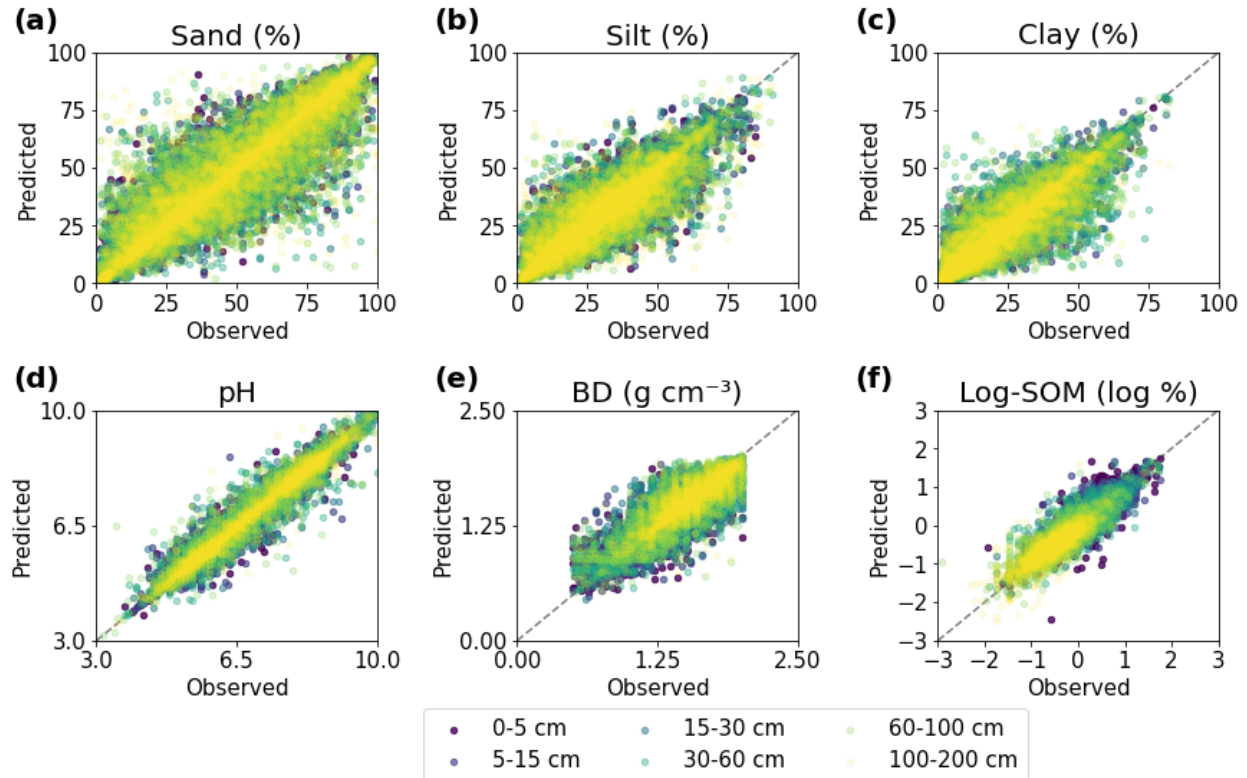
509

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

514 other predicted soil properties. And this pattern of reduced accuracy persists throughout
515 all depths.

516

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



528

529 **Figure 5: Evaluating posterior predictions with observations for six soil properties: (a)**

530 **sand, (b) silt, (c) clay, (d) pH, (e) bulk density (BD), and (f) log-scaled soil organic**

531 **matter (SOM). The left side shows scatter plots of posterior predictions versus**

532 **observations across six depth intervals, with each depth represented by a distinct**

533 **color. The dashed black line represents perfect prediction.**

534

535 **3.2 Comparison of Prior and Posterior Soil Predictions**

536 Prior and posterior predictions of soil properties are compared against co-located

537 observations to assess the added value of residual correction. The radar plots in Figure 6

538 illustrate the improvements achieved through the residual correction method using three

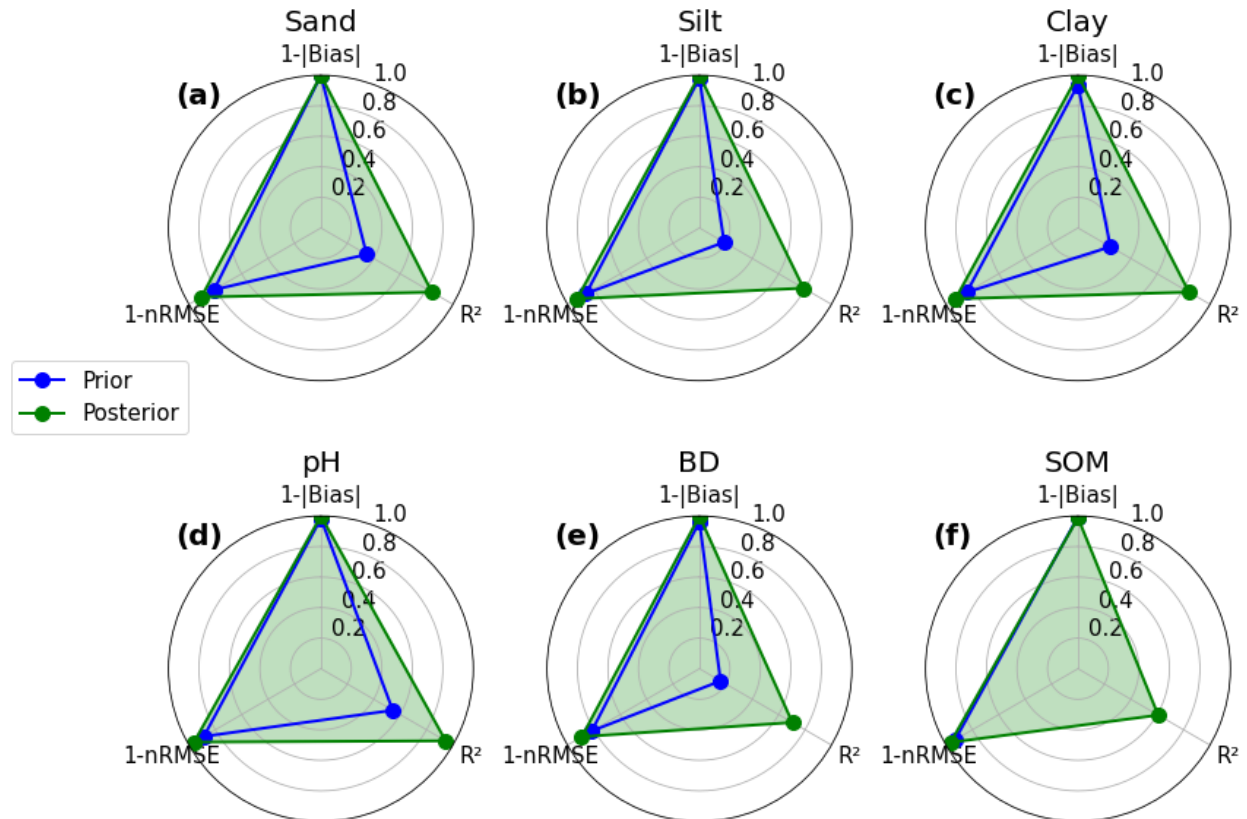
539 normalized metrics: 1-normalized absolute bias ($1-|\text{Bias}|$), coefficient of determination

540 (R^2), and 1-normalized RMSE by ranges of soil variability ($1-n\text{RMSE}$). These metrics are

541 computed with values of soil properties, instead of on their residuals. Values in Figure 6
542 closer to the outer edge of each plot indicate better model performance. Overall, all soil
543 properties maintain reasonable normalized bias and nRMSE (with nRMSE values
544 consistently less than 0.02 for both prior and posterior predictions). However, the prior
545 predictions tend to underestimate the variability of soil properties. As a result, the
546 normalized metrics for prior and posterior predictions are similar, while the R^2 values show
547 some differences.

548
549 For all soil properties, posterior predictions consistently outperform prior predictions
550 across all metrics. For particle size fractions, R^2 values show the largest improvements:
551 sand increases from 0.35 to 0.84, silt from 0.19 to 0.79, and clay from 0.25 to 0.84. The
552 nRMSE metric also shows improvements. Sand decreases from 0.19 to 0.09, silt from 0.14
553 to 0.07, and clay from 0.16 to 0.07, showing reductions in prediction errors using the
554 residual correction.

555
556 Aggregating data from all depths, Figure 6 shows the degree of improvement across
557 different soil properties. Prior pH predictions already demonstrate reasonable accuracy,
558 with an R^2 of 0.54 and nRMSE of 0.11. After the residual correction, these metrics improve
559 to 0.94 for R^2 and 0.04 for nRMSE. Bulk density and SOM show the biggest gains. For bulk
560 density, the R^2 increasing from 0.16 to 0.70 and nRMSE reducing from 0.18 to 0.11. Prior
561 SOM are underfitted with a low R^2 value. With the residual correction, the posterior SOM
562 show a positive R^2 of 0.61. The nRMSE for SOM also improves from 0.07 to 0.04.



563

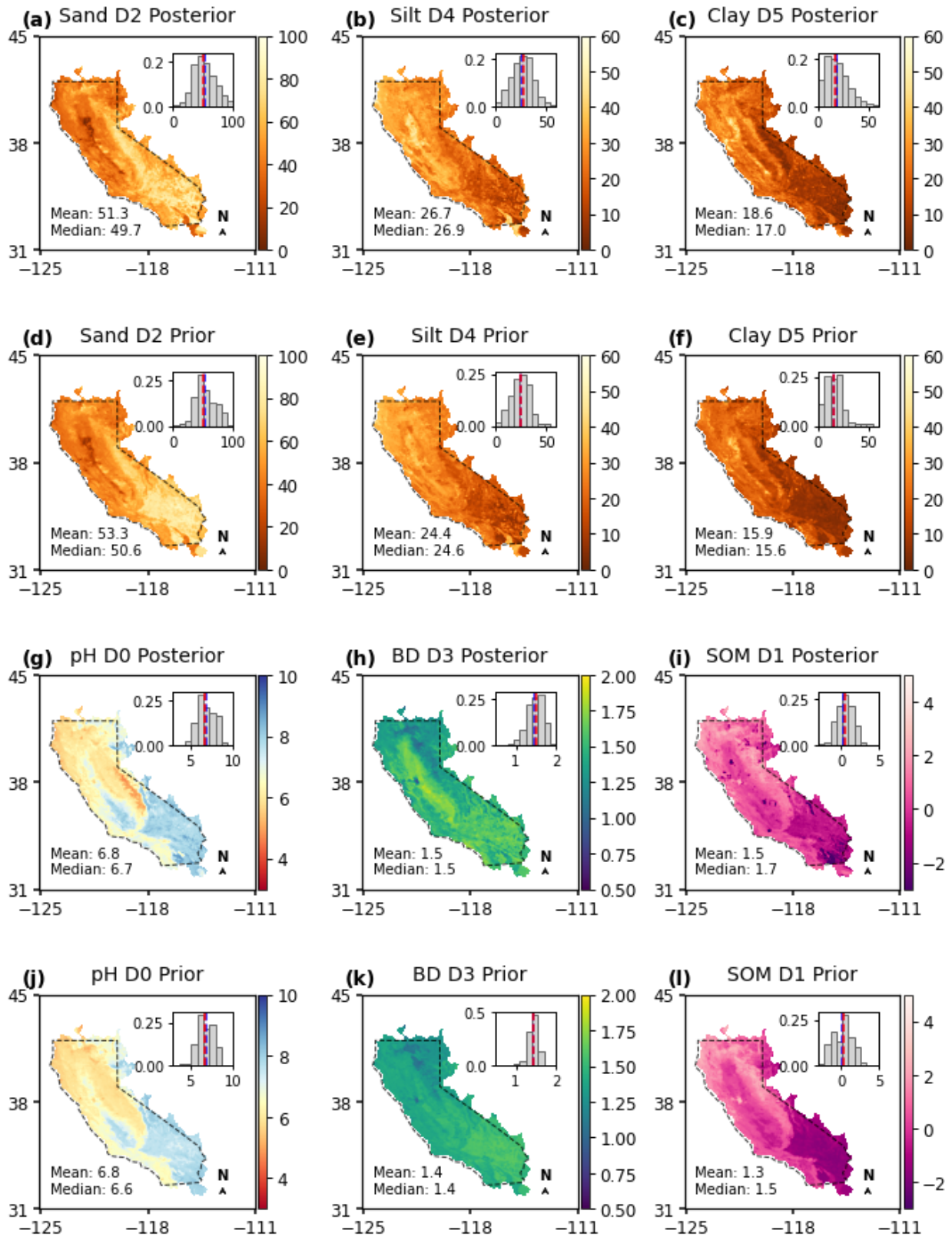
564 **Figure 6: Radar plots comparing the performance metrics of prior and posterior**
 565 **predictions for six soil properties: (a) sand, (b) silt, (c) clay, (d) pH, (e) oven-dry bulk**
 566 **density (BD), and (f) soil organic matter (SOM). Each plot presents three metrics: 1-**
 567 **normalized absolute bias ($1-|\text{Bias}|$), coefficient of determination (R^2), and 1-**
 568 **normalized RMSE by ranges of soil variability ($1-n\text{RMSE}$). Prior predictions are shown**
 569 **in blue, and posterior predictions in green. All metrics are scaled from 0 to 1, where**
 570 **values closer to the outer edge of the plot indicate better model performance. The**
 571 **green shaded area highlights the improvement achieved by the posterior predictions**
 572 **over prior estimates.**

573

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

580

581 For SOM and bulk density, the prior predictions often underestimate the observed
582 variation. Figure 7 shows that the residual correction processes add noticeable spatial
583 variations between prior and posterior soil maps. Prior bulk density values are often
584 clustered around 1.5 g/cm^3 , whereas the posterior histogram presents a broader range,
585 spanning from 1.25 g/cm^3 to 1.6 g/cm^3 , capturing more heterogeneity of bulk density.
586 Similarly, the residual correction adds soil heterogeneity to SOM. The posterior SOM can
587 delineate water bodies, where SOM content is abruptly lower than the surrounding areas.
588 Additionally, the posterior SOM maps present hill features in the desert areas.



589

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

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

599

600 Soil profiles used for evaluating residual correction are grouped according to their
601 corresponding pixel's land use classification from the National Land Cover Database
602 (NLCD). Figure 8 presents selected vertical soil profiles of sand content, oven-dry bulk
603 density, and SOM across three land use categories: forest, cultivated crops, and wetland.
604 The number of samples varies by land use, with forests having the most, cultivated crops
605 approximately half as many, and wetlands the fewest across California. To ensure a
606 balanced visualization, a similar number of profiles are selected from each category. Sand
607 content is chosen due to its broader range of variation (0-100%) compared to silt and clay
608 (< 60% range). SOM and bulk density, which show relatively lower performance metrics,
609 are included to assess the model's 'lower-bound performance'. These vertical profiles
610 were not used during model training.

611

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

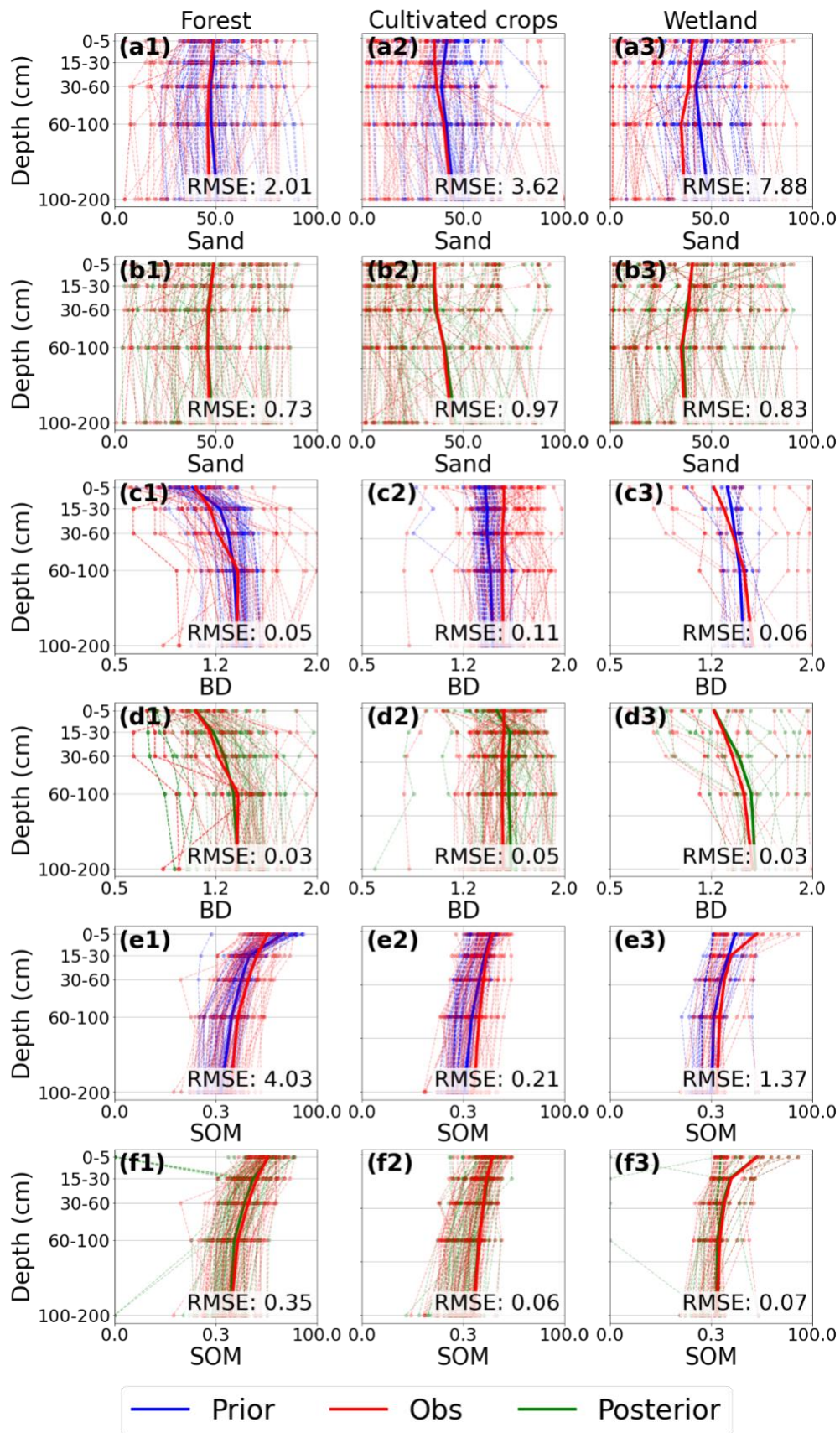
620

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

628

629 Dashed lines in Figure 8 represent individual soil profiles. Prior predictions often
630 underestimated the variability in soil properties, struggling to capture extreme values. After
631 the residual correction, the posterior predictions are better able to approximate these
632 extremes. However, the correction process sometimes introduces additional noise. For
633 example, some low SOM values (such as 0.001 g/cm^3) were generated during residual

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



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

645

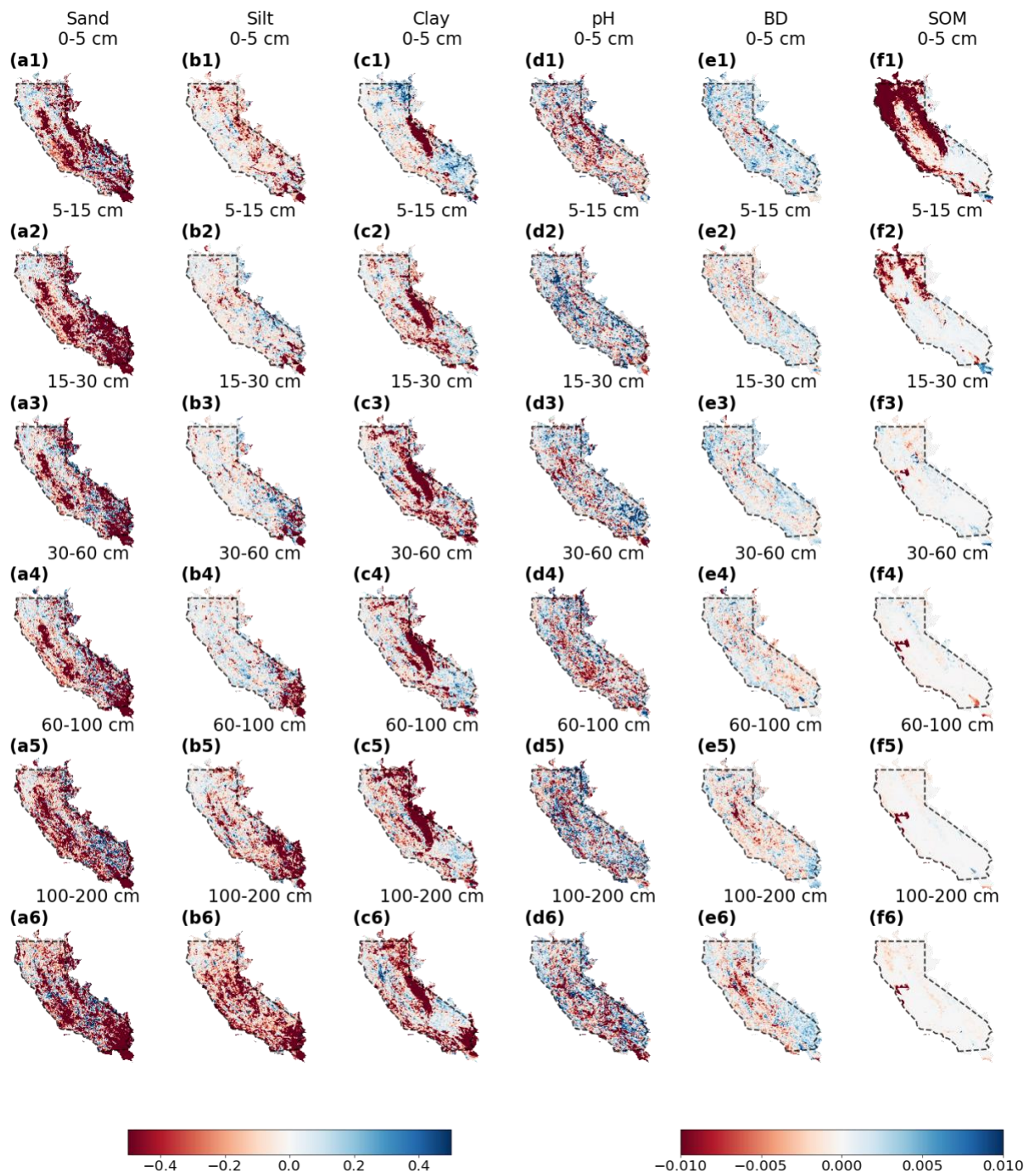
646 **3.3 Uncertainty Analysis**

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

654

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

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



665

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

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

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

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

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

675

676 **4 Discussion**

677 **4.1 Limitations in Soil Profile Data**

678 The effectiveness of residual correction depends on the spatial and vertical distribution of
679 soil profiles used to calculate residuals. In regions with sparse sampling, such as
680 California's desert areas (Figure 1), the limited number of profiles leads to interpolating the
681 entire area using limited observations. If soil heterogeneity is not captured by these limited
682 samples, the residual correction would overlook it. For soil texture, most data collected by
683 staff working on multiple projects under the National Institute of Food and Agriculture
684 (NIFA) and the Sustainable Agricultural Systems (SAS) programs range from the surface to
685 1.1 meters deep (additional field measurements used in this work). We use spline
686 interpolation to predict soil texture data beyond 1.1-m depths. It assumes vertical
687 continuity in soil properties, which may not reflect abrupt changes in subsurface layers.

688

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

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

698

699 **4.2 Computational Challenges**

700 The iterative residual correction process on distributions requires computational
701 resources, particularly when applied to large-extent or high-resolution datasets. This
702 process involves adjusting multiple values for each pixel, as each pixel represents a
703 distribution of soil properties. This process can be approached in two ways. The first
704 method involves correcting the residual values for each pixel, adding these residuals to
705 update the posterior values of soil properties, and then converting these updated values to
706 generate a posterior distribution of soil properties. The second method first converts all
707 pixel values into the same histogram bins and then corrects the shape of these histogram
708 bins for each pixel. Thus, the number of values retained per pixel affects computational
709 expense. Based on our experience, using method two, especially for soil texture, requires
710 100-bin histograms. Using method one with 20 most probable prior property values for
711 residual correction can achieve comparable results while reducing memory usage.

712

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

721

722 **4.3 Temporal and Spatial Constraints**

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

729

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

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

737
738 (2) Independent dataset evaluation: Use independent datasets to evaluate the
739 model. CONUS-wide instrumental network, such as the U.S. Climate Reference
740 Network and the National Ecological Observatory Network, provide independent
741 soil data. However, these datasets have limitations as they were collected with
742 clustering to certain landscapes, potentially introducing bias in the evaluation.

743

744 **4.4 Similar Studies**

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

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

758
759 Development of the following DSM products incorporates quantitative models in their
760 methodology. POLARIS produces probabilistic soil property maps using machine learning
761 and the DSMART algorithm (Chaney et al., 2016, 2019; Odgers et al., 2015), while the
762 uncertainties in the DSMART algorithm can propagate into POLARIS. SOLUS integrates
763 legacy soil data with georeferenced field observations and employs linear adjusted
764 Random Forest to predict soil properties (Nauman et al., 2024). SOLUS hierarchizes soil
765 data with different qualities into its training dataset, giving more attention to georeferenced
766 observations. However, since it also uses resampled soil data derived from polygon-based
767 soil map units, this process may introduce additional uncertainties into the final product.
768 The pHRF with IRC follows a different approach. Unlike most DSM methods that directly
769 predict soil properties from input data, this approach works in two steps: first, it generates
770 a prior estimate of soil taxa and property values, then iteratively adjusts these estimates to
771 improve model performance. In future work, the pHRF with IRC method will be applied on
772 large scale and assessed with more soil properties to evaluate its generalizability.

773

774 **5 Conclusion**

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

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

788

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

795

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

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

809

810 **Data Availability**

811 Data will be made available on request. Code is available on
812 https://github.com/emmaxu43/IRC_CA/tree/main.

813

814 **Author Contributions**

815 Chengcheng Xu and Nathaniel Chaney designed the study and developed the
816 methodology. Chengcheng Xu wrote the original draft and wrote the codes to produce the
817 methodology and analyses. Nathaniel Chaney supervised the work, provided resources
818 and funding, and helped guide the research direction. Elia Scudiero provided funding,
819 project management, co-supervision. Elia Scudiero and Ray Anderson provided soil
820 property samples from California that were used as part of the input dataset. Chengcheng

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

823

824 **Competing Interests**

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

826

827 **Acknowledgements**

828 This research was supported by the Agriculture and Food Research Initiative Competitive
829 Grant no. 2020-69012-31914 from the USDA National Institute of Food and Agriculture. The
830 authors want to thank Dr. Todd Skaggs for his and his teams' support for gathering input
831 data for this work. His and Dr. Ray Anderson's efforts are supported by USDA-ARS, Office
832 of National Programs (projects 2036-61000-019-000-D and 2036-61000-019-006-R). The
833 U.S. Department of Agriculture prohibits discrimination in all its programs and activities on
834 the basis of race, color, national origin, age, disability, and where applicable, sex, marital
835 status, familial status, parental status, religion, sexual orientation, genetic information,
836 political beliefs, reprisal, or because all or part of an individual's income is derived from
837 any public assistance program (not all prohibited bases apply to all programs). Persons
838 with disabilities who require alternative means for communication of program information
839 (braille, large print, audiotape, etc.) should contact USDA's TARGET Center at (202) 720-
840 2600 (voice and TDD). To file a complaint of discrimination, write to USDA, Director, Office
841 of Civil Rights, 1400 Independence Avenue, S.W., Washington, D.C. 20250-9410, or call

842 (800) 795-3272 (voice) or (202) 720-6382 (TDD). USDA is an equal opportunity provider and
843 employer.

844

845 **Financial Support**

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

848

849 **References**

850 Arrouays, D., McKenzie, N., Hempel, J., Forges, A. R. de, and McBratney, A. B.:
851 GlobalSoilMap: Basis of the global spatial soil information system, CRC Press, 496 pp.,
852 2014.

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

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

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

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

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

868 Corwin, D. L. and Scudiero, E.: Field-scale apparent soil electrical conductivity, Soil
869 Science Society of America Journal, 84, 1405–1441, <https://doi.org/10.1002/saj2.20153>,
870 2020.

871 Grunwald, S., Thompson, J. A., and Boettinger, J. L.: Digital Soil Mapping and Modeling at
872 Continental Scales: Finding Solutions for Global Issues, *Soil Science Society of America*
873 *Journal*, 75, 1201–1213, <https://doi.org/10.2136/SSSAJ2011.0025>, 2011.

874 Haghverdi, A., Najarchi, M., öztürk, H. S., and Durner, W.: Studying unimodal, bimodal, PDI
875 and bimodal-PDI variants of multiple soil water retention models: I. Direct model fit using
876 the extended evaporation and dewpoint methods, *Water (Switzerland)*, 12,
877 <https://doi.org/10.3390/w12030900>, 2020.

878 Hartemink, A. E., Hempel, J., Lagacherie, P., McBratney, A., McKenzie, N., MacMillan, R. A.,
879 Minasny, B., Montanarella, L., de Mendonça Santos, M. L., Sanchez, P., Walsh, M., and
880 Zhang, G.-L.: GlobalSoilMap.net – A New Digital Soil Map of the World, in: *Digital Soil*
881 *Mapping: Bridging Research, Environmental Application, and Operation*, edited by:
882 Boettinger, J. L., Howell, D. W., Moore, A. C., Hartemink, A. E., and Kienast-Brown, S.,
883 Springer Netherlands, Dordrecht, 423–428, [https://doi.org/10.1007/978-90-481-8863-](https://doi.org/10.1007/978-90-481-8863-5_33)
884 [5_33](https://doi.org/10.1007/978-90-481-8863-5_33), 2010.

885 Hengl, T., Heuvelink, G. B., and Stein, A.: A generic framework for spatial prediction of soil
886 variables based on regression-kriging, *Geoderma*, 120, 75–93, 2004.

887 Hengl, T., De Jesus, J. M., Heuvelink, G. B. M., Gonzalez, M. R., Kilibarda, M., Blagotić, A.,
888 Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas,
889 R., MacMillan, R. A., Batjes, N. H., Leenaars, J. G. B., Ribeiro, E., Wheeler, I., Mantel, S.,
890 and Kempen, B.: SoilGrids250m: Global gridded soil information based on machine
891 learning, *PLoS ONE*, <https://doi.org/10.1371/journal.pone.0169748>, 2017.

892 Jiang, Q., Fu, Q., and Wang, Z.: Delineating site-specific irrigation management zones,
893 *Irrigation and Drainage*, 60, 464–472, <https://doi.org/10.1002/ird.588>, 2011.

894 Lesch, S., Rhoades, J., and Corwin, D.: ESAP-95 version 2.01 R: User manual and tutorial
895 guide, *Research Rpt*, 146, 17, 2000.

896 Lesch, S. M.: Sensor-directed response surface sampling designs for characterizing spatial
897 variation in soil properties, *Computers and Electronics in Agriculture*, 46, 153–179,
898 <https://doi.org/10.1016/j.compag.2004.11.004>, 2005.

899 Li, N., Zhao, X., Wang, J., Sefton, M., and Triantafilis, J.: Digital soil mapping based site-
900 specific nutrient management in a sugarcane field in Burdekin, *Geoderma*, 340, 38–48,
901 <https://doi.org/10.1016/j.geoderma.2018.12.033>, 2019.

902 McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping,
903 *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.

904 Minasny, B. and McBratney, A. B.: A conditioned Latin hypercube method for sampling in
905 the presence of ancillary information, *Computers & geosciences*, 32, 1378–1388, 2006.

906 Mueller, T. G., Pierce, F. J., Schabenberger, O., and Warncke, D. D.: Map Quality for Site-
907 Specific Fertility Management, *Soil Science Society of America Journal*, 65, 1547–1558,
908 <https://doi.org/10.2136/sssaj2001.6551547x>, 2001.

909 National Cooperative Soil Survey: NCSS Soil Characterization Database (Lab Data Mart),
910 2018.

911 Nauman, T. W., Kienast-Brown, S., Roecker, S. M., Brungard, C., White, D., Philippe, J., and
912 Thompson, J. A.: Soil landscapes of the United States (SOLUS): Developing predictive soil
913 property maps of the conterminous United States using hybrid training sets, *Soil Science
914 Society of America Journal*, 88, 2046–2065, <https://doi.org/10.1002/saj2.20769>, 2024.

915 Nussbaum, M., Zimmermann, S., Walthert, L., and Baltensweiler, A.: Benefits of
916 hierarchical predictions for digital soil mapping—An approach to map bimodal soil pH,
917 *Geoderma*, 437, 116579, <https://doi.org/10.1016/j.geoderma.2023.116579>, 2023.

918 Odgers, N. P., McBratney, A. B., and Minasny, B.: Digital soil property mapping and
919 uncertainty estimation using soil class probability rasters, *Geoderma*, 237,
920 <https://doi.org/10.1016/j.geoderma.2014.09.009>, 2015.

921 Oliver, M. A. and Webster, R.: A tutorial guide to geostatistics: Computing and modelling
922 variograms and kriging, *CATENA*, 113, 56–69,
923 <https://doi.org/10.1016/j.catena.2013.09.006>, 2014.

924 Ortuani, B., Chiaradia, E. A., Priori, S., L'Abate, G., Canone, D., Comunian, A., Giudici, M.,
925 Mele, M., and Facchi, A.: Mapping Soil Water Capacity Through EMI Survey to Delineate
926 Site-Specific Management Units Within an Irrigated Field, *Soil Science*, 181, 252,
927 <https://doi.org/10.1097/SS.000000000000159>, 2016.

928 Poggio, L., De Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E., and
929 Rossiter, D.: SoilGrids 2.0: Producing soil information for the globe with quantified spatial
930 uncertainty, *SOIL*, 7, 217–240, <https://doi.org/10.5194/SOIL-7-217-2021>, 2021.

931 Powers, J. S., Corre, M. D., Twine, T. E., and Veldkamp, E.: Geographic bias of field
932 observations of soil carbon stocks with tropical land-use changes precludes spatial
933 extrapolation, *Proceedings of the National Academy of Sciences*, 108, 6318–6322,
934 <https://doi.org/10.1073/pnas.1016774108>, 2011.

935 Ramcharan, A., Hengl, T., Nauman, T., Brungard, C., Waltman, S., Wills, S., and Thompson,
936 J.: Soil Property and Class Maps of the Conterminous United States at 100-Meter Spatial
937 Resolution, *Soil Science Society of America Journal*, 82, 186–201,
938 <https://doi.org/10.2136/sssaj2017.04.0122>, 2018.

939 Rossiter, D. G., Poggio, L., Beaudette, D., and Libohova, Z.: How well does digital soil
940 mapping represent soil geography? An investigation from the USA, *SOIL*, 8, 559–586,
941 <https://doi.org/10.5194/soil-8-559-2022>, 2022.

- 942 Schmidinger, J. and Heuvelink, G. B. M.: Validation of uncertainty predictions in digital soil
943 mapping, *Geoderma*, 437, 116585, <https://doi.org/10.1016/j.geoderma.2023.116585>,
944 2023.
- 945 Scudiero, E., Corwin, D. L., Markley, P. T., Pourreza, A., Rounsville, T., Bughici, T., and
946 Skaggs, T. H.: A system for concurrent on-the-go soil apparent electrical conductivity and
947 gamma-ray sensing in micro-irrigated orchards, *Soil and Tillage Research*, 235, 105899,
948 2024.
- 949 Sharififar, A., Sarmadian, F., Malone, B. P., and Minasny, B.: Addressing the issue of digital
950 mapping of soil classes with imbalanced class observations, *Geoderma*, 350, 84–92,
951 <https://doi.org/10.1016/j.geoderma.2019.05.016>, 2019.
- 952 Shi, G., Sun, W., Shangguan, W., Wei, Z., Yuan, H., Zhang, Y., Liang, H., Li, L., Sun, X., Li, D.,
953 Huang, F., Li, Q., and Dai, Y.: A China dataset of soil properties for land surface modeling
954 (version 2), <https://doi.org/10.5194/essd-2024-299>, 29 August 2024.
- 955 Soil, K.: Survey laboratory methods manual, Soil Survey Investigations Report, 1996.
- 956 Soil Survey Staff: Kellogg Soil Survey Laboratory methods manual, U.S. Department of
957 Agriculture, Natural Resources Conservation Service, Lincoln, Nebraska, 2014.
- 958 Soil survey staff: Gridded National Soil Survey Geographic (gNATSGO) Database for the
959 Conterminous United States, 2023. Natural Resources Conservation Service, United
960 States Department of Agriculture.
- 961 Soil Survey Staff: Gridded Soil Survey Geographic (gSSURGO) Database for the
962 Conterminous United States, 2025. Natural Resources Conservation Service, United
963 States Department of Agriculture.
- 964 Soil Survey Staff, Natural Resources Conservation Service, and United States Department
965 of Agriculture: Soil Survey Geographic (SSURGO) Database for the CONUS, 2023. Natural
966 Resources Conservation Service, United States Department of Agriculture.
- 967 Sylvain, J.-D., Anctil, F., and Thiffault, É.: Using bias correction and ensemble modelling for
968 predictive mapping and related uncertainty: A case study in digital soil mapping,
969 *Geoderma*, 403, 115153, <https://doi.org/10.1016/j.geoderma.2021.115153>, 2021.
- 970 Takoutsing, B., Heuvelink, G. B. M., Stoorvogel, J. J., Shepherd, K. D., and Aynekulu, E.:
971 Accounting for analytical and proximal soil sensing errors in digital soil mapping, *European
972 Journal of Soil Science*, 73, e13226, <https://doi.org/10.1111/ejss.13226>, 2022.
- 973 Vereecken, H., Schnepf, A., Hopmans, J. W., Javaux, M., Or, D., Roose, T., Vanderborght, J.,
974 Young, M. H., Amelung, W., Aitkenhead, M., Allison, S. D., Assouline, S., Baveye, P., Berli,
975 M., Brüggemann, N., Finke, P., Flury, M., Gaiser, T., Govers, G., Ghezzehei, T., Hallett, P.,
976 Hendricks Franssen, H. J., Heppell, J., Horn, R., Huisman, J. A., Jacques, D., Jonard, F.,

977 Kollet, S., Lafolie, F., Lamorski, K., Leitner, D., McBratney, A., Minasny, B., Montzka, C.,
978 Nowak, W., Pachepsky, Y., Padarian, J., Romano, N., Roth, K., Rothfuss, Y., Rowe, E. C.,
979 Schwen, A., Šimůnek, J., Tiktak, A., Van Dam, J., van der Zee, S. E. A. T. M., Vogel, H. J.,
980 Vrugt, J. A., Wöhling, T., and Young, I. M.: Modeling Soil Processes: Review, Key
981 Challenges, and New Perspectives, *Vadose Zone Journal*, 15, vzt2015.09.0131,
982 <https://doi.org/10.2136/vzt2015.09.0131>, 2016.

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

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

991 Xu, C., Huang, J., Hartemink, A. E., and Chaney, N. W.: Pruned hierarchical Random Forest
992 framework for digital soil mapping: Evaluation using NEON soil properties, *Geoderma*, 459,
993 117392, <https://doi.org/10.1016/j.geoderma.2025.117392>, 2025.

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

996