

UC Merced

UC Merced Previously Published Works

Title

Soil Science-Informed Machine Learning

Permalink

<https://escholarship.org/uc/item/8qf2n8zs>

Authors

Minasny, Budiman

Bandai, Toshiyuki

Ghezzehei, Teamrat A

et al.

Publication Date

2024-12-01

DOI

10.1016/j.geoderma.2024.117094

Peer reviewed

1 **Soil Science-Informed Machine Learning**

2

3 Budiman Minasny¹, Toshiyuki Bandai², Teamrat A. Ghezzehei³; Yin-Chung Huang¹,
4 Yuxin Ma⁴, Alex. B. McBratney¹, Wartini Ng¹, Sarem Norouzi⁵, Jose Padarian¹,
5 Rudiyanto⁶, Amin Sharififar¹, Quentin Styc¹, Marliana Widyastuti¹

6

7

8

9 ¹ School of Life and Environmental Sciences, The University of Sydney, NSW 2006,
10 Australia.

11 ² Earth and Environmental Sciences Area, Lawrence Berkeley National Laboratory,
12 Berkeley, CA 94720, USA.

13 ³ Life & Environmental Sciences Department, University of California, Merced, CA
14 95343, USA.

15 ⁴ New South Wales Department of Climate Change, Energy, the Environment and
16 Water, Parramatta, NSW 2150, Australia.

17 ⁵ Department of Agroecology, Aarhus University, 8830 Tjele, Denmark.

18 ⁶ Faculty of Fisheries and Food Science, Universiti Malaysia Terengganu, 21030 Kuala
19 Nerus, Terengganu, Malaysia.

20

21 **Abstract**

22 Machine learning (ML) applications in soil science have significantly increased over
23 the past two decades, reflecting a growing trend towards data-driven research
24 addressing soil security. This extensive application has mainly focused on enhancing
25 predictions of soil properties, particularly soil organic carbon, and improving the
26 accuracy of digital soil mapping (DSM). Despite these advancements, the
27 application of ML in soil science faces challenges related to data scarcity and the
28 interpretability of ML models. There is a need for a shift towards Soil Science-
29 Informed ML (SoilML) models that use the power of ML but also incorporate soil
30 science knowledge in the training process to make predictions more reliable and
31 generalisable. This paper proposes methodologies for embedding ML models with
32 soil science knowledge to overcome current limitations. Incorporating soil science
33 knowledge into ML models involves using observational priors to enhance training
34 datasets, designing model structures which reflect soil science principles, and
35 supervising model training with soil science-informed loss functions. The informed
36 loss functions include observational constraints, coherency rules such as
37 regularisation to avoid overfitting, and prior or soil-knowledge constraints that
38 incorporate existing information about the parameters or outputs. By way of
39 illustration, we present examples from four fields: digital soil mapping, soil
40 spectroscopy, pedotransfer functions, and dynamic soil property models. We discuss
41 the potential to integrate process-based models for improved prediction, the use of
42 physics-informed neural networks, limitations, and the issue of overparametrisation.
43 These approaches improve the relevance of ML predictions in soil science and
44 enhance the models' ability to generalise across different scenarios while
45 maintaining soil science principles, transparency and reliability.

46 **1. Introduction**

47 The 2024 Nobel Prize in Physics was awarded to researchers who utilised physics-
48 based tools to develop methods that advance machine learning through artificial
49 neural networks. Over the past two decades, the use of machine learning (ML) in
50 soil research has surged. In 2023, an average of 8 papers per day were published on
51 topics related to “machine learning” and “soil” (Scopus, June 2024). These
52 advancements, highlight the growing importance of ML in various scientific fields,
53 including soil research. (See Box 1 for the definition of ML.) Past reviews show that
54 the application of ML in soil science spans various areas, including soil organic
55 carbon (SOC), hydrology, contamination, remote sensing, erosion, ML methods and
56 modelling, spectroscopy, and crops (Li et al., 2024; Padarian et al., 2020b). In
57 particular, the application of ML is extensive in digital soil mapping (DSM) studies.
58 The review by Wadoux et al. (2020) on ML in DSM indicated that most studies
59 emphasise predicting soil properties (in particular SOC) and improving prediction
60 accuracy. However, only a few studies account for existing soil knowledge in the
61 modelling processes.

62 Undoubtedly, machine learning has revolutionised the processing of large soil
63 databases, finding patterns which are difficult to uncover using traditional statistical
64 models (Heung et al., 2016). Soil observational data, collected via field and
65 laboratory techniques and numerous sensors, provide extensive datasets that
66 conventional statistical models may not efficiently handle (Safanelli et al., 2021;
67 Tziolas et al., 2020). ML models excel in discovering patterns within spatiotemporal
68 soil data, which are often challenging for process-based models to address. In
69 addition, ML facilitates the generation of detailed soil information, scaling from field-
70 level observations to global insights (Helfenstein et al., 2024; Padarian et al., 2022b;
71 Poggio et al., 2021; Rosin et al., 2023).

72 While ML can replicate observed patterns in training data, it often falls short of
73 explaining observed phenomena, and the learned patterns are usually not
74 generalisable. ML models require substantial volumes of data, yet soil data are
75 limited and sparse. The efficacy of ML models is constrained by the quantity and
76 quality of training data, hindering their ability to predict “unseen” phenomena.
77 (Read Box 2 “Six Dangers of ML in Soil Science”.) There is an ongoing discussion on

78 incorporating soil knowledge in ML models and the interpretability of the calibrated
79 ML models (Ma et al., 2019; Wadoux et al., 2020). There is growing interest in
80 applying interpretable ML models to explain how models predict certain attributes,
81 addressing the “black box” issue. Weindorf and Chakraborty (2024) argued for a
82 balance between ML modeling with human insight and knowledge for
83 contextualising findings, and ensuring the completeness, validity, and interpretation
84 of AI-generated results. Increasingly, there is a call to incorporate fundamental
85 domain knowledge and physical rules into ML models to enhance their reliability and
86 accuracy by providing theoretical constraints and informative priors (von Rueden et
87 al., 2023). Concurrently, there is a push to model soil biogeochemical processes
88 using physical rules, foundational to numerous achievements in computational
89 physics and chemistry (Tang et al., 2024). In physics, hydrology, and related fields,
90 there is increasing interest in physics-informed machine learning, aiming to guide
91 ML models towards solutions that are physically plausible (Karniadakis et al., 2021;
92 Kashinath et al., 2021). Notably, physics-informed ML models have been
93 investigated to model soil water movement (Bandai and Ghezzehei, 2022).

94 We propose Soil Science-Informed ML (SoilML), which integrates soil-specific
95 knowledge—including pedology, physical, chemical, and biological, processes into
96 ML models, expanding the scope beyond Physics-Informed ML (PIML). SoilML
97 prioritises modelling soil systems, accounting for interactions such as water cycling,
98 soil-water-plant dynamics, and biogeochemical transformations. This approach aims
99 to enhance model interpretability, improve predictions in data-scarce environments,
100 and ensure that outputs are consistent with real-world soil behavior.

101 The paper is structured as follows. Section 2 demonstrates ways for incorporating
102 soil science knowledge or principles in ML models, moving beyond merely
103 identifying important predictors. Subsequently in Section 3, we present specific
104 examples from four key fields: pedotransfer functions, digital soil mapping, soil
105 spectroscopy, and modelling soil properties in space and time. We discuss the
106 capabilities and limitations of conventional approaches and explore the potential to
107 integrate soil science knowledge under the SoilML framework, followed by
108 implications (Kashinath et al., 2021). Section 4 provides a discussion of SoilML
109 models to address the unique 3D structure and interactions in soil systems, the

110 issue of overfitting, enhancing model interpretability, reliability, and predictive
111 accuracy for soil science applications.

112

113

114

115

116

117

Box 1. Definitions

Artificial Intelligence (AI) is the field of study focused on designing and developing intelligent machines capable of performing tasks which mimic human intelligence.

Machine Learning (ML) is a subset of AI that uses algorithms to perform specific tasks without explicit instructions. The models learn and make predictions based on patterns and inferences derived from data, focusing on prediction accuracy. In this context, we do not consider statistical models such as linear regression and partial least squares regression as ML as they rely on predefined functional assumptions.

Deep learning is a subset of ML that involves a type of algorithm called artificial neural networks. These neural networks are designed to recognise patterns in data by processing data through multiple layers of processing units.

Physical rules are fundamental principles that describe how physical systems behave and interact in the natural world based on scientific observations, experiments, theories, and mathematical models.

Mechanistic or process-based models are mathematical models that describe one or multiple processes based on the underlying mechanisms and interactions

among system components. Based on the principles of physics, chemistry, biology, and related sciences, combined with empirical relationships, these models aim to represent how different components of a system work together to produce observed behaviours.

118

119

120

121

Box 2. Six Dangers of Machine Learning in Soil Science

(1) Data science of soil materials, ML models without soil science context

ML modelling often follows a workflow of processes that apparently do not require in-depth soil knowledge. A simple search on Google Scholar for "machine learning in soil classification" yields numerous papers, primarily from the computer science field, that apply ML techniques to predict soil types based on properties or images of soil. These studies often treat soil merely as a material with inadequately informed labels. This approach can result in information that lacks practical relevance and does not contribute to a deeper understanding of soil science or soil.

(2) Unscrutinised machine-learned soil prediction models

Defining the objective of an ML modelling exercise is essential. If the goal is to achieve the highest accuracy for a specific problem, interpretability may not be a priority. Considering the complexity of natural phenomena and human limitations in understanding complex relationships, demanding complete transparency from ML models may not be feasible. Nonetheless, it is crucial to ensure that the model provides a valid generalisation of the phenomenon being studied. Overreliance on automated outputs without sufficient scrutiny could lead to misinterpretations. Many papers in soil science literature use ML modelling without attempting interpretability, raising questions about the utility of such work for advancing soil

science knowledge.

(3) Lack of transparency in proprietary soil models

The utility of ML spans beyond research into commercial domains, such as predicting soil properties using near-infrared spectroscopy or an online platform for predicting soil properties using remote sensing images. This is particularly relevant in agronomic, soil contamination and soil carbon accounting applications. While research settings often require transparency regarding methods and data, commercial entities tend to be secretive to maintain a competitive edge, raising concerns about the reliability of their model predictions. Consequently, soil prediction becomes proprietary. Addressing these concerns may require implementing methods like uncertainty assessment, independent validation tests, and reporting on the diversity of soil types used in training the models.

(4) Stagnation of theoretical advancement

The focus on applying ML to predict and model soil properties and processes could overshadow the need to develop new theoretical frameworks in soil science. Without ongoing theoretical advancements, soil science may become overly dependent on data analytics without fostering innovative ideas. Theory generates hypotheses and generally leads to efficient experimentation and data generation. It's not clear at this stage that this is the case for ML-generated prediction models.

(5) Doing too much with too little

There is a tendency to produce regional or global maps of various interesting soil characteristics without acknowledging the limited data and the risk of extrapolating these models in areas where data is sparse or of poor quality, leading to unreliable or misleading results. Another tendency is solely relying on ML models to infer controls of soil properties prediction. This overextension can compromise the integrity of soil science research and its applications. Sensible guidelines are required for the data density required for such predictions; for example, global maps based on several hundred observations are probably questionable, whereas those based on tens of thousands of observations inspire more confidence.

(6) Decline in direct soil observations and human fieldwork

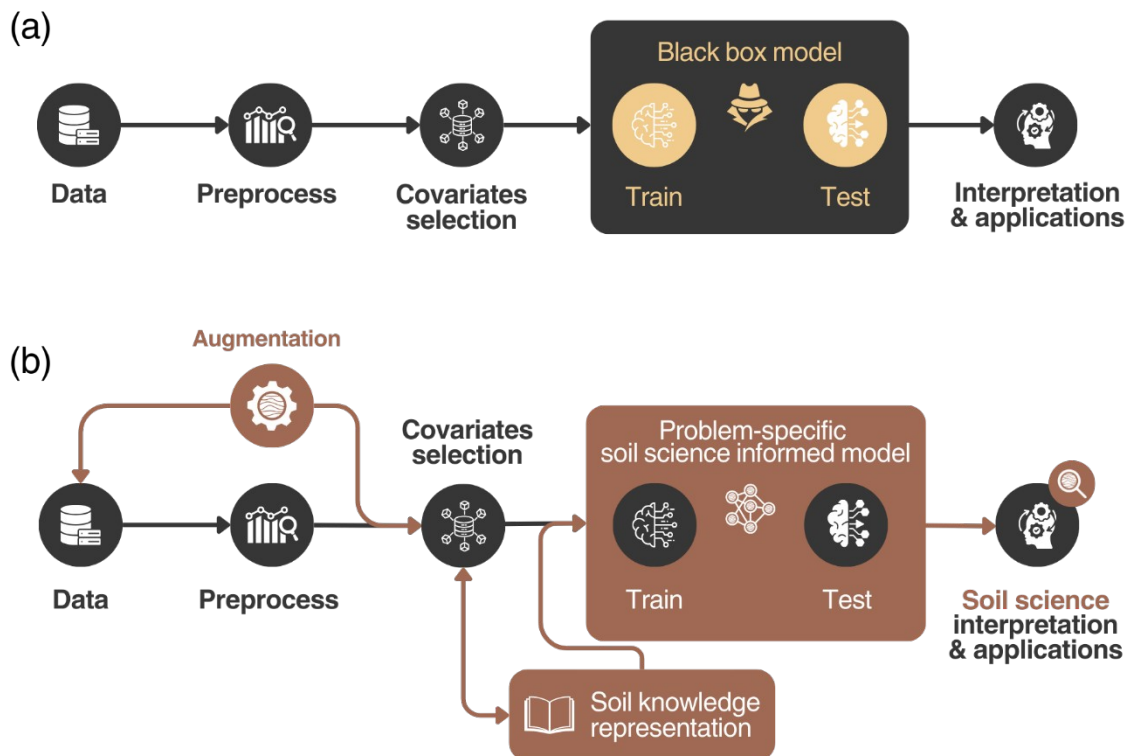
Overreliance on ML might lead to overconfidence and decreased fieldwork and the gathering of new observations of soil, which are important for understanding soil in its natural conditions and accurately interpreting data and models. This shift could reduce the practical understanding of soil conditions and processes, diminishing the empirical grounding of soil science. Since soil is dynamic and responds to human forcings, continued widespread real-world observation is essential. Often, modelling and prediction can be better improved by accruing new observations at key locations rather than through incremental improvements of new ML methods.

122

123

124 **2. Incorporating soil science knowledge in ML models**

125



126

127 Figure 1. (a) The general workflow for ML applications in soil science. (b) the
128 iterative and problem-specific nature of soil science-informed models. Conventional
129 steps are in black circles, with steps involving soil science supervision in brown
130 circles.

131

132 The steps for conducting an ML analysis on soil data typically follow the procedure
133 depicted in Figure 1(a). The process starts with data collection and pre-processing to
134 check for outliers, selection of covariates or predictors, followed by training of an ML
135 model, which could include tuning the hyperparameters. The model is then tested
136 on a proportion of the dataset which was not used in the training process. Finally,
137 the model is interpreted using procedures such as variables of importance or
138 Shapley values. (In ML, Shapley values quantify the relative contribution of each

139 predictor to a model's output, for more details, see Padarian et al. (2020a) and
140 Wadoux and Molnar (2022)). The entire procedure follows a standard ML workflow.

141 Most studies have utilised supervised learning, where ML algorithms are tasked with
142 predicting labels based on a set of features. However, the concept of 'supervision'
143 should extend beyond mere labels to encompass a broader prior knowledge
144 framework. This knowledge is often embodied in functions or sets of rules that may
145 depend on specific "labels". Such supervision incorporates domain-specific insights
146 that guide the learning process, enabling the algorithm to make more informed and
147 contextually appropriate predictions. ML models should not only learn from
148 observation data points but also integrate structured forms of knowledge
149 effectively, thereby enhancing their predictive accuracy and relevance to soil
150 science (see Figure 1(b), soil science supervision steps).

151 There is a lack of discussion on how to effectively incorporate soil science
152 knowledge or physical rules in ML models. Here, we argue that ML models need to
153 be iteratively designed and problem-specific, and they should be supervised to
154 predict patterns conforming to soil phenomena. SoilML could deliver predictive
155 models grounded in soil science, which not only achieve higher prediction accuracy
156 but also enhance the models' ability to generalise predictions. Additionally, SoilML
157 could improve transparency, thereby increasing the plausibility and reliability of
158 these models (Kashinath et al., 2021; Wadoux et al., 2020).

159

160 **2.1.1 Source of knowledge**

161 The fields of Informed Machine Learning and Physics-Informed Machine Learning
162 have emerged to address the empirical nature of current ML models by
163 incorporating prior knowledge into the training process. von Rueden et al. (2023)
164 discuss the source of knowledge, its representation, and how it is integrated into ML
165 algorithms (Figure 2). The source of knowledge can be in three forms: specific
166 scientific knowledge (in our case, soil science), general knowledge, and expert
167 knowledge. Soil science knowledge could include laws or equations, principles, and
168 rules. General knowledge is often intuitive and implicitly validated by human
169 reasoning or empirical studies. Expert knowledge tends to be based on experience,
170 for example, the relationship between certain soil properties and their covariates

171 (Lark et al., 2007) or mental models embedded in soil maps and their legends (Bui,
172 2004; Qi and Zhu, 2003).

173 These sources of knowledge can be represented as algebraic forms (e.g. Philip's
174 infiltration equation) or differential equations (Richardson-Richards equation),
175 simulation results, pseudo-observations, soil maps, rules, or probabilistic
176 relationships (Hudson, 1992). In turn, these forms of knowledge can be integrated
177 into the ML workflow through training data, model structure design, learning
178 algorithms, and final evaluation. Current approaches to incorporating soil science
179 knowledge in ML models involve several strategies. For example, in DSM, covariates
180 may be selected to conform to soil-forming factors or *scorpan* covariates (See
181 section 3.1). In soil moisture dynamics modelling, predictors could be selected
182 based on components of a water-balance model. Another idea is incorporating
183 pseudo-observations based on expert opinion on soil properties in areas which lack
184 observations such as in high-elevation areas or extreme environments. Finally,
185 interpretative tools such as Shapley value could interpret how predictors contribute
186 to ML predictions, aligning predictions with existing knowledge (Padarian et al.,
187 2020a; Wadoux and Molnar, 2022).

188

189 **2.1.2 Incorporating soil science knowledge in ML models**

190 Karniadakis et al. (2021) advocate three ways of incorporating soil science
191 knowledge in ML models: observational priors, model structure design, and learning
192 guidance (Figure 2).

193 **(1) Observational priors**

194 This approach involves augmenting training data to reflect underlying knowledge
195 about the subject. Expert knowledge is mostly represented in DSM studies,
196 including the addition of synthetic or pseudo-observations to the training data. DSM
197 commonly relies on legacy soil data derived from laboratory measurements, which
198 can be limited in spatial coverage. Field observations such as hand texture can
199 provide a dense and complementary source of soil data, capturing variability across
200 the landscape that laboratory data may miss. Eymard et al. (2024) demonstrated
201 that integrating field observations of soil texture, even with potential biases, can

202 improve DSM predictions by identifying unique landscape features not represented
203 in laboratory datasets, ultimately enhancing both model accuracy and the
204 understanding of soil processes.

205 In addition, soil survey can be spatially biased due to preferential sampling patterns,
206 and may have gaps in coverage due to inaccessible areas, such as steep terrain or
207 remote regions. For example, Koch et al. (2019) used 13,000 boreholes to map the
208 depth to the redox layer across Denmark using random forest regression kriging,
209 but found that lowland areas were underrepresented. To address this, synthetic
210 observations were added in these regions based on hydrogeological knowledge,
211 improving lowland representation. Similarly, outputs from soil process-based models
212 can be used to fill temporal gaps in observations, which will be discussed in Section
213 3.4. See also Box 3 on reducing overparametrisation in ML models via data
214 augmentation.

215

216 **(2)Model structure design**

217 The architecture of ML models should be designed to ensure that their predictions
218 are consistent with established soil science principles. This involves selecting
219 appropriate model types, designing input and output layers and connections that
220 can process and interpret soil data, and implementing mechanisms that incorporate
221 domain-specific knowledge into the learning process. For example, ML structure
222 needs to accommodate soil profile information. In the case of predicting soil at
223 multiple depths within a profile, a multitasking ML model that predicts soil
224 properties at multiple depths simultaneously would be preferable to creating
225 independent soil depth prediction functions (Padarian et al., 2019b). In another
226 example, conventional maps are updated or disaggregated using ML models with
227 expert knowledge inputs, such as defining soil-landscape conditions in which a
228 particular soil type could occur (Holmes et al., 2015; Lamichhane et al., 2021;
229 Odgers et al., 2014; van Zijl et al., 2019).

230

231 **(3)Learning guidance**

232 The training of ML models can be directed using loss functions and constraints to
233 ensure that the solutions align with soil science processes. Typically, ML models are

234 trained to minimise a loss or cost function; commonly, this involves adjusting the
235 model parameters to minimise the mean squared error between the observed and
236 predicted values.

237 Following Tang et al. (2024) we could define the loss function of an ML model as:

238

239 $L = \text{observational constraints} + \text{coherency rules} + \text{prior constraints} \quad (1)$

240

241 **Observational constraints** are usually defined as the mean squared error
242 between observed and predicted values for continuous variables or classification
243 error for categorical variables. For example, an ML model predicting the parameters
244 of a soil water retention function would minimise the difference between measured
245 and observed water retention at defined pressure heads.

246 **Coherency rules**, also known as regularisation or penalty function, aim to
247 constrain the parameters to obey physical processes related to model parameters,
248 thus avoiding overfitting. For instance, in water retention prediction, the relationship
249 between water content and pressure head must adhere to the monotonicity of the
250 curve (van Genuchten, 1980).

251 **Prior or knowledge-based constraints**, involve incorporating soil science and
252 general knowledge or assumptions about the parameters or outputs, guiding the
253 model towards more plausible solutions. For example, specifying ranges within
254 which certain parameters must lie based on prior studies or expert knowledge, and
255 imposing non-negativity constraints on parameters or outputs (e.g., ensuring that
256 soil moisture content or soil thickness cannot be negative).

257 The three terms of the loss function can be weighted differentially, depending on
258 the problem being solved.

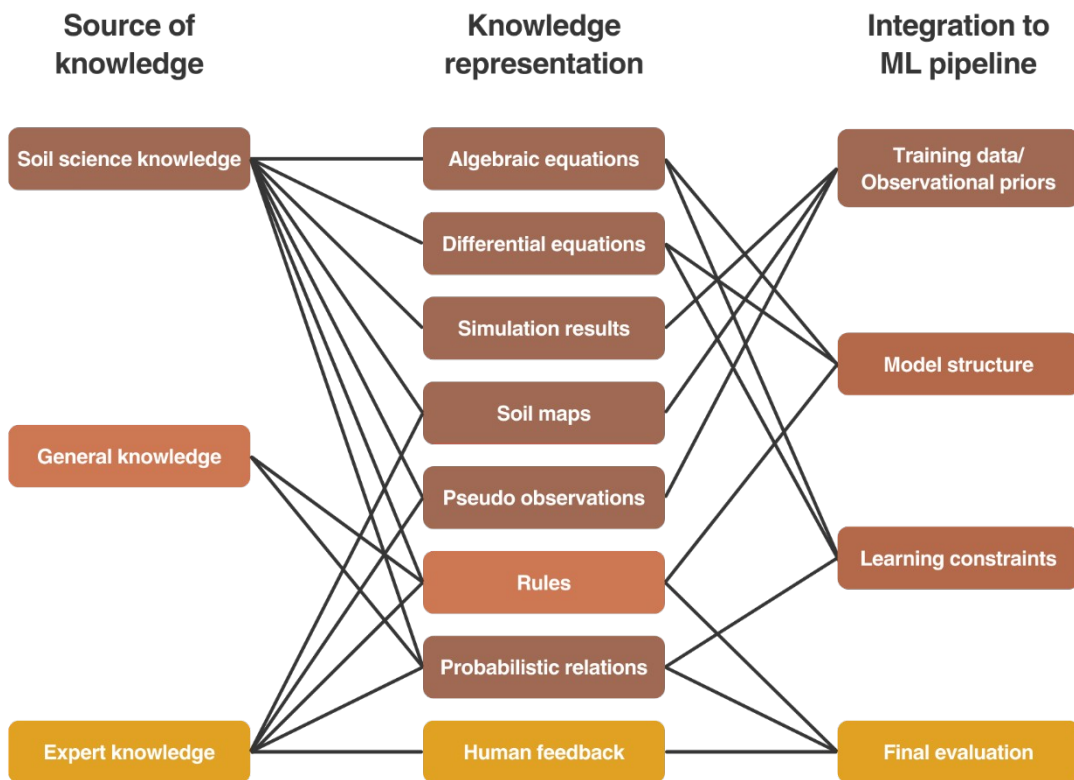
259 The three approaches of knowledge incorporation outlined above are not standalone
260 but could be combined to incorporate prior soil information and model constraints.
261 Finally, soil scientists should evaluate the final outputs of the models in terms of the
262 feasibility of the prediction or maps to evaluate against soil science knowledge or
263 principles. For example, soil scientists could identify the congruency of soil-
264 landscapes maps created by DSM or select digital maps of soil classes and

265 properties for implementing land suitability rules (Bui et al., 2020; Holmes et al.,
266 2021).

267 Model structure design and learning guidance are typically applied together by
268 modifying the ML input-output architecture and loss functions to be minimised. This
269 approach requires a flexible ML framework that allows the structure and model loss
270 function to be customised. ML models with a fixed structure, such as tree models,
271 e.g. Cubist or random forest, may not be well-suited for such applications.
272 Nevertheless, efforts could be made to modify the algorithms such as the spatial
273 random forest model by Talebi et al. (2022).

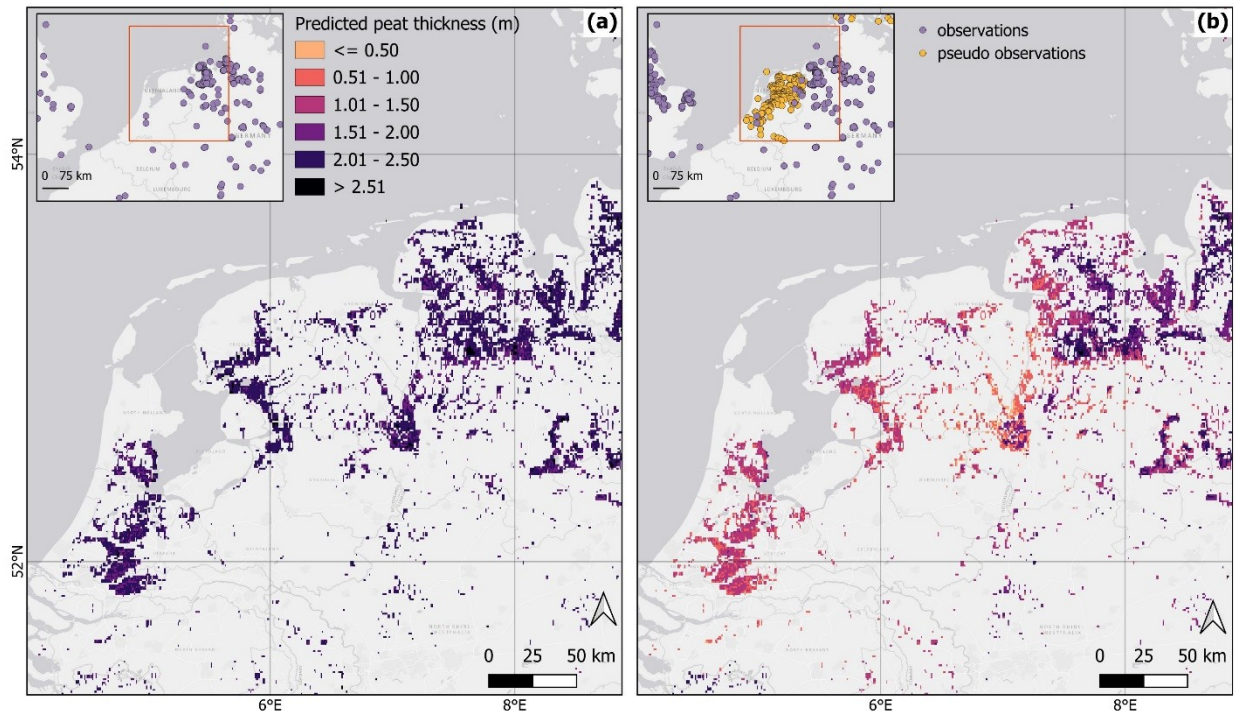
274 A flexible ML framework that can accommodate these requirements includes neural
275 networks with a generic input layer, one or several hidden layers and an output
276 layer. The layers consist of multiple units connected via weights, allowing the model
277 to learn a variety of functions. The structure of inputs and outputs can be modified
278 to fit different dimensions of soil prediction, such as a 1-D, 2-D or 3-D. Additionally,
279 convolutional layers could be added for filtering purposes, and custom objective or
280 loss functions could be defined to align with specific goals. In the next section, we
281 will explore examples of these models in greater detail.

282



283

284 Figure 2. Soil Science-Informed ML, pathways to supervise ML models with soil
 285 science knowledge (adapted from von Rueden et al. (2023)).



286

287 Figure 3. A comparison of peat thickness prediction in European peatlands (a)
 288 without and (b) with observational priors, the addition of data points from national
 289 peat thickness maps (adapted from Widyastuti et al. (2024)).

290

291 **3. Applications of Soil Science-Informed Machine Learning (SoilML)**

292 In this section, we introduce examples representing the application of SoilML
 293 through various forms of knowledge and their incorporation in ML models in several
 294 soil science domains, including digital soil mapping, soil spectroscopy, pedotransfer
 295 functions, and modelling dynamic soil properties. All these examples address soil
 296 security in terms of biomass production, carbon sequestration, and water cycling.

297

298 **3.1 Digital Soil Mapping**

299 Digital soil mapping (DSM) is a process of creating soil maps using spatial covariates
 300 that are combined with field observations, expressed as the "scorpan" model
 301 (McBratney et al., 2003)

302
$$S = f(s, c, o, r, p, a, n), \quad (2)$$

303 where S represents soil classes or attributes. This model provides empirical
304 quantitative descriptions of relationships between soil and other spatially-
305 referenced factors: soil (s), climate (c), organisms (o), topography (r), parent
306 material (p), age (a), and spatial position (n).

307 In DSM, the procedures involved collecting geo-referenced soil observations that are
308 intersected with environmental (scorpan) covariates. A spatial soil prediction
309 function is built to relate observed soil properties of interest to these environmental
310 covariables using ML models. The calibrated spatial soil prediction function can then
311 predict and map soil properties across the area (Arrouays et al., 2020).

312 As discussed in Wadoux et al. (2020), examples of incorporating soil science
313 knowledge in DSM procedures include experts selection of scorpan covariates which
314 conform to the soil-forming processes of the region to be mapped. Pseudo
315 observations could also be added in areas that lack field observations. In the global
316 peat thickness mapping study by Widyastuti et al. (2024), many regions of the world
317 lacked direct field observations. Incorporating pseudo-observations derived from
318 national peat thickness maps can help guide the model. Figure 3 shows an example
319 of peat thickness prediction using a random forest model. Initially, the random
320 forest model was trained only with available field observations, which resulted in
321 the peat thickness values being overpredicted by 2-3 m over the Netherlands and
322 Germany. Incorporating 500 points from peat maps of Sweden, the Netherlands, and
323 Denmark reduced the mean predicted thickness by over half (mean = 1.08 m),
324 resulting in a more accurate and realistic map (Figure 3).

325

326 *3.1.1 Case study: Contextual information for soil mapping using convolutional* 327 *neural networks*

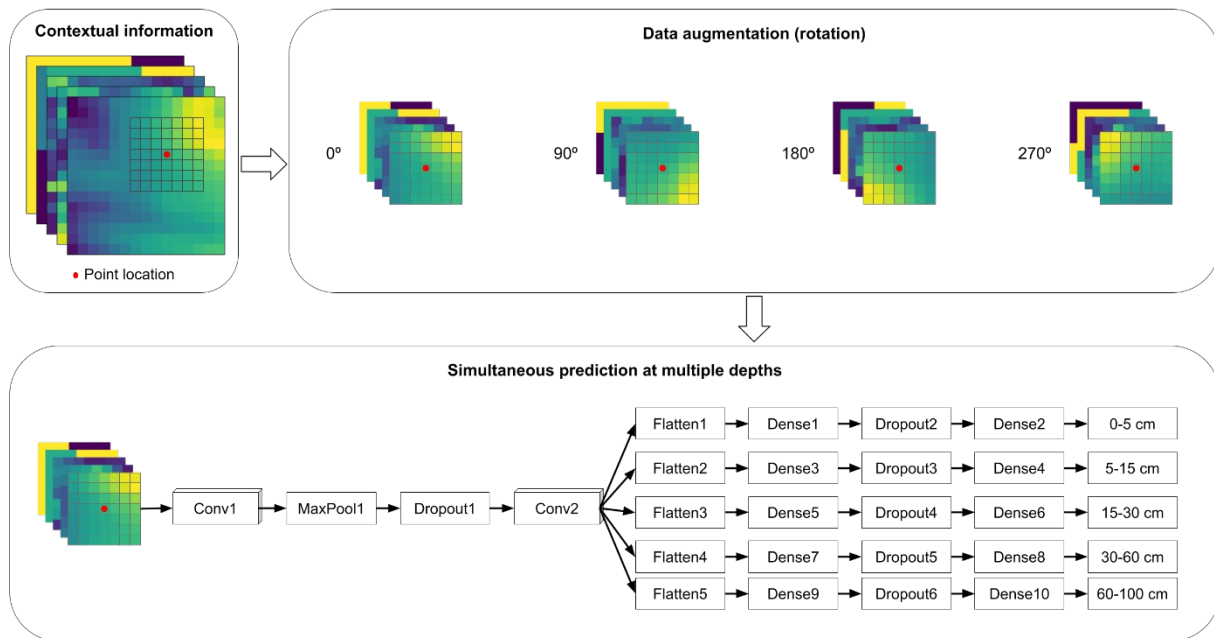
328 *Conventional approach:* DSM models typically use point observations intersected
329 with pixel-wise spatial covariates for calibration. Ideally, contextual information
330 around the observations should be included as covariates. Studies include relative
331 elevation around a point to provide contextual information. For example, Behrens et
332 al. (2010) used differences in elevation from observation points to each of the
333 surrounding neighbourhoods as predictors to capture the relative position of the

334 observation point on the landscape. Other approaches include calculating terrain
335 attributes at various neighbourhood window sizes (Miller et al., 2015).

336 *SoilML*: Padarian et al. (2019b) demonstrated that a convolutional neural network
337 (CNN) model using images of covariates (terrain and climate variables) can
338 effectively explore spatial relationships between a point observation and its
339 neighbouring pixels (Figure 4). The model also includes a 3-D stack of images as
340 input, data augmentation to reduce overfitting, and simultaneous prediction of
341 multiple depths.

342 Using a soil mapping example in Chile, the CNN model was trained to
343 simultaneously predict SOC at multiple depths across the country. To increase data
344 representation, data augmentation was employed to generate new samples by
345 modifying the original data without changing its meaning. This included rotating a
346 3-D array by 90, 180, and 270 degrees. This step also acted as regularisation,
347 reducing model variance and overfitting, and induced rotation invariance by
348 ensuring the model responds similarly to rotated data, such as a soil profile next to
349 a gully. The results showed that the CNN model reduced the error by 30% compared
350 to conventional techniques that only used point information of covariates. For
351 country-wide mapping at a 100 m resolution, a neighbourhood size of 3 to 9 pixels
352 proved more effective than using a single point or larger neighbourhood sizes.
353 Additionally, the CNN model produced less prediction uncertainty and more
354 accurately predicted soil carbon at deeper layers.

355 *Implications*: The CNN framework is designed to accept images as input, capturing
356 information about the observation and its spatial context. Its convolutional layers
357 apply various filters, in the case of a DEM, it effectively mimics the calculation of
358 terrain attributes across different window sizes (Taghizadeh-Mehrjardi et al., 2020).
359 This contrasts with other ML models that require algorithm modifications to handle
360 spatial data. For example, Talebi et al. (2022) developed a spatial random forest
361 model that uses local spatial covariates, which were transformed into vectorised
362 spatial patterns, as predictors. In addition, regularisation, or the addition of a
363 penalty function to the loss function (Eq. 1), could constrain the model to follow
364 certain soil-landscape rules, e.g. a penalty could be added to the loss function when
365 soil thickness on the top of the hill is predicted to be larger than on the lower slope.



366

367 Figure 4. A deep learning framework for digital soil mapping incorporating
 368 contextual information and data augmentation for training a CNN model to predict
 369 soil properties at multiple depths.

370

371 3.1.2 Case study: 3-D soil mapping

372 *Conventional approaches:* The topic of mapping soil properties across space and
 373 depth has gained wide interest. However, soil profiles are usually observed via
 374 horizons, which vary in thickness and depth. In DSM, the variation of soil properties
 375 down a profile is often harmonised using the equal-area spline depth function
 376 approach. Soil observations at various depth intervals are first harmonised to pre-
 377 determined depth intervals. To create maps of soil at these defined depth intervals,
 378 models are trained to predict soil properties at several depth intervals
 379 simultaneously using either neural networks or other ML models capable of
 380 multivariate outputs.

381 Other studies propose that soil properties at any depth can be mapped using a
 382 model that incorporates depth along with spatial covariates as predictor variables,
 383 creating a '3D' model. However, ML models consider depth as one of the covariates,

384 indifferent to spatial covariates. Due to the limited depth inputs, tree models such
385 as random forest are sensitive to the training data and tend to predict soil
386 properties at depths as stepped values (Ma et al., 2021).

387 *SoilML*: We propose designing a neural network to predict soil properties at regular
388 depth intervals to alleviate such a problem. The model would take spatial covariates
389 as inputs, and the training of the model disaggregates and predicts soil properties
390 at all depths simultaneously.

391 For example, 59 soil cores, varying in depth between 85 and 130 cm were used in a
392 study by (Fajardo et al., 2016). The cores had SOC measurements via visible-near-
393 infrared (vis-NIR) and shortwave infrared (SWIR) spectroscopy (wavelengths
394 between 350 and 2500 nm) at every 2 cm down to a depth of 1 m. To simulate
395 horizon sampling, SOC observations were grouped by soil horizons (S). The spatial
396 prediction model used the following covariates as inputs: terrain (elevation, wetness
397 index, mid-slope position, altitude above channel network), remote sensing images
398 (Vis-NIR and SWIR bands), and predicted SOC every 2 cm from the surface to 1 m.
399 Simulating soil observations by layers, the training data of SOC observations were
400 grouped by soil horizons. Thus, the loss function for the model is:

$$401 \quad L = \sum_{i=1}^n \left(\sum_{j=1}^{m(i)} (S_{ij} - \hat{s}_{ij})^2 \right) \quad (3)$$

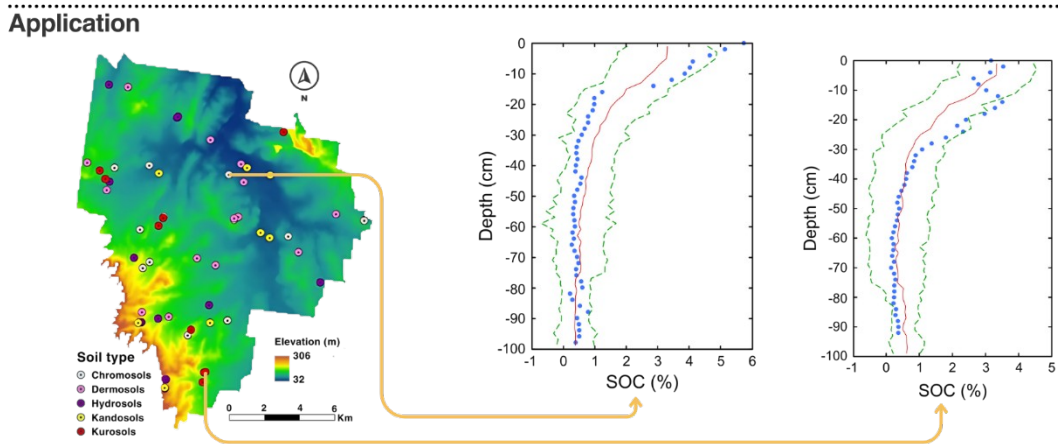
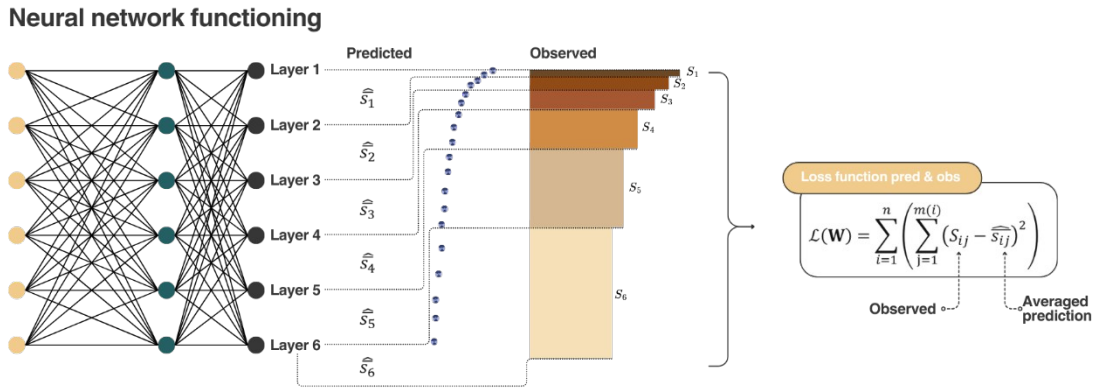
402 where n is the number of soil cores, m is the number of layer observations per core,
403 and S is the observed SOC value per layer of observations. Note that the model
404 predicts SOC at specific points and \hat{s}_{ij} refers to the aggregated or averaged
405 predicted value corresponding to each observed layer. Figure 5 shows an example of
406 the predicted SOC values across the profile.

407 *Implications and prospects*: Although numerous studies have incorporated depth as
408 a covariate to generate 3D maps, it is important to be cautious about combining
409 spatial covariates (covering geographical areas with grid spacing ranging from
410 approximately 1 to 1000 m) with depth, which varies from about 0.01 to 2 m. ML
411 models may struggle to distinguish the significant differences in scale and
412 continuity between these types of measurements. Formal geostatistical approaches
413 which predict in 3D by disaggregating the bulk depth measurements (Orton et al.,

414 2020), or using the Gaussian process regression (Wang et al., 2024), provide more
415 robust solutions.

416 While soil properties at various depth intervals have been extensively mapped over
417 the years, there remains a gap in mapping the distribution patterns of soil horizons
418 and representing soil as a 3-D continuum. Soil is fundamentally a three-dimensional
419 body composed of distinct horizons. Various studies have employed techniques such
420 as electromagnetic induction and interpolation based on cone penetrometer
421 resistance to map soil layers (Grunwald et al., 2001) and the thickness of soil
422 horizons (Chaplot et al., 2010). For instance, Mendonça Santos et al. (2000) mapped
423 the thickness of each of the 12 horizons in a Swiss floodplain in two dimensions,
424 then stacked these results to represent a three-dimensional volume. Similarly,
425 Gastaldi et al. (2012) combined logistic regression with ordinary regression to first
426 model the occurrence of each horizon and subsequently their thickness.
427 Advancements in ML, particularly neural networks, now offer the potential to model
428 soil as a profile of horizons and predict each horizon's thickness and composition as
429 observed.

430



431

432 Figure 5. An example of a neural network model that predicts point soil observations
 433 along the soil profile depth using environmental covariates. The neural networks
 434 model was trained using soil profile data, with the loss function minimising the
 435 difference between the averaged predicted layer values and the observed soil layer
 436 values. The figure below shows the prediction of SOC across the Hunter Valley in
 437 NSW, Australia. The observed (blue dots), the predicted (red line), and the
 438 prediction interval (green dashed-line).

439

440 3.1.3 Case study: Soil class mapping incorporating taxonomic distance

441 *Conventional approaches:* Digital soil class mapping typically begins with the
 442 description of soil profiles and allocating the profiles to soil classes according to an
 443 established soil classification system. This process continues with correlating the
 444 observed soil classes with co-located covariates at each observation site. Most ML
 445 training in supervised classification involves minimising classification errors:

446
$$L = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c I(y_i \neq k) \quad (4)$$

447 where $i = 1, 2, \dots, n$ is the number of observations, and $k = 1, 2, \dots, c$ is the number
 448 of classes, $I(y_i \neq k)$ is an indicator when observed class y_i is not equal to class k . This
 449 error criterion assumes that the errors across all classes are of equal importance.
 450 However, this is not valid for soil classes and does not allow for situations where
 451 some errors are more important than others.

452 *SoilML*: Taxonomic distance between soil classes can be incorporated into a
 453 supervised classification routine. Minasny and McBratney (2007) calculated the
 454 taxonomic distance between soil classes based on a central concept, e.g. to define
 455 a modal soil profile for each soil class. The taxonomic distance matrix between soil
 456 classes can be represented as \mathbf{D} , with $D_{j,k}$, represent the distance between class j
 457 and class k . The supervised classification loss function could be defined as the
 458 average misclassification cost:

459
$$L_a = \frac{1}{n} \sum_{i=1}^n D(C_i, \hat{C}_i) \quad (5)$$

460 where L_a is the average taxonomic distance error, $D(C, \hat{C})$ is the taxonomic distance
 461 between observed class C and predicted class \hat{C} . By using classification trees that
 462 minimise the taxonomic distance over misclassification error, the methodology is
 463 refined to model soil class relationships.

464 *Implications*: Defining taxonomic distance extends beyond predicting soil classes to
 465 encompass the development of soil classification units. While early numerical soil
 466 classification methods in the 1950s were constrained by limited data and
 467 technology, modern advances now allow taxonomic distance calculations to explore
 468 correlations between national systems and global frameworks like the World
 469 Reference Base or USDA Soil Taxonomy. For example, Michéli et al. (2016) used
 470 taxonomic distance to differentiate USDA soil great groups, demonstrating its utility
 471 in objectively refining classification criteria. Similarly, Hughes et al. (2017) showed
 472 that taxonomic distance calculations can aid in translating soil classes across
 473 various classification systems, enhancing global comparability and consistency.

474 Laborczi et al. (2019) compared topsoil (0-30 cm) texture classes in Hungary using
475 two methods: directly compiled maps of clay, silt, and sand content for the 0-30 cm
476 depth, and synthesised maps derived from the thickness-weighted average of the
477 0-5, 5-15, 15-30 cm layers. While the soil texture class maps produced by both
478 methods are similar, taxonomic distances between the two maps reveal more
479 pronounced discrepancies in certain regions. Significant differences are observed
480 particularly in hilly and mountainous areas, which could pose challenges in erosion
481 and sedimentation modelling and the prediction of flash floods. Additionally,
482 inaccuracies in mapping salt-affected and hydromorphic soils could impact water
483 management and irrigation planning. Nevertheless, DSM of soil classes still rarely
484 considered the taxonomic distance in the ML training workflow.

485

486 **3.2 Soil Spectroscopy**

487 Soil can reflect, scatter, or emit electromagnetic radiation, resulting in a unique
488 spectral signature. Soil responds uniquely to infrared radiation, making infrared
489 spectrometers suitable for soil analysis because they can measure rapidly, cost-
490 effectively, and non-destructively. An infrared spectrometer can predict multiple soil
491 properties from a single-spectrum measurement. However, soil is a complex mixture
492 of mineral and organic constituents, it is challenging to assign specific spectral
493 features to particular physical, chemical, or biological components. Therefore,
494 empirical multivariate calibration techniques are commonly employed to predict soil
495 properties by relating spectra data to observed soil characteristics (Chen et al.,
496 2023b; Hutengs et al., 2021; Vohland et al., 2022).

497 *3.2.1 Case study: Physical model for soil spectra response to moisture and hydraulic* 498 *properties*

499 One major factor affecting soil reflectance is the presence of water (Lobell and
500 Asner, 2002). Wet soil typically reflects less light than dry soil. This sensitivity of soil
501 reflectance to moisture allows for the rapid estimation of soil water content through
502 vis-NIR and SWIR reflectance measurements (Liu et al., 2002).

503 *Conventional approaches:* Various empirical models have been developed to relate
504 soil reflectance to soil water content in the Vis-NIR-SWIR spectra (Babaeian et al.,

505 2019). These models include partial least squares regression (Bogrekci and S. Lee,
506 2006; Castaldi et al., 2015), principal component regression (Chang et al., 2001),
507 and ML models (Hassan-Esfahani et al., 2015; Zaman et al., 2012). While these
508 models are effective, they require extensive databases for calibration and their
509 applicability is restricted to the specific soil conditions under which they were
510 developed, as moisture response to NIR radiation depends on soil types and
511 constituents (Babaeian et al., 2019).

512 *SoilML*: Radiative transfer models can effectively describe diffuse infrared radiation
513 in soil. The Kubelka and Munk (KM) model (Kubelka and Munk, 1931) is a two-flux
514 radiative transfer model that describes light transfer through a particulate medium,
515 characterised by absorption (κ) and scattering (s) coefficients. The model uses a set
516 of differential equations to account for light travelling in two opposing directions and
517 yields reflectance and transmittance as a function of k , s , and depth. The optical
518 depth is assumed to be infinite for soil, and therefore the transmission becomes
519 negligible.

520 Sadeghi et al. (2015) applied the KM model to explore the relationship between soil
521 water content and reflectance. They proposed that the optical properties (k and s)
522 of soil can be expressed by a linear volume averaging of the optical properties of its
523 constituents, i.e., solid particles, water, and air. Based on this approach, they
524 derived a physically based and linear equation that explicitly expresses SWIR to
525 water content:

$$\frac{\theta}{\theta_s} = \frac{r - r_d}{r_s - r_d} \quad (6)$$

526 where θ is the volumetric water content ($\text{m}^3 \text{m}^{-3}$), θ_s is the saturated water content
527 ($\text{m}^3 \text{m}^{-3}$), and r is the transformed reflectance. The parameters r_d , and r_s are the
528 transformed reflectance of soil in dry and saturated states, respectively.
529 Transformed reflectance (r) can be calculated from the measured reflectance (R) as
530 follows:

$$r = \frac{(1 - R)^2}{2R} \quad (7)$$

531 Norouzi et al. (2022) hypothesised that the two distinct forms of soil water, i.e.,
532 capillary and adsorbed water, impact soil reflectance differently (Figure 6). Building

533 on this hypothesis, they considered different optical properties for capillary and
 534 adsorbed water and derived a new model to describe the relationship between soil
 535 reflectance and water content:

$$r = r_d + \overbrace{c_a \theta_a^{p_a}}^{r_a} + \overbrace{c_c \theta_c^{p_c}}^{r_c} \quad (8)$$

536 The total transformed reflectance of wet soil (r) can be decomposed into three
 537 components: r_d , r_a , and r_c corresponding to the dry soil, adsorptive water, and
 538 capillary water, respectively. Parameters c_a , p_a , c_c , p_c are the optical properties
 539 related to adsorbed and capillary water. In this equation, θ_a and θ_c are the
 540 volumetric water contents of adsorbed and capillary water that can be derived from
 541 the soil retention curve. Norouzi et al. (2022) used the model by Lebeau and Konrad
 542 (2010) for soil water retention curve to partition the total water content (θ) into its
 543 components θ_a and θ_c :

$$\theta = \theta_c + \theta_a \quad (9)$$

544 where the capillary component is modelled based on Kosugi (1996):

$$\theta_c = \frac{1}{2} \theta_s \operatorname{erfc} \left[\frac{\ln(h/h_m)}{\sqrt{2} \sigma} \right] \quad (10)$$

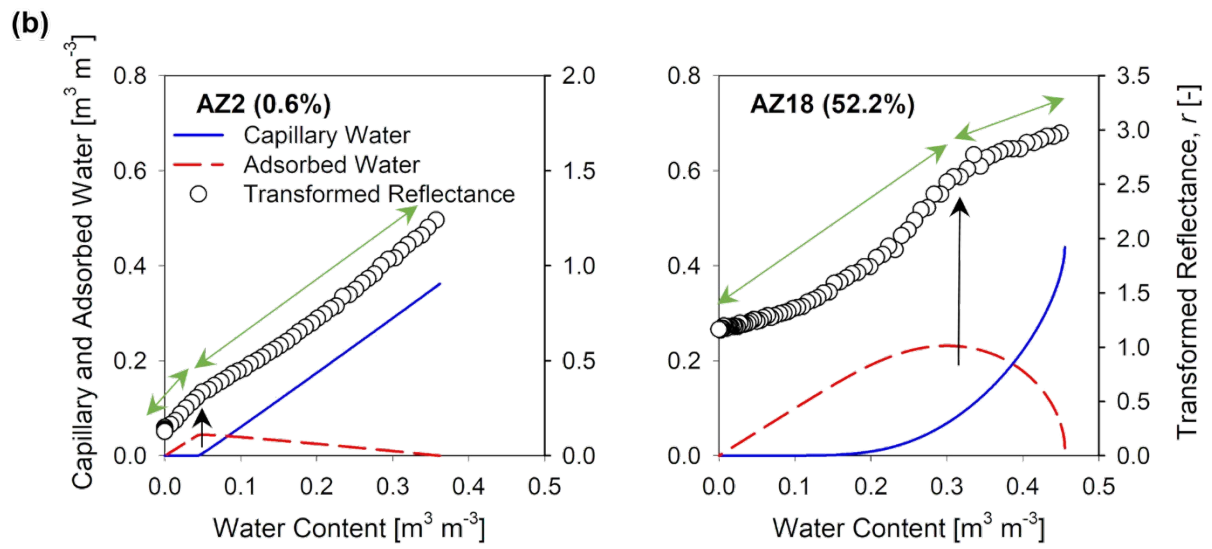
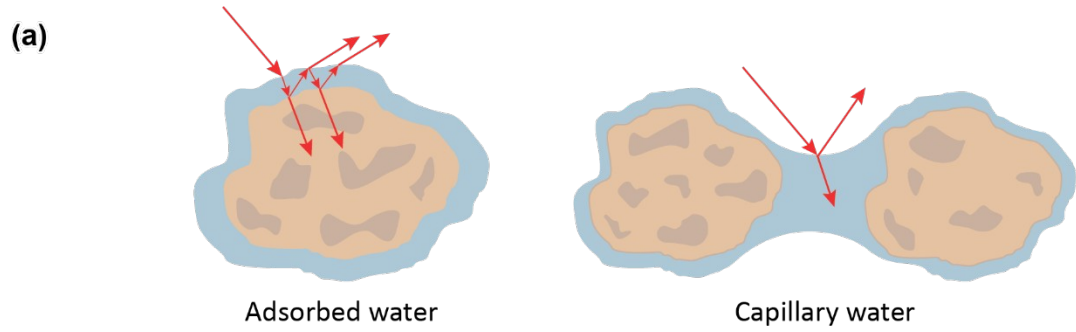
545 where h is the pressure head, θ_s is the saturated volumetric water content, and erfc
 546 denotes the complementary error function; h_m , σ , and θ_o are fitting parameters.

547 The adsorptive component is represented using the Campbell and Shiozawa (1992)
 548 model for extremely low matric potentials, which linearly diminishes as the amount
 549 of capillary water increases:

$$\theta_a = \theta_o \left(1 - \frac{\ln|h|}{\ln|h_d|} \right) \left(1 - \frac{\theta_c}{\theta_s} \right) \quad (11)$$

550 where h_d is the pressure head at oven dryness and generally corresponds to a finite
 551 value of -10^7 cm (Campbell and Shiozawa, 1992). Equation (8), in combination with
 552 Equations (10) and (11), directly connects soil reflectance to the soil water retention
 553 curve.

554



555

556 Figure 6. (a) An illustration of reflections and refractions of light beams (indicated
557 by red arrows) as they interact with adsorbed and capillary water. (b) Effects of
558 capillary and adsorbed water (left axis) on transformed reflectance measurements
559 at 2210 nm (right axis) for two Arizona soils with 0.6% and 52.2% clay content. The
560 black arrows indicate the points where the reflectance slope (shown by green
561 arrows) changes sharply, marking the transition from capillary to adsorbed water
562 regimes.

563

564 Norouzi et al. (2022) demonstrated that soil reflectance is influenced not only by the
565 amount of water in the soil but also by the structure of water, specifically the
566 capillary and adsorbed components. They showed that when the soil water
567 retention curve is known, Equation (8) accurately describes the relationship
568 between soil moisture and reflectance. As shown in Figure 6b, a noticeable change
569 in reflectance at 2100 nm occurred at a specific point marked by a black arrow. The
570 slope of the reflectance, marked by a green arrow, changes before and after this
571 point, indicating that the shift corresponds to the maximum water content for
572 adsorbed water, as seen when compared to the water components on the left axis.
573 This transition signifies where capillary water recedes, and adsorbed water becomes
574 the dominant component at the surface. This reflectance change is highly
575 dependent on soil texture, occurring at lower water content for coarse-textured soils
576 (e.g., at $\sim 0.05 \text{ m}^3 \text{ m}^{-3}$ for AZ2) and at a higher water content for fine-textured soils
577 (e.g., at $\sim 0.3 \text{ m}^3 \text{ m}^{-3}$ for AZ18).

578 This also means that Equation (8) can be inverted to derive the soil water retention
579 curve using NIR spectra of soil reflectance and moisture content measured during
580 an evaporation experiment, optimising retention curve parameters and optical
581 properties to match observed NIR reflectance. Validation with 21 soils of varying
582 textures and mineralogy demonstrated accurate retrieval of the entire retention
583 curve, from saturation to oven dryness (Norouzi et al., 2023).

584 *Implications and prospect:* Considering that the SWIR of a drying thin soil sample
585 from saturation to air-dry can be measured within a few hours, the physics-based
586 approach proposed by Norouzi et al. (2023) can be an efficient method for

587 measuring the soil water retention curve, which often takes several weeks.
588 Although the above example is not an ML model, the water retention function could
589 be defined as a neural network. This model can be further constrained by
590 Richardson-Richards' equation to align with water content and time measurements
591 collected during the evaporation experiment. Moreover, the radiation transfer model
592 can integrate factors such as soil water content, particle size, and organic matter
593 effect. It has the potential to predict soil texture and organic matter content by
594 calibrating optical, absorption, and scattering coefficients (Wu et al., 2023). In
595 combination with physics-informed ML, these models can improve both the
596 predictability and interpretability of soil spectra models.

597

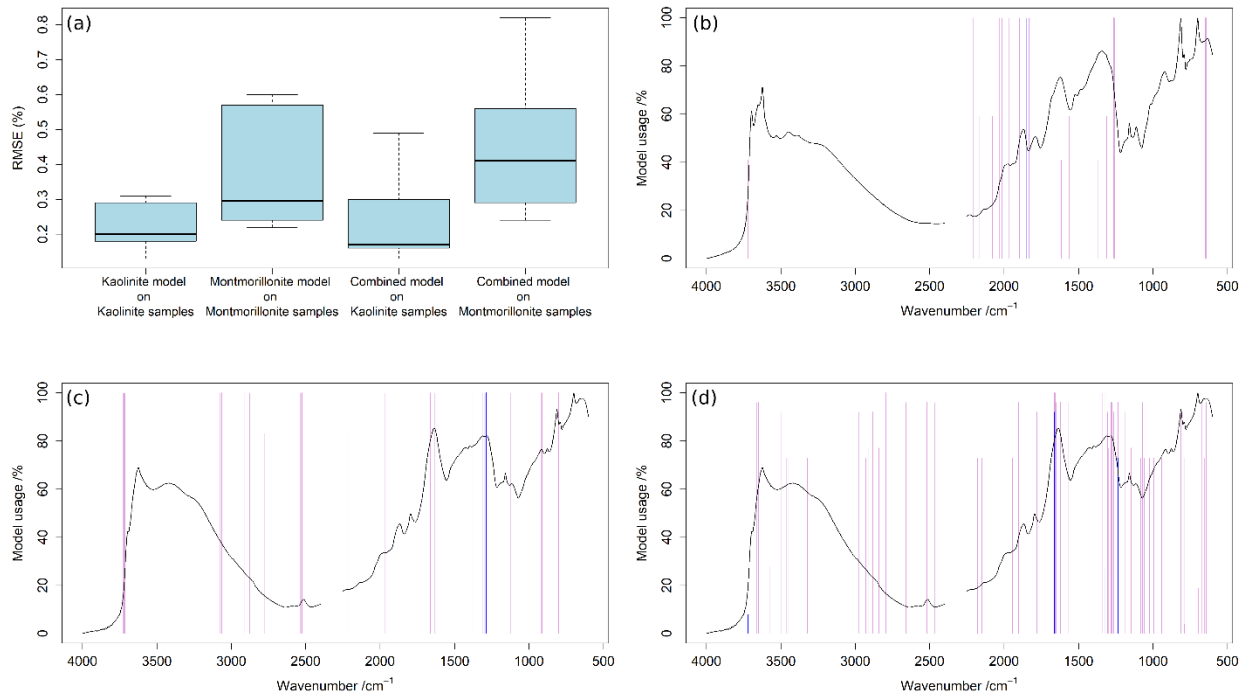
598 *3.2.2 Case study: Building soil-based spectra functions*

599 *Conventional approaches:* Soil spectra are typically pre-processed with smoothing or
600 transformation to remove noise and serve as inputs of regression models or ML
601 algorithms. The models are trained to minimise root mean square error (RMSE) and
602 maximize the coefficient of determination (R^2). The working of the models can be
603 explained by the importance or usage of variables in the model. For example,
604 variable importance in projection (VIP) score is used in partial least square
605 regression modelling to help identify which wavelength is mostly related to the soil
606 property, and the model usage rate can be used to evaluate Cubist models (Chen et
607 al., 2023a; Seidel et al., 2022). However, there is no information input from soil
608 science knowledge when training the model, and the prediction result will merely
609 depend on the relationship between the spectra features and the soil properties. In
610 this way, soil science knowledge only serves the purpose of explaining outcomes,
611 rather than being directly involved in model building.

612 *SoilML:* Prior soil information, e.g. morphological and mineralogical characteristics,
613 can help divide the samples into homogeneous groups before modelling, and
614 models will therefore be trained based on soils with shared properties. By
615 comparing the effects of models trained on (1) all samples and (2) sample sets
616 divided by prior information, this case study demonstrated the possibility of
617 including soil knowledge in the modelling process.

618 In this case study, 370 Bt horizon soil samples with 0–5% carbon content were
619 extracted from the Kellogg Soil Survey Laboratory (KSSL) dataset (Soil Survey Staff,
620 2014). X-ray diffraction analysis revealed kaolinite and montmorillonite as the
621 dominant clay minerals. Modelling of total carbon content was conducted separately
622 on 185 kaolinite-dominant samples and 185 montmorillonite-dominant samples with
623 mid-infrared spectra and the Cubist regression tree model. Spectra were pre-
624 processed with Savitzky-Golay smoothing, SNV transformation, and trimming off the
625 CO₂ peak. Samples were randomly divided into a 70% calibration set and a 30%
626 validation set for the training and testing of Cubist models, and this process was
627 performed 10 times to get a distribution of results.

628 Results show that individual models based on dominant mineralogical components
629 were more accurate than the total model (Figure 7a). The model created using all
630 samples tended to have higher spread in the boxplot, which indicated less
631 robustness than models from the pre-divided training set. The kaolinite model
632 mainly used wavenumbers around 2000 cm⁻¹ (Figure 7b), and the montmorillonite
633 model more relied on multiple wavenumbers across the spectrum (Figure 7c). The
634 combined model, on the other hand, utilised more conditions and variables than the
635 individual models (Figure 7d), which might be due to the heterogeneity of soils
636 dominated by different mineralogical characteristics. By including prior soil
637 information in modelling, grouping the samples based on their mineralogical
638 component improved the performance of models and enabled clearer differentiation
639 of wavelengths used in the models.



640

641 Figure 7. (a) Boxplots of root mean square error (RMSE) results from ten repetitions
 642 of Cubist modelling using various input data. (b)-(d) Variable importance of Cubist
 643 models to predict total carbon content with mid-infrared spectra using: (b) kaolinite-
 644 dominant samples for calibration, (c) montmorillonite-dominant samples for
 645 calibration, and (c) combined kaolinite-dominant and montmorillonite-dominant
 646 samples for calibration. Black lines are the mean spectra of calibration samples,
 647 purple vertical lines are the variables used as predictors, while the blue vertical
 648 lines are the conditions of Cubist models.

649

650 *Implications:* Analysing soil spectra alongside soil science knowledge enables the
 651 identification of specific soil components, enhancing the effectiveness of statistical
 652 methods and improving the understanding of soil properties and processes. By
 653 grouping soils based on pedological information, such as soil order or soil horizon or
 654 mineralogical component, researchers can refine models to achieve more accurate
 655 and interpretable predictions. This approach encourages soil scientists to look
 656 beyond mere prediction accuracy and develop a deeper understanding of the soil.
 657 Incorporating soil knowledge can involve pre- or post-machine learning calibration,
 658 such as inspecting spectra or grouping soils by mineralogy to guide the models.
 659 After modelling, verifying that predictions align with soil science principles is crucial,

660 ensuring that ML applications do not overshadow the fundamental soil
661 understanding (Ma et al., 2023).

662

663 **3.3 Pedotransfer functions**

664 Pedotransfer functions (PTFs) translate basic soil data into more complex, labour-
665 intensive, and costly soil properties (Weber et al., 2024). They serve as predictive
666 tools for estimating certain soil properties from easily measured or available data,
667 thereby bridging the gap between available and required data. A prominent
668 application of PTFs is in predicting the soil water retention curve, which describes
669 the soil water content (θ), i.e., the volume of water per volume of soil under
670 equilibrium at a given pressure head (h). Since measuring a soil water retention
671 curve is time-consuming, PTFs offer a practical alternative by estimating it based on
672 soil physical properties such as texture, bulk density, and SOC (Bagnall et al., 2022;
673 Weber et al., 2024). PTFs are also used in various applications, including assessing
674 irrigation, drainage, and evapotranspiration, enhancing DSM, and providing inputs
675 for process-based simulation models to evaluate soil functions.

676 *3.3.1 Case study: Prediction of water retention and hydraulic conductivity curves*

677 *Conventional approaches:* The structure of a PTF typically involves using ML models
678 to relate predictors (input data, such as soil texture and bulk density) to a
679 predictand (output, such as water content at field capacity). In the context of
680 predicting water retention curves using neural networks, there are three main
681 model configurations (Figure 8):

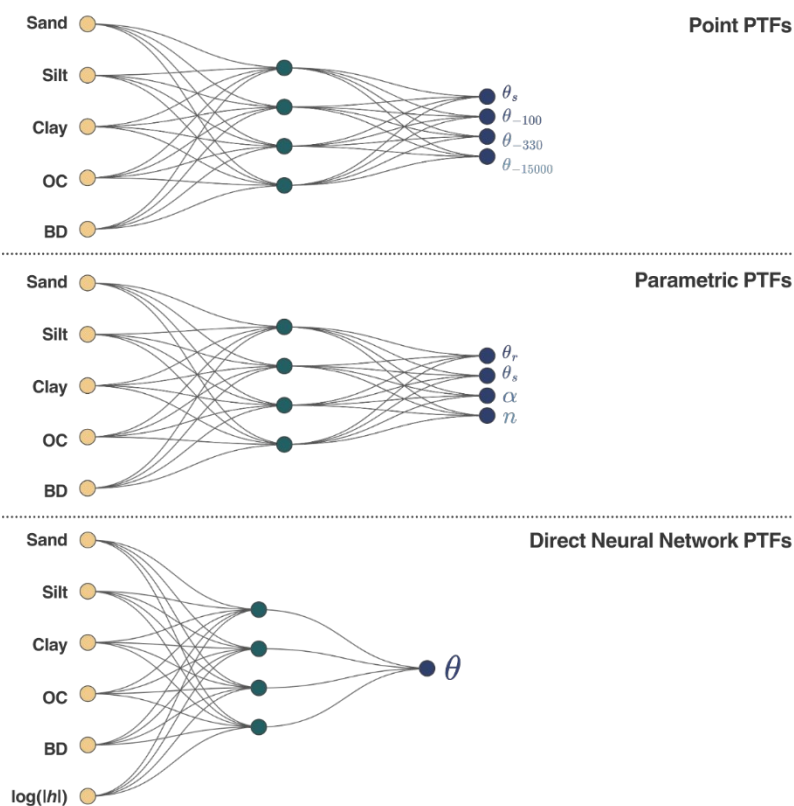
682 Point PTFs: a point PTF predicts water contents (θ) at specific pressure heads (h)
683 from basic soil properties such as sand, silt, clay, SOC, and bulk density. They
684 require a training dataset that includes measurements of water retention at the
685 specified pressure heads along with the basic soil properties.

686 Parametric PTFs: This configuration uses a hydraulic model capable of representing
687 the data, focusing on predicting parameters of the hydraulic model. The output
688 parameters are then used to form a continuous function describing the relationship
689 between the dependent variable (θ) and the independent variable (h). This method

690 is favoured for its ability to provide a continuous prediction curve and is commonly
 691 employed in water retention modelling.

692 Direct Neural Network PTFs: In this setup, neural networks are directly applied to
 693 model water retention. The pressure head, along with basic soil properties, are used
 694 as inputs, allowing the model to learn a non-specific form of the soil water retention
 695 curve.

696



697

698 Figure 8. Three configurations of PTFs predicting soil retention: point, parametric,
 699 and direct neural networks (based on Haghverdi et al. (2012)).

700

701 In parametric PTFs, the van Genuchten equation (van Genuchten, 1980) is
 702 commonly used to model the water retention curve:

703
$$\theta(h) = \theta_r + (\theta_s - \theta_r) S_e(h)$$

704 $S_e(h) = \left(\frac{\theta - \theta_r}{\theta_s - \theta_r} \right)^{\frac{1}{m}} \left(1 + |ah|^n \right)^{-m}$ (12)

705 where the water content (θ) as a function of pressure head (h) is described by four
 706 parameters: θ_r , residual water content; θ_s , saturated water content; a , the inverse of
 707 air-entry pressure; and n , curve shape factor, with m defined as $m = 1 - 1/n$. The van
 708 Genuchten model can be combined with the capillary theory model of Mualem
 709 (1976) to predict the unsaturated hydraulic conductivity curve, known as the
 710 Mualem-van Genuchten model.

711 Creating a parametric PTF first involves fitting the van Genuchten equation to
 712 observations to estimate the parameter vector $\phi = [\theta_r, \theta_s, a, n]$. This is followed by
 713 forming relationships between basic soil properties (sand, clay, bulk density) and
 714 the parameters using neural networks (or other ML models) by minimising the
 715 following loss function:

716 $L = \sum_{i=1}^n \left(\sum_{k=1}^p (\phi_{ik} - \hat{\phi}_{ik})^2 \right)$, (13)

717 where n is the number of observations, and p is the number of parameters to be
 718 estimated. The drawback of this approach is that the predicted parameters here do
 719 not necessarily bear a physical relationship.

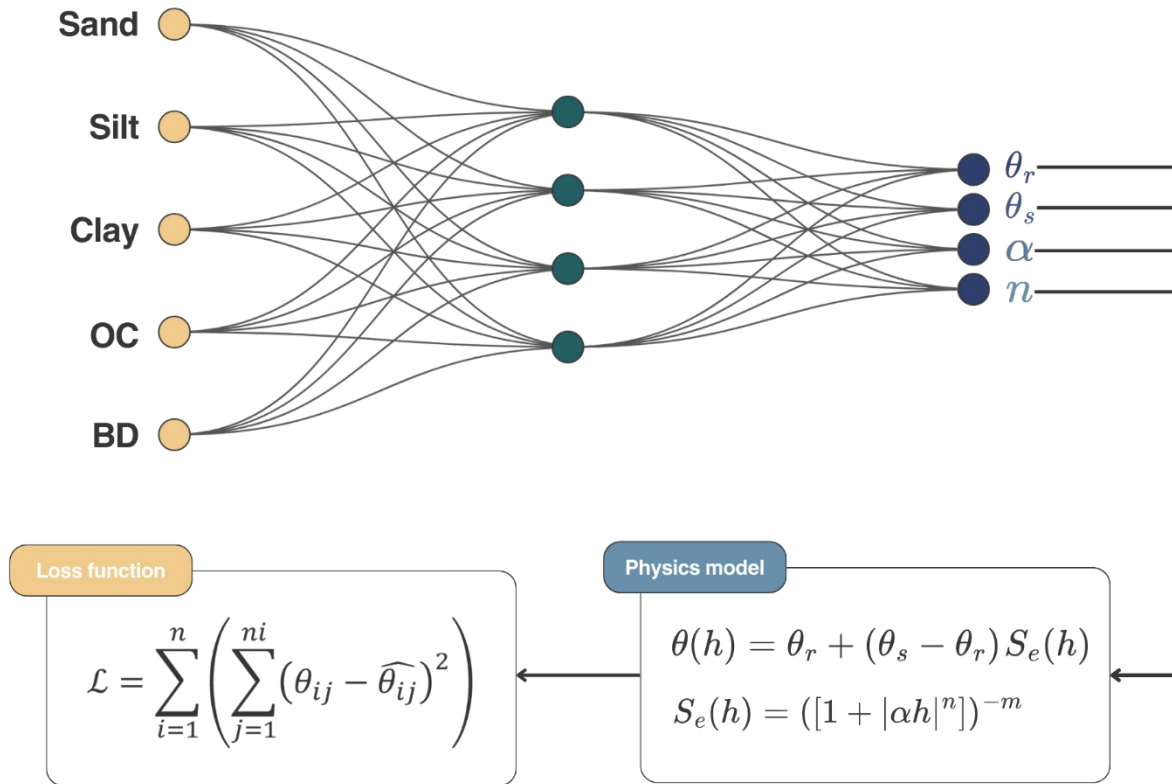
720 *SoilML*: We can incorporate the van Genuchten function in training the ML model by
 721 requiring the estimated parameters to predict the observed water retention $[\theta(h)]$
 722 rather than predicting each parameter of the van Genuchten equation
 723 independently. Minasny and McBratney (2002) used neural networks to predict the
 724 parameters of the van Genuchten function using soil properties (sand, clay, bulk
 725 density). The neural networks model predicted the van Genuchten parameters
 726 $[\hat{\theta}_r, \hat{\theta}_s, \hat{a}, \hat{n}]$ but was trained to minimise the difference between the observed and
 727 predicted water content:

728 $L_\theta = \sum \left(\theta(h) - \hat{\theta}(h \vee \hat{\theta}_r, \hat{\theta}_s, \hat{a}, \hat{n}) \right)^2$ (14)

729 In this case, the ML model is constrained to predict parameters that fit the water
 730 retention data (Figure 9). This led to more realistic prediction values and a more
 731 accurate estimation of the water retention relationship. Using soil water retention

732 data from Australia, Minasny and McBratney (2002) demonstrated that the PTFs
 733 trained using Eq. (14) predicted water retention much better compared to models
 734 that were trained using Eq. (13). In addition, the parameters of the van Genuchten
 735 model were better constrained according to theoretical expected values.

736



737

738 Figure 9. A physics-informed pedotransfer function for predicting a water retention
 739 curve function.

740

741 Recent research, such as that by Peters et al. (2024) and Weber et al. (2020), have
 742 highlighted the shortcomings of the Mualem-van Genuchten model (1976; 1980),
 743 particularly under dry conditions. With the residual water content θ_r as a fitted
 744 parameter, the model implies that the water content would never be lower than
 745 that value. It also focuses on hydraulic conductivity driven by capillarity and fails to

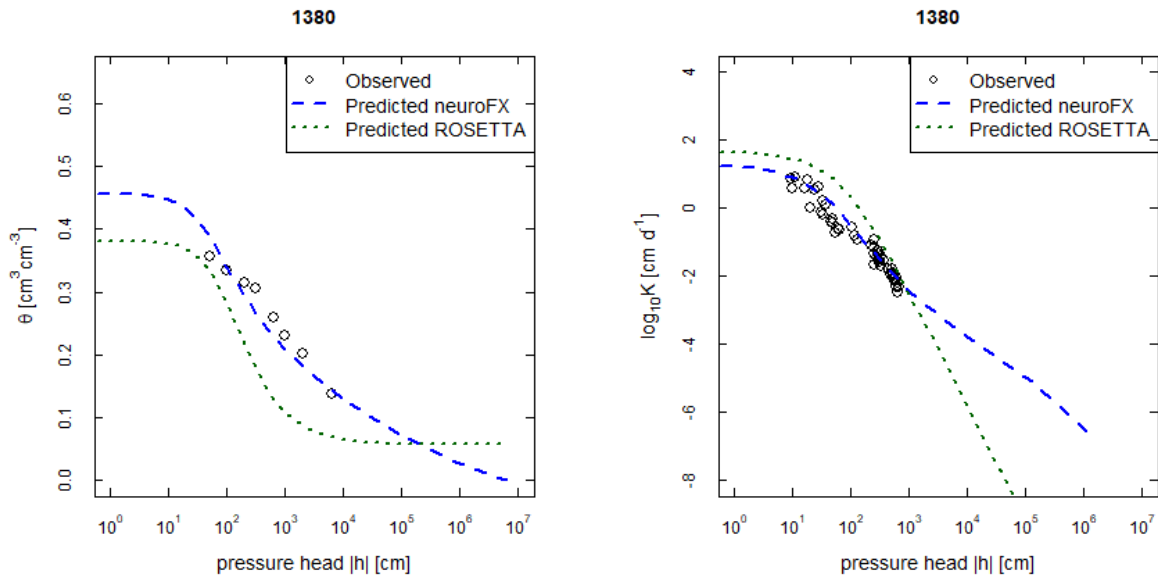
746 provide a reliable description of water retention and conductivity across the
747 complete range of soil water content levels.

748 Ruidiyanto et al. (2021) employed a comprehensive water retention and hydraulic
749 conductivity model, referred to as the FXW model. It is based on the Fredlund and
750 Xing (1994) water retention model and the hydraulic conductivity model of Wang et
751 al. (2018). The FXW model can calculate the retention and hydraulic conductivity
752 curves across the entire range of matric heads, from saturation to complete
753 dryness. The water retention follows a series of functions that aim to scale the
754 water content from saturation θ_s to complete dryness ($\theta=0$) at a defined pressure
755 head h_0 at -6.3×10^6 cm, as shown in Eq. (A1) – (A6).

756 While more complicated, the number of parameters of the water retention curve [
757 θ_s, α, n, m] is the same as the van Genuchten model. Ruidiyanto et al. (2021)
758 developed a PTF called neuroFX that predicts parameters of the water retention
759 curve using a neural network that takes sand, silt, clay, and bulk density as inputs.
760 The loss function of the neural networks is defined in terms of measured versus
761 predicted water content (Eq. A7).

762 Once the water retention parameters PTF was predicted [$\hat{\theta}_s, \hat{\alpha}, \hat{n}, \hat{m}$], the predicted
763 parameters were used to calculate the effective saturation S_{ek} . The FXW parameters
764 for hydraulic conductivity [$\log(K_s), L$] were then estimated using another neural
765 network function trained to minimise the difference in hydraulic conductivity values
766 (Eq. A8).

767 These PTFs were shown to describe water retention and hydraulic conductivity more
768 accurately than conventional PTFs. In sandy to loamy soils, conventional PTFs
769 trained to predict the Mualem-van Genuchten parameters (ROSETTA) show an
770 under-prediction of hydraulic conductivity in the dry range by several orders of
771 magnitude (Zhang and Schaap, 2017). Moreover, ROSETTA produced non-zero water
772 content at the dry end. The neuroFX PTF fits both water retention and hydraulic
773 conductivity data well across the entire range of water contents (Figure 10).



774

775 Figure 10. An example of water retention and hydraulic conductivity curve of a
 776 sandy loam predicted with PTF using neuroFX compared to the conventional Rosetta
 777 model. The Rosetta model uses the Mualem-van Genuchten model, where the water
 778 content does not reach zero as the soil is drying and the conductivity drops rapidly
 779 at dry potentials.

780

781 In the direct neural networks PTFs, the neural networks are trained to model the
 782 water retention function directly (Haghverdi et al., 2012). Since the neural network
 783 learns the shape of the retention curve solely from measurements, the performance
 784 of such PTFs is highly dependent on the quality, density, and distribution of the soil
 785 water retention curve measurements within the training set (Haghverdi et al.,
 786 2014). Norouzi et al. (2024) addressed this issue by imposing physical constraints
 787 on the relationship between the pressure head in the input layer and water content
 788 in the output layer. Specifically, four constraints were imposed: a monotonically
 789 decreasing constraint between $\log(|h|)$ and water content (θ), enforcing linear
 790 behaviour at the dry end of the retention curve, setting a specified range for the
 791 pressure head at zero water content (h_0), and enforcing a constant water content
 792 constraint above air-entry pressure. The loss function used for training the neural
 793 network is given as:

$$J = \frac{\lambda_1}{N_{wet}} \sum_{i=1}^{N_{wet}} [\hat{\theta}^{(i)} - \theta^{(i)}]^2 + \frac{\lambda_2}{N_{dry}} \sum_{i=1}^{N_{dry}} [\hat{\theta}^{(i)} - \theta^{(i)}]^2 + \frac{\lambda_3}{S_1} \sum_{i=1}^{S_1} \dots \dots \dots \dots \dots \quad (1)$$

794

795 where pF is defined as the logarithm of the absolute value of the pressure head in
 796 centimetres. The first two terms focus on the mean squared error (MSE) between
 797 predicted and measured volumetric water contents, differentiated by wet-end ($pF \leq$
 798 4.2) and dry-end ($pF > 4.2$) conditions. The parameters N_{wet} and N_{dry} are defined as
 799 the number of training examples from the wet-end and dry-end, respectively. The
 800 next four terms ensure the model adheres to physical laws. In particular, the third
 801 term enforces linearity at the dry end by setting the second derivatives in that
 802 region to zero for set 1 of the residual points (S_1). These residual points are specific
 803 combinations of data points generated within the input space (including sand, silt,
 804 clay, SOC, bulk density, and pF) that are used to enforce physical laws. The fourth
 805 and fifth terms bound the range of pressure head at zero water content using sets 2
 806 and 3 of residual points, and the last term forces the water content to remain
 807 constant above the air-entry pressure using set 4 of the residual points. The
 808 monotonicity constraint is enforced by constructing inherently monotonic neural
 809 network architectures (Runje and Shankaranarayana, 2023). The Lambdas (λ) are
 810 weights that determine the relative contribution of each term in the loss function.
 811 The resulting neural network PTF is capable of predicting a non-specific form of the
 812 soil water retention curve from saturation to dryness and is differentiable with
 813 respect to the pressure head.

814 *Implications and prospect:* Overall, rather than predicting parameters of a soil
 815 function independently, incorporating the physical model in the training process can
 816 guide and constrain the ML models to predict physically-based values more
 817 accurately. Additional criteria could be added to the loss function to impose physical
 818 constraints. For example, the predicted soil water retention curve could be
 819 constrained to satisfy a realistic soil evaporation characteristic length calculated
 820 from the same water retention parameters. The characteristic length values must
 821 be in a realistic range (e.g. < 1 m) due to the limitation of capillary continuity of an
 822 evaporating soil surface (Or, 2020). This approach could also be used to map soil
 823 water retention curves. If we have observations of water retention data over an
 824 area, we could predict the water retention parameters from spatial covariates by

825 minimising the observed water retention data using Eq. (A7) or Eq. (15). Yang et al.
826 (2015) provided an example of this using Bayesian hierarchical models.

827

828 **3.4 Modelling soil properties in space and time**

829 Modelling dynamic soil properties is crucial for understanding how soil changes over
830 time and for improving land management practices. Static soil properties, which are
831 assumed to be relatively constant over time, are mapped based on their spatial
832 relationships with the landscape. In contrast, dynamic properties, such as soil water
833 content, SOC, and nutrient availability, vary with time due to environmental and
834 anthropogenic factors. Some properties change more rapidly than others (e.g. soil
835 temperature versus soil pH), making it important to gauge the timescale of their
836 prediction. Process-based models are effective in accounting for major soil
837 processes within specific soil profiles or layers, but they require calibration to local
838 conditions. The spatial application of these models can be challenging due to limited
839 data for model initialisation and parameterisation and significant computational
840 demands. On the other hand, ML models excel in spatial modelling but lack the
841 capability to simulate processes.

842 To model the dynamic soil properties in space and time using the SoilML framework,
843 several techniques can be used:

844 - Residual models: This approach involves using a ML model to predict the residuals
845 of a physical-based model. It involves learning the errors in the physical-based
846 model prediction as compared to observations and using this information to correct
847 the predictions of the physical model (Willard et al., 2020). This residual modelling
848 approach only learns what components are missing from the physical model and
849 does not incorporate any informed knowledge.

850 - Meta or surrogate models: This approach involves using ML models to emulate
851 physical-based models. This involves generating scenarios of various input soil and
852 climate variables and running them through a simulation model to obtain simulation
853 results that can be used as training data. An ML model is then trained to model the
854 output as a function of these inputs (Perlman et al., 2014). The ML could identify the

855 sensitivity of the physical model and key variables influencing the model output,
856 identifying under- or overrepresented inputs (Luo et al., 2019).

857 - Hybrid models (combination of ML and process-based models): Integrating ML
858 models with process-based models can significantly improve the capacity to model
859 soil properties across space and time. Soil data are often spatially well represented
860 but temporally sparse. One approach involves using the outputs of the process-
861 based model as additional training data (both spatially and temporally) for the ML
862 model (Ma et al., 2019). Zhang et al. (2023) and Zhang et al. (2024) demonstrated
863 that incorporating process-based model outputs as supplementary training data for
864 ML models leads to higher prediction accuracy for soil carbon than using standalone
865 ML models. Another approach is to use the output of process-based models over an
866 area as a dynamic covariate, in combination with other static covariates, in an ML
867 model. Xie et al. (2022) integrated the predictions of process-based soil carbon as a
868 dynamic covariate into a DSM model and found improved prediction accuracy
869 compared to both the standalone ML models and process-based models.

870 - Physics-Informed Neural Networks (PINN): Neural networks are used as function
871 approximators by embedding physical laws in the learning process. The physical
872 laws can be described using partial differential equations (PDE), allowing neural
873 networks to model complex behaviours and dynamics accurately. PINN uses the
874 backpropagation method of neural networks to calculate the partial derivative of the
875 differential equations, and thus, the neural network solution adheres to the physical
876 equations and observations during training (Bandai and Ghezzehei, 2021; Cai et al.,
877 2021). Applications of PINN in soil studies include retrieval of soil moisture using
878 GNS reflectometry (Kilane, 2024) and soil water and heat flow (Wang et al.,
879 2023b).

880

881 *3.4.1 Case study: Modelling soil water flow using Physics-Informed Neural Network* 882 *(PINN)*

883 *Conventional approaches:* Soil moisture dynamics can be described by the
884 Richardson-Richards' equation (Richards, 1931; Richardson, 1922), which is based
885 on the conservation of mass and the Darcy-Buckingham law (Buckingham, 1907).
886 The Richards equation incorporates water retention curves and hydraulic

887 conductivity functions to encode macroscopic soil hydraulic properties on the scale
888 of interest. Commonly, parametric models (e.g., the Mualem- van Genuchten model)
889 are used to represent the soil hydraulic functions. Their parameters are estimated
890 via inverse modelling, where the parameters are adjusted by repeatedly solving
891 Richards' equation to match the model output with the data. A widely used
892 software, HYDRUS, has such an inverse modelling capability, where a finite element
893 method solves Richards' equation (Šimůnek et al., 2016). A limitation of this
894 approach is the inflexibility of the parametric models used to represent the soil
895 hydraulic functions. For example, if the Mualem-van Genuchten model is used as
896 the parametric model, inverse modelling would fail if the target soil's water
897 retention curve exhibits a bimodal shape.

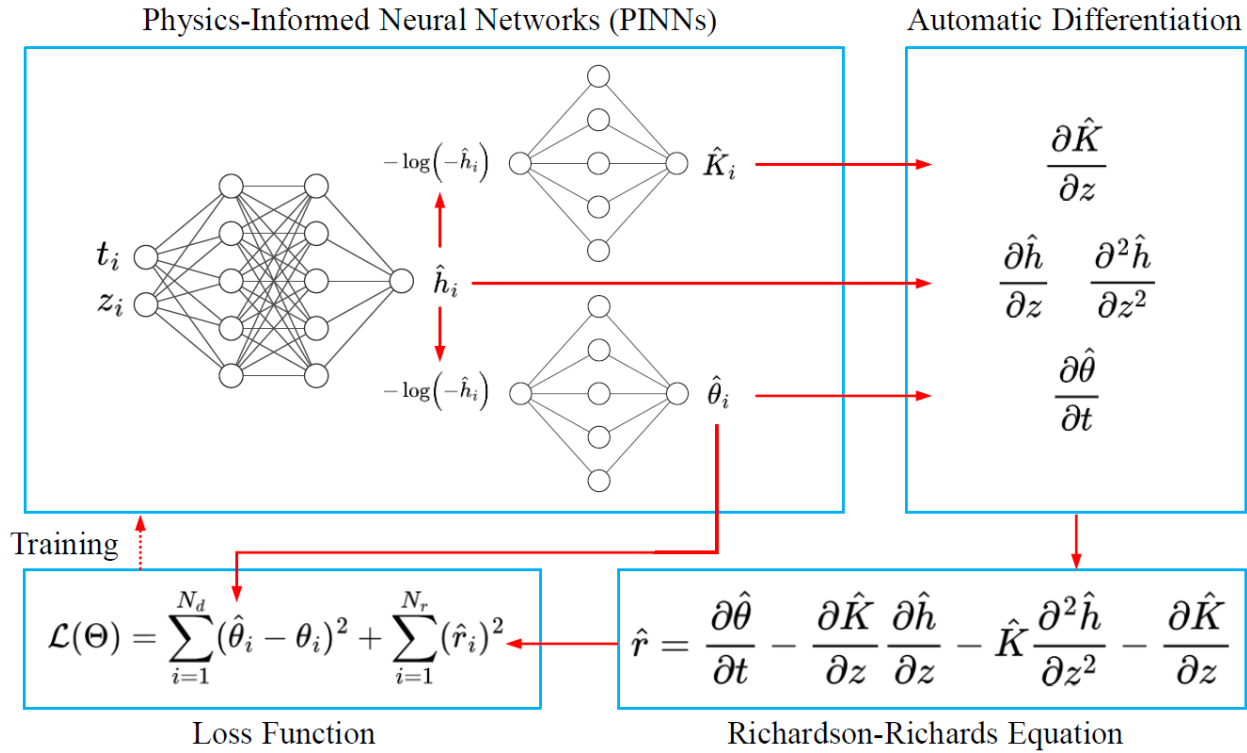
898 *SoilML*: Several studies have proposed a physics-informed neural network (PINN)
899 approach for inverse modelling based on Richards' equation to improve the
900 capability to analyse soil moisture data. In the original PINNs proposed by Raissi et
901 al. (2019), fully-connected neural networks are used to represent the solution to
902 partial differential equations as a function of the temporal and spatial coordinates.
903 The neural networks are trained to minimise a loss function consisting of an
904 observation constraint term and a PDE residual term (see Eq. 1). The PDE residual
905 term can be computed by automatic differentiation (Baydin et al., 2018), which is
906 implemented in the neural networks framework. Tartakovsky et al. (2020) employed
907 PINNs to estimate the hydraulic conductivity function for a time-independent two-
908 dimensional Richards' equation. Subsequently, Bandai and Ghezzehei (2021)
909 developed PINNs for the time-dependent one-dimensional Richards' equation to
910 estimate both water retention curves and hydraulic conductivity functions. In their
911 framework (Figure 11), two monotonically constrained neural networks (Daniels and
912 Velikova, 2010) are used to represent the soil hydraulic functions.

913 Through numerical experiments, they demonstrated that the PINNs framework has
914 the potential to estimate soil hydraulic functions without initial and boundary
915 conditions. Furthermore, several studies have improved upon their PINNs
916 framework. To improve the stability of PINNs against sparse and noisy data, Depina
917 et al. (2021) replaced the monotonic neural networks with the Mualem-van
918 Genuchten model and estimated the model's parameters via a global optimisation
919 algorithm. They validated their approach using both synthetic data and laboratory

920 infiltration experimental data. The model has been extended to model layered soils
921 (Bandai and Ghezzehei, 2022). Recently, the PINNs approach has been extended to
922 multi-physics problems in vadose zone hydrology, such as solute transport in
923 unsaturated soils (Haruzi and Moreno, 2023) and coupled heat and water transport
924 (Wang et al., 2023b).

925 While PINNs have shown some success using synthetic data and laboratory
926 experimental data, it remains difficult to train PINNs and obtain consistent results
927 with field-observed soil moisture data due to the limited amount and accuracy of
928 the data. Additionally, training PINNs for long temporal domains is challenging
929 because the original formulation of PINNs does not encode the temporal causality of
930 dynamical systems. Although many methods have been proposed to alleviate those
931 issues, as discussed in Wang et al. (2023a), they have not yet been applied to soil
932 processes.

933 As an alternative approach, Bandai et al. (2024) proposed a hybrid method where
934 the Richards' equation is solved using a traditional numerical method (i.e., finite
935 volume method with Backward Euler method), and neural networks are embedded
936 in the numerical model to represent soil hydraulic functions. This approach
937 leverages the flexibility of neural networks to represent unknown functions in
938 physics-based models (e.g., soil hydraulic functions) while ensuring the basic
939 physics encoded in Richards' equation is maintained. This contrasts with the PINNs
940 approaches, where Richards' equation is enforced in a soft manner as a loss term in
941 the loss function, and therefore, the basic physics laws, such as the conservation of
942 mass, Buckingham-Darcy law, and temporal causality, are not guaranteed. They
943 demonstrated that their neural network approach better fit infiltration experimental
944 data than using the Mualem- van Genuchten model.



945

946

947 Figure 11. Physics-informed neural networks for the Richardson-Richards equation.
 948 The temporal and space coordinates (t and z) are fed into a fully connected neural
 949 network (a) to calculate the water potential h , which is then further converted into
 950 the hydraulic conductivity K and the volumetric water content θ by two monotonic
 951 neural networks ((b) and (c)), respectively. The three neural networks are trained
 952 simultaneously by minimising the loss function consisting of the data misfit term
 953 and the residual of the Richardson-Richards equation (Bandai and Ghezzehei, 2021).

954

955 3.4.2 Case study: Soil temperature modelling

956 Soil temperature is influenced by various soil properties, such as thermal
 957 conductivity and heat capacity, which are affected by factors like bulk density,
 958 moisture content, and organic matter (Jury and Horton, 2004). While ML models are
 959 often used to predict soil temperature based on historical data, these models can be
 960 limited by their reliance on specific observation periods and may not fully capture
 961 the underlying causes of temperature variations (Lembrechts et al., 2022). Although

962 ML algorithms can identify nonlinear relationships between soil temperature and air
963 temperature along with other climate variables, they are generally incapable of
964 showing the physical processes involved. Another drawback of uninformed ML
965 models is their reliance on large amounts of data for reliable calibration and lack
966 generalisability. Time series of soil temperature data are commonly recorded at
967 meteorological stations, providing complete temporal coverage but sparse spatial
968 coverage.

969 Soil temperature is governed by soil thermal properties, which vary with soil
970 moisture level (assuming constant organic matter and bulk density over time)
971 (Ochsner et al., 2001). The soil heat capacity and thermal conductivity determine
972 the heat flow rate and, consequently, temperature change and fluctuations over
973 time. Considering a static moisture level (i.e., at field capacity or wilting point), soil
974 heat flow rate can vary due to spatial variations in soil bulk density, organic matter
975 content, texture, and mineralogy, as these variables affect soil thermal diffusivity
976 even under constant solar radiation and other environmental conditions.

977 Soil heat capacity and thermal conductivity can be used to calculate soil
978 temperature in space and time using a physical rules-informed model such as the
979 standard heat flow equation. The volumetric heat capacity of soil is defined as the
980 amount of heat required to raise a unit volume of soil by one degree of
981 temperature. As soil is a composite of air, water, and solid materials, soil heat
982 capacity is described by the heat capacities of all the constituents, weighted by
983 their volumetric fractions. Thus, volumetric soil heat capacity can be expressed as:

$$984 \quad C_{soil} = X_a C_a + X_w C_w + \sum_{j=1}^N X_{sj} C_{sj} \quad (16)$$

985 where X refers to the volume fraction, C is volumetric heat capacity, and the
986 subscripts a , w , and sj refer to air, water, and solid constituent j (for N different solid
987 materials in the soil). Soil thermal conductivity quantifies the rate at which heat
988 energy is conducted through a unit area of soil under a unit temperature gradient in
989 a direction perpendicular to the area. While soil thermal conductivity can be directly
990 measured, it can also be estimated using PTFs (e.g. He et al., 2020; Wessolek et al.,
991 2023; Zhang and Wang, 2017).

992 The amount of thermal energy that moves through an area of soil in a unit of time is
993 known as soil heat flux or heat flux density. The ability of a soil to conduct heat
994 determines how fast its temperature changes during the day or between seasons.
995 The magnitude of this heat energy is a component of the soil surface energy
996 balance, which varies with surface cover, moisture content, and solar irradiance.
997 Heat energy is transported through soil by several mechanisms, including
998 conduction, convection of heat by flowing liquid water and moving air, convection of
999 latent heat, and radiation. However, the most important heat transfer in soil is by
1000 conduction, which refers to the heat transported by molecular collisions. The
1001 conductive heat flux for a pure substance in one dimension is described by Fourier's
1002 law:

$$1003 \quad J_{HC} = \lambda \frac{dT}{dz} \quad (17)$$

1004 where J_{HC} is the amount of thermal energy, λ is thermal conductivity, T is
1005 temperature, and z is soil depth. Combined with the continuity equation, we have
1006 the general heat transport equation that describes the change in temperature with
1007 time (Carslaw and Jaeger, 1959):

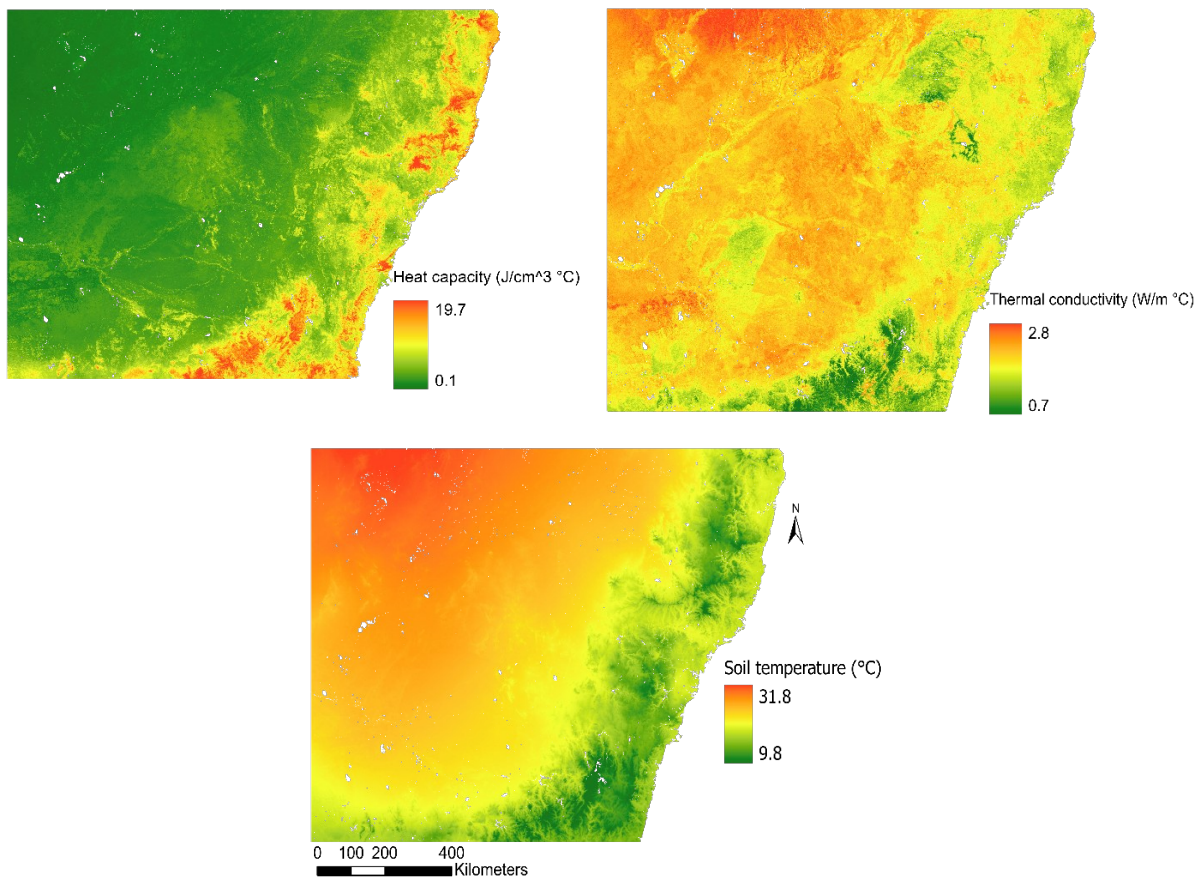
$$1008 \quad C \frac{\partial T}{\partial t} = \lambda \frac{\partial^2 T}{\partial z^2} \quad (18)$$

1009 Figure 12 demonstrates that using physics-based equations, we could calculate soil
1010 temperature from air temperature directly by considering soil thermal properties.
1011 Since heat transport governs soil temperature, a danger of blind use of ML in soil
1012 temperature modelling is the lack of physical rules in the prediction. For example,
1013 the boundary conditions (at the soil surface and different depths), temporal
1014 fluctuations of surface temperature, and heat flow variation to varying depths due
1015 to differences in soil thermal properties can significantly affect soil temperature
1016 estimation (Cichota et al., 2004; Gao et al., 2017; Ouzzane et al., 2014).

1017 *Implications and prospects:* The nature of soil temperature dynamics and its spatio-
1018 temporal variations often prevent ML models from recognising the underlying
1019 phenomena and understanding why they vary at different scales. While efforts have
1020 been made to incorporate physical knowledge into ML to make it interpretable (e.g.
1021 Abimbola et al., 2021; Li et al., 2022), the specific physical rules driving soil

1022 temperature within ML are still not well recognised. Xie et al. (2024) derived a PINN
1023 model for soil temperature prediction based on the heat transport equation (Eq. 18).
1024 Trained using historical soil temperature data, they were able to predict one-
1025 dimensional soil temperature at multiple soil depths accurately. ML models
1026 incorporating remote sensing products (such as land cover and land surface
1027 temperature) can help determine related soil and environmental properties that are
1028 not considered in the physical model. These properties include land cover, the rate
1029 of heat transport from the atmosphere, and other factors influencing soil
1030 temperature. Future work should focus on combining physical rules within ML
1031 algorithms to improve the accuracy and reliability of soil temperature predictions.

1032



1033

1034 Figure 12. The figure shows an example of soil heat capacity and thermal
1035 conductivity for New South Wales state, Australia, at 5–15 cm soil depth layer with a
1036 moisture level of 60% of the field capacity. The lower map is the February 2022

1037 average temperature map for the 5–15cm soil depth, obtained using a steady-state
1038 analytical method.

1039

1040 *3.4.3 Prospective case study of modelling soil carbon dynamics*

1041 *Conventional approaches:* There is increasing interest in modelling SOC dynamics
1042 and SOC changes to infer carbon emission or sequestration from the atmosphere.
1043 Both physical rules-based and ML models are used to estimate SOC changes across
1044 regions. Physical rules-based models typically define several SOC pools, each
1045 characterised by a mean residence time derived from first-order kinetics.
1046 Conversely, ML models are often trained on sparse temporal soil data along with
1047 both static and dynamic covariates (Sun et al., 2021; Yang et al., 2022). Static
1048 covariates include soil characteristics, topography, and long-term climate patterns.
1049 Dynamic covariates often involve temporal data like land use or land cover and
1050 vegetation indices from remote sensing images. There are also cumulative temporal
1051 indices such as total rainfall since a specific period, and years since land cover
1052 changes to reflect temporal dynamics (Padarian et al., 2022b). These dynamic
1053 covariates track SOC dynamics from one state to another to infer SOC changes.
1054 However, the time scale of the changes in the remotely-sensed images is not
1055 aligned with the SOC dynamic changes, and fails to explain the underlying
1056 processes behind SOC changes. For example, a change in land cover from forest to
1057 cropping field will cause SOC to decline rapidly but this process could take several
1058 years to reach a steady state. Similarly, a change from cropping field to pasture
1059 could accumulate SOC slowly and take several years to achieve equilibrium. In
1060 addition, SOC changes can occur over shorter periods, such as crop rotation (Fang
1061 et al., 2018), or longer periods, such as decomposition. Furthermore, surface
1062 conditions detected by remote sensing may not adequately represent subsurface
1063 processes. Integrating these processes into models is crucial for a more
1064 comprehensive understanding of SOC changes.

1065 *SoiML:* Physical rules-based SOC models often struggle to accurately resolve spatial
1066 information because decomposition constants may vary with soil types and
1067 topography. Conversely, ML models, lacking process-based insights, tend to produce
1068 abrupt changes in SOC when there is a shift in land use from one period to the next.

1069 Physical rules-based approaches can be particularly useful in addressing the
1070 problem of limited observational data in soil carbon dynamics modelling. We can
1071 create an ML model that allows SOC dynamics to follow physical rules. Here, we
1072 provide a framework for integrating various soil properties and environmental
1073 factors, offering a way to enhance model reliability in predicting soil carbon
1074 changes.

1075 SOC is observed in space and time $C_{x,t}$ and can be modelled as a mass balance of
1076 production and input (I), decomposition (k). The evolution of the organic carbon
1077 content (C) for a particular soil depth through space and time can therefore be
1078 expressed as (Andrén and Kätterer, 1997):

1079
$$\frac{dC}{dt} = I - k_1 C_1 + h k_1 C_1 - k_2 C_2 \quad (19)$$

1080 Soil carbon change over time ($\frac{dC}{dt}$) can be modelled as consisting of a fast pool (C_1)
1081 and a slow pool (C_2). Carbon input I enters the fast pool with a rate constant of k_1 ,
1082 which in turn becomes humified and mineral-associated at a rate of h into the slow
1083 pool, which has a rate constant of k_2 . The input I depends on the types of organic
1084 matter, above and below-ground, climate, soil type, depth, and management
1085 practices. The humification and decomposition constants vary in space and time,
1086 influenced by temperature, clay content, and moisture levels.

1087 First, we build a neural network model which will predict parameters $\phi = [h, k_1, k_2]$
1088 from soil characteristics and factors related to climate, topography, vegetation and
1089 human activities:

1090
$$\phi = f(\text{soil}, \text{topography}, \text{climate}, \text{vegetation}) \quad (20)$$

1091 A second neural network can be constructed to predict soil C in space and time that
1092 conforms to the C dynamics equation. The neural networks would incorporate static
1093 inputs such as soil texture, topography, long-term mean rainfall, and temperature,
1094 along with dynamic inputs such as land use and vegetation indices, together with
1095 output from the first neural network:

1096
$$C_{x,t} = f(\text{soil}, \text{topography}, \text{climate}, \text{vegetation}, I_{x,t}, k_{1x}, k_{2x}, h_x) \quad (21)$$

1097 where the input I is a function of:

1098 $I_{x,t} = f(\text{soil}, \text{topography}, \text{climate}, \text{vegetation})$ (22)

1099 The network can be trained based on observed C concentration, with a loss
 1100 function:

1101 $L = w_1 \sum (C_{x,t} - \hat{C}_{x,t})^2 + \lambda w_2 \sum \left(\frac{dC}{dt} + \frac{d\hat{C}}{dt} \right)^2$ (23)

1102 The first term accounts for the sparsely observed carbon in space and time, and the
 1103 second term provides a constraint that the model will adhere to the dynamics
 1104 defined in Equation (19), based on a series of long-term experimental data or
 1105 simulated data. The terms w_1 and w_2 refer to the weights for the first and second
 1106 terms.

1107

1108 **4. Discussion, Assumptions and Limitations**

1109 We have demonstrated a variety of forms of soil science knowledge and how they
 1110 can be incorporated into the ML training process. Here, we discuss aspects of ML
 1111 that can improve soil science understanding.

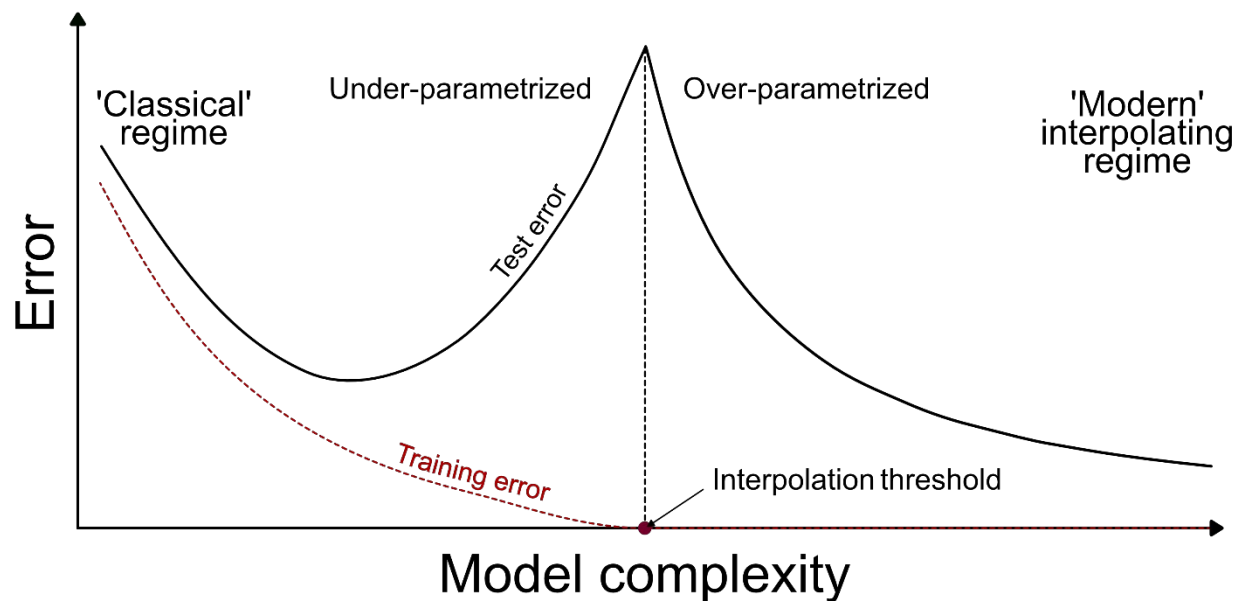
1112 *Soil is a unique 3-D volume* and ML models should be soil science-informed. It is
 1113 important to design ML model applications that specifically accommodate the
 1114 multidimensional nature of soil (Poggio and Gimona, 2014). Soil is not just a simple
 1115 substrate but a three-dimensional body with unique properties varying by depth.
 1116 Thus, when developing ML models, they need to be soil science-informed,
 1117 incorporating architecture, variables, and data layers that reflect the unique
 1118 characteristics of the soil.

1119 *Modify the ML models to suit our needs, not modify our data to suit ML needs.* This
 1120 means a shift in how we approach ML development. Traditionally, much of the focus
 1121 in ML has been on adjusting, filtering, or transforming soil data to fit the
 1122 requirements of existing algorithms and models. This approach could lead to loss of
 1123 information or oversimplification of soil data. Instead, we should adapt ML models to
 1124 work with soil data, such as modifying and regularising their loss functions. This soil-
 1125 centric approach in model development ensures that the technology serves the
 1126 specific needs of its applications, rather than forcing data into predefined,
 1127 potentially limiting frameworks.

1128 *Overparametrisation and interpolation.* Classical statistics promotes Occam's razor
1129 principle, which suggests selecting the hypothesis with the fewest assumptions
1130 among competing hypotheses. This translates to the preference for models with
1131 fewer parameters, as they are easier to interpret and less likely to overfit the data.
1132 However, ML models are usually overparametrised, having many more parameters
1133 compared to the size of the training data (Belkin, 2021). Some models (such as
1134 random forest and boosting methods) are designed to perfectly fit (interpolate) the
1135 training data, which is usually noisy. Interpolating noisy data using ML models,
1136 traditionally associated with detrimental overfitting, has been demonstrated to
1137 perform well on test data (Belkin, 2021).

1138 Belkin et al. (2019) proposed the double descent phenomenon in ML where the
1139 error, when plotted against model complexity, shows a two-phase behaviour
1140 contrary to the traditional U-shaped bias-variance trade-off (Figure 13). Initially, as
1141 model complexity increases, the test error decreases up to an interpolation
1142 threshold where the model perfectly fits the training data, causing the test error to
1143 peak. Unexpectedly, if the complexity continues to increase beyond this point, the
1144 test error decreases again, leading to a second descent. This phenomenon has been
1145 observed across various ML models. However, the ability of ML modes to interpolate
1146 does not necessarily mean the model is more accurate and generalisable. Practices
1147 that focus on regularising the training rather than achieving a perfect fit are being
1148 advocated (See Box 3). Currently, there is still a lack of metrics to quantify an ML
1149 model's complexity with respect to its ability to generalise (Dar et al., 2021). Soil
1150 science data are usually relatively small compared to disciplines such as image or
1151 language processing. There is still a lack of understanding of the interpolation effect
1152 of ML models trained on a relatively small dataset (e.g. less than 100 observations).
1153 Overparametrisation and the double descent phenomena do not necessarily lead to
1154 improved accuracy. Thus, an independent validation dataset is needed to evaluate
1155 the generalisability of the ML models.

1156



1157

1158 Figure 13. The double-descent error curve with “classical” and “modern
 1159 interpolating” regimes, showing training error (red dashed line) and test error (solid
 1160 line) as a function of model complexity. The left curve is the classical U-shaped risk
 1161 curve arising from the bias-variance trade-off. The right curve, separated by the
 1162 interpolation threshold, represents ML models with interpolation or zero training
 1163 error. From Belkin et al. (2019).

1164

1165 *Uncertainty analysis.* This analysis helps acknowledge the limits of models. In soil
 1166 science, these challenges are magnified due to the limited, sparse, and often
 1167 heterogeneous nature of soil data (Libohova et al., 2019). In ML, uncertainties can
 1168 stem from uncertain data and incomplete knowledge. Although some studies report
 1169 prediction intervals and confidence levels, a comprehensive approach to uncertainty
 1170 quantification remains a challenge, emphasizing the need for better methods to
 1171 evaluate and communicate the reliability of soil property predictions. Uncertainty
 1172 quantification is essential for assessing prediction reliability, especially under
 1173 unseen scenarios. Bayesian approaches are the standard method; however, the
 1174 computational demands of modelling the full posterior distribution are very high.
 1175 This challenge can be mitigated by using dropout techniques to approximate the
 1176 posterior distribution. For example, Padarian et al. (2022a) used the Monte Carlo

1177 dropout method for the prediction of SOC using vis-NIR-SWIR data, and
1178 demonstrated its capability to identify large uncertainty when new data presented is
1179 different from the data used during training.

1180 *Interpretability.* Soil scientists are also interested in using ML models to have a
1181 better understanding of soil processes and formulate hypotheses (Padarian et al.,
1182 2020a; Wadoux and Molnar, 2022). The prevailing assumption in ML models is that,
1183 with sufficiently covariate capturing spatial dependence relationships, the spatial
1184 patterns of soils can be predicted, and the drivers of those spatial patterns can be
1185 identified (Bui, 2004; Bui et al., 2020). Interpretable ML models help clarify the
1186 significance of specific predictors in estimating soil properties, tackling the "black
1187 box" nature of many ML algorithms (Roscher et al., 2020). *Post hoc* analysis—
1188 including interpreting, visualising, and evaluating ML predictions—can determine
1189 how well a model aligns with established soil science knowledge. However, while ML
1190 primarily aims to minimise prediction error, soil scientists are more interested in
1191 uncovering underlying processes. Interpretability relies on human judgment, and
1192 just because a model highlights certain predictors as important doesn't necessarily
1193 mean they cause the observed effects or offer new insights. Techniques like Shapley
1194 values, while useful for generating hypotheses, could lead to biased conclusions if
1195 not carefully handled. Therefore, interpretability should not be conflated as a
1196 verification of a model's generalisability. It also should not be viewed as a
1197 confirmation of a model's accuracy but must be integrated with domain knowledge
1198 to validate and enhance predictions thoroughly.

1199 *Dynamic soil properties prediction in space and time.* Soil data are often
1200 incomplete, noisy, and sparse in space and time. ML models, especially tree models,
1201 often struggle to interpolate these sparse data effectively. To overcome these
1202 limitations, it is beneficial to define the model's structure consistent with soil
1203 science principles, incorporate prior knowledge of soil in space (and time), and
1204 define a loss function that obeys physical principles. ML models are also often used
1205 to predict soil carbon fate under future climate scenarios, yet ML models usually
1206 perform poorly when used for extrapolation or predicting unseen or rare events.

1207 Finally, SoilML should ensure *Reproducibility*, which involves a systematic approach
1208 to documenting, sharing, and verifying the processes and results of analysis. This

1209 includes making the codes, and methodologies accessible to other researchers so
1210 that they can replicate the findings. Thorough evaluations are also crucial to identify
1211 any biases inherent in the model, data, or methodology. This might involve testing
1212 the model under various conditions, using diverse datasets to check for consistency,
1213 and scrutinising the assumptions underlying the model's predictions.

1214

1215 **Box 3. Reducing overparametrisation in ML models**

1216 As ML is a data-hungry model, efforts are being made to reduce
1217 overparametrisation by either increasing the number of observations or simplifying
1218 the model architecture. Some of the approaches include (Dar et al., 2021):

1219 - Data augmentation, generating additional synthetic data can increase the
1220 diversity of training data, which can help reduce overparametrisation. Data
1221 augmentation could include applying various covariates transformations (such as
1222 rotations) to existing data or based on prior information (see Figure 4).

1223 - Transfer learning, involves using a pre-trained deep neural network (DNN) on
1224 a related problem to improve training efficiency and performance on a target
1225 problem with less data (Padarian et al., 2019a). Transfer learning is done by
1226 transferring and fine-tuning one or more layers from the source DNN. This approach
1227 can mitigate overparametrisation by leveraging the learned parameters from the
1228 source task.

1229 - Pruning models, pruning highly parametrised ML models into less complex
1230 forms can improve the trade-off between generalisation performance and
1231 computational requirements. This approach is also useful for applications with
1232 limited storage space, computation time, and energy consumption.

1233

1234 **5. Outlook**

1235

1236 SoilML have been applied in four key areas: digital soil mapping, soil spectroscopy,
1237 pedotransfer functions, and dynamic soil property modeling. These applications

1238 demonstrate how SoilML enhances model accuracy, improves interpretability, and
1239 preserves the principles of soil science.

1240 ML approaches have successfully produced digital soil maps of continents and the
1241 world, but there is still a lack of modelling soil processes. Soil processes are
1242 commonly predicted statically using ML models without considering temporal
1243 processes. Currently most ML-derived outputs of soil maps are used as inputs of
1244 baseline data for assessing future conditions using process-based models. While
1245 physics-informed ML models are growing in environmental and earth modelling,
1246 their application in soil science is still minimal. With an increasing demand for
1247 quantifying soil functions, there is a need and potential to upscale our physical and
1248 chemical process models to a larger extent. SoilML has the potential to accelerate
1249 advancements in soil science by integrating soil-specific knowledge into the ML
1250 process. This can be achieved through the use of observational priors, tailored
1251 model structures, and informed loss functions that incorporate physical constraints
1252 and coherency rules.

1253 There are still practical challenges of SoilML include high computational demands,
1254 the need for soil-specific priors, and difficulties in integrating multi-source data with
1255 varying spatial and temporal resolutions. Effective collaboration among the
1256 communities, including process-based modellers, pedometricians, remote sensing
1257 experts, and data scientists, is essential to advance the growth of SoilML for soil
1258 security assessment.

1259

1260 **Acknowledgements**

1261 BM acknowledges funding from ARC Discovery Project Forecasting Soil Conditions,
1262 DP200102542.

1263

1264 **Declaration of Generative AI in the writing process**

1265 In the preparation of this paper, ChatGPT 4o, was employed to assist with improving
1266 the clarity, grammar, and overall quality of the English language. Following its use,
1267 the authors thoroughly reviewed and edited the content to ensure accuracy and
1268 alignment with the intended meaning. The AI's involvement was limited to language

1269 refinement, and no content generation, interpretation of data, or intellectual
1270 contributions were made by the AI. The responsibility for the scientific content, and
1271 conclusions of this paper lies entirely with the authors.

1272

1273

1274 **Appendix**

1275 The following are the equations of the water retention and hydraulic conductivity
 1276 curves according to the FXW model.

1277 The water retention curve according to Fredlund and Xing (1994):

1278 $\theta(h) = \theta_s S_e(h)$ (A1)

1279 where θ_s is the saturated water content and $S_e(h)$ is the effective saturation,
 1280 calculated as:

1281 $S_e(h) = C_f(h) \Gamma(h)$ (A2)

1282 With $C_f(h) = \left[1 - \frac{\ln\left(1 + \frac{h}{h_r}\right)}{\ln\left(1 + \frac{h_0}{h_r}\right)} \right]$ and $\Gamma(h) = \left(\ln[\exp(1) + |ah|^n] \right)^{-m}$ (A3)

1283 The unsaturated hydraulic function, $K(h)$, based on Wang et al. (2018) is defined as:

1284 $K(h) = K_s K_r(h)$ (A4)

1285 where K_s is the saturated hydraulic conductivity and

1286 $K_r(h) = S_{ek}^L \gamma^2$ (A5)

1287 $S_{ek}(h) = \frac{\Gamma(h) - \Gamma(h_0)}{1 - \Gamma(h_0)}$ and $\gamma = \left[1 - \left(1 - \Gamma^{\frac{1}{m}} \right)^{1 - \frac{1}{n}} \right]^2$ (A6)

1288 The neuroFX model minimises the following functions:

1289 Min: $L_\theta = \sum \left(\theta(h) - \hat{\theta}(h \vee \hat{\theta}_s, \hat{a}, \hat{n}, \hat{m}) \right)^2$ (A7)

1290 Min: $L_K = \sum \left(\log K(h) - \log \hat{K}(h \vee \hat{\theta}_s, \hat{a}, \hat{n}, \hat{m}, \hat{K}_s, L) \right)^2$ (A8)

1291

1292

1293

1294

1295

1296 **References**

- 1297 Abimbola, O.P., Meyer, G.E., Mittelstet, A.R., Rudnick, D.R., Franz, T.E., 2021.
1298 Knowledge-guided machine learning for improving daily soil temperature
1299 prediction across the United States. *Vadose Zone Journal* 20(5), e20151.
- 1300 Andrén, O., Kätterer, T., 1997. ICBM: The introductory carbon balance model for
1301 exploration of soil carbon balances. *Ecological Applications* 7(4), 1226-1236.
- 1302 Babaeian, E., Sadeghi, M., Jones, S.B., Montzka, C., Vereecken, H., Tuller, M., 2019.
1303 Ground, proximal, and satellite remote sensing of soil moisture. *Reviews of*
1304 *Geophysics* 57(2), 530-616.
- 1305 Bagnall, D.K., Morgan, C.L.S., Cope, M., Bean, G.M., Cappellazzi, S., Greub, K.,
1306 Liptzin, D., Norris, C.L., Rieke, E., Tracy, P. et al., 2022. Carbon-sensitive
1307 pedotransfer functions for plant available water. *Soil Science Society of*
1308 *America Journal* 86(3), 612-629.
- 1309 Bandai, T., Ghezzehei, T.A., 2021. Physics-informed neural networks with
1310 monotonicity constraints for Richardson-Richards equation: Estimation of
1311 constitutive relationships and soil water flux density from volumetric water
1312 content measurements. *Water Resources Research* 57(2), e2020WR027642.
- 1313 Bandai, T., Ghezzehei, T.A., 2022. Forward and inverse modeling of water flow in
1314 unsaturated soils with discontinuous hydraulic conductivities using physics-
1315 informed neural networks with domain decomposition. *Hydrol. Earth Syst. Sci.*
1316 26(16), 4469-4495.
- 1317 Bandai, T., Ghezzehei, T.A., Jiang, P., Kidger, P., Chen, X., Steefel, C.I., 2024. Learning
1318 constitutive relations from soil moisture data via physically constrained
1319 neural networks. *Water Resources Research* 60(7), e2024WR037318.
- 1320 Baydin, A.G., Pearlmutter, B.A., Radul, A.A., Siskind, J.M., 2018. Automatic
1321 differentiation in machine learning: a survey. *J. Mach. Learn. Res.* 18(1), 5595-
1322 5637.
- 1323 Behrens, T., Schmidt, K., Zhu, A.X., Scholten, T., 2010. The ConMap approach for
1324 terrain-based digital soil mapping. *European Journal of Soil Science* 61(1),
1325 133-143.
- 1326 Belkin, M., 2021. Fit without fear: remarkable mathematical phenomena of deep
1327 learning through the prism of interpolation. *Acta Numerica* 30, 203 - 248.
- 1328 Belkin, M., Hsu, D., Ma, S., Mandal, S., 2019. Reconciling modern machine-learning
1329 practice and the classical bias-variance trade-off. *Proceedings of the National*
1330 *Academy of Sciences* 116(32), 15849-15854.
- 1331 Bogrekci, I., S. Lee, W., 2006. Effects of soil moisture content on absorbance spectra
1332 of sandy soils in sensing phosphorus concentrations using UV-VIS-NIR
1333 spectroscopy. *Transactions of the ASABE* 49(4), 1175-1180.
- 1334 Buckingham, E., 1907. *Studies on the Movement of Soil Moisture*. US Department of
1335 Agriculture, Bureau of Soils, Washington DC.
- 1336 Bui, E.N., 2004. Soil survey as a knowledge system. *Geoderma* 120(1), 17-26.
- 1337 Bui, E.N., Searle, R.D., Wilson, P.R., Philip, S.R., Thomas, M., Brough, D., Harms, B.,
1338 Hill, J.V., Holmes, K., Smolinski, H.J. et al., 2020. Soil surveyor knowledge in
1339 digital soil mapping and assessment in Australia. *Geoderma Regional* 22,
1340 e00299.
- 1341 Cai, S., Wang, Z., Wang, S., Perdikaris, P., Karniadakis, G.E., 2021. Physics-informed
1342 neural networks for heat transfer problems. *Journal of Heat Transfer* 143(6).
- 1343 Campbell, G., Shiozawa, S., 1992. Prediction of hydraulic properties of soils using
1344 particle-size distribution and bulk density data, International workshop on

1345 indirect methods for estimating the hydraulic properties of unsaturated soils.
 1346 California: University of California, pp. 317-328.
 1347 Carslaw, H.S., Jaeger, J.C., 1959. Conduction of Heat in Solids. Second Ed. ed. Oxford
 1348 University Press, Oxford.
 1349 Castaldi, F., Palombo, A., Pascucci, S., Pignatti, S., Santini, F., Casa, R., 2015.
 1350 Reducing the influence of soil moisture on the estimation of clay from
 1351 hyperspectral data: a case study using simulated PRISMA data. *Remote Sens-*
 1352 *Basel* 7(11), 15561-15582.
 1353 Chang, C.-W., Laird, D.A., Mausbach, M.J., Hurburgh, C.R., 2001. Near-infrared
 1354 reflectance spectroscopy-principal components regression analyses of soil
 1355 properties. *Soil Science Society of America Journal* 65(2), 480-490.
 1356 Chaplot, V., Lorentz, S., Podwojewski, P., Jewitt, G., 2010. Digital mapping of A-
 1357 horizon thickness using the correlation between various soil properties and
 1358 soil apparent electrical resistivity. *Geoderma* 157(3), 154-164.
 1359 Chen, S., Saby, N.P.A., Martin, M.P., Barthès, B.G., Gomez, C., Shi, Z., Arrouays, D.,
 1360 2023a. Integrating additional spectroscopically inferred soil data improves the
 1361 accuracy of digital soil mapping. *Geoderma* 433, 116467.
 1362 Chen, S., Xue, J., Shi, Z., 2023b. Spectral-guided ensemble modelling for soil
 1363 spectroscopic prediction. *Geoderma* 437, 116594.
 1364 Cichota, R., Elias, E.A., de Jong van Lier, Q., 2004. Testing a finite-difference model
 1365 for soil heat transfer by comparing numerical and analytical solutions.
 1366 *Environmental Modelling & Software* 19(5), 495-506.
 1367 Daniels, H., Velikova, M., 2010. Monotone and partially monotone neural networks.
 1368 *IEEE Transactions on Neural Networks* 21(6), 906-917.
 1369 Dar, Y., Muthukumar, V., Baraniuk, R.G., 2021. A farewell to the bias-variance
 1370 tradeoff? An overview of the theory of overparameterized machine learning.
 1371 arXiv preprint.
 1372 Depina, I., Jain, S., Mar Valsson, S., Gotovac, H., 2021. Application of physics-
 1373 informed neural networks to inverse problems in unsaturated groundwater
 1374 flow. *Georisk: Assessment and Management of Risk for Engineered Systems*
 1375 *and Geohazards* 16(1), 21-36.
 1376 Eymard, A., Richer-de-Forges, A.C., Martelet, G., Tissoux, H., Bialkowski, A.,
 1377 Dalmasso, M., Chrétien, F., Belletier, D., Ledemé, G., Laloua, D. et al., 2024.
 1378 Exploring the untapped potential of hand-feel soil texture data for enhancing
 1379 digital soil mapping: Revealing hidden spatial patterns from field
 1380 observations. *Geoderma* 441, 116769.
 1381 Fajardo, M., McBratney, A., Whelan, B., 2016. Fuzzy clustering of Vis-NIR spectra for
 1382 the objective recognition of soil morphological horizons in soil profiles.
 1383 *Geoderma* 263, 244-253.
 1384 Fang, J., Yu, G., Liu, L., Hu, S., Chapin, F.S., 2018. Climate change, human impacts,
 1385 and carbon sequestration in China. *Proceedings of the National Academy of*
 1386 *Sciences* 115(16), 4015-4020.
 1387 Fredlund, D.G., Xing, A., 1994. Equations for the soil-water characteristic curve.
 1388 *Canadian Geotechnical Journal* 31(4), 521-532.
 1389 Gao, Z., Russell, E.S., Missik, J.E.C., Huang, M., Chen, X., Strickland, C.E., Clayton, R.,
 1390 Arntzen, E., Ma, Y., Liu, H., 2017. A novel approach to evaluate soil heat flux
 1391 calculation: An analytical review of nine methods. *Journal of Geophysical*
 1392 *Research: Atmospheres* 122(13), 6934-6949.

- 1393 Gastaldi, G., Minasny, B., McBratney, A., 2012. Mapping the occurrence and
1394 thickness of soil horizons within soil profiles, *Digital Soil Assessments and*
1395 *Beyond* CRC Press, Balkema London, pp. 145-148.
- 1396 Grunwald, S., Lowery, B., Rooney, D.J., McSweeney, K., 2001. Profile cone
1397 penetrometer data used to distinguish between soil materials. *Soil and Tillage*
1398 *Research* 62(1), 27-40.
- 1399 Haghverdi, A., Cornelis, W.M., Ghahraman, B., 2012. A pseudo-continuous neural
1400 network approach for developing water retention pedotransfer functions with
1401 limited data. *Journal of Hydrology* 442-443, 46-54.
- 1402 Haghverdi, A., Öztürk, H.S., Cornelis, W.M., 2014. Revisiting the pseudo continuous
1403 pedotransfer function concept: Impact of data quality and data mining
1404 method. *Geoderma* 226-227, 31-38.
- 1405 Haruzi, P., Moreno, Z., 2023. Modeling water flow and solute transport in
1406 unsaturated soils using physics-informed neural networks trained with
1407 geoelectrical data. *Water Resources Research* 59(6), e2023WR034538.
- 1408 Hassan-Esfahani, L., Torres-Rua, A., Jensen, A., McKee, M., 2015. Assessment of
1409 surface soil moisture using high-resolution multi-spectral imagery and
1410 artificial neural networks. *Remote Sens-Basel* 7(3), 2627-2646.
- 1411 He, H., He, D., Jin, J., Smits, K.M., Dyck, M., Wu, Q., Si, B., Lv, J., 2020. Room for
1412 improvement: A review and evaluation of 24 soil thermal conductivity
1413 parameterization schemes commonly used in land-surface, hydrological, and
1414 soil-vegetation-atmosphere transfer models. *Earth-Science Reviews* 211,
1415 103419.
- 1416 Helfenstein, A., Mulder, V.L., Hack-ten Broeke, M.J.D., van Doorn, M., Teuling, K.,
1417 Walvoort, D.J.J., Heuvelink, G.B.M., 2024. BIS-4D: mapping soil properties and
1418 their uncertainties at 25m resolution in the Netherlands. *Earth Syst.*
1419 *Sci. Data* 16(6), 2941-2970.
- 1420 Heung, B., Ho, H.C., Zhang, J., Knudby, A., Bulmer, C.E., Schmidt, M.G., 2016. An
1421 overview and comparison of machine-learning techniques for classification
1422 purposes in digital soil mapping. *Geoderma* 265, 62-77.
- 1423 Holmes, K.W., Griffin, E.A., Odgers, N.P., 2015. Large-area spatial disaggregation of a
1424 mosaic of conventional soil maps: evaluation over Western Australia. *Soil*
1425 *Research* 53(8), 865-880.
- 1426 Holmes, K.W., Griffin, E.A., van Gool, D., 2021. Digital soil mapping of coarse
1427 fragments in southwest Australia: Targeting simple features yields detailed
1428 maps. *Geoderma* 404, 115282.
- 1429 Hudson, B.D., 1992. The soil survey as paradigm-based science. *Soil Science Society*
1430 *of America Journal* 56(3), 836-841.
- 1431 Hughes, P., McBratney, A.B., Huang, J., Minasny, B., Micheli, E., Hempel, J., 2017.
1432 Comparisons between USDA Soil Taxonomy and the Australian Soil
1433 Classification System I: Data harmonization, calculation of taxonomic
1434 distance and inter-taxa variation. *Geoderma* 307, 198-209.
- 1435 Hutengs, C., Eisenhauer, N., Schädler, M., Lochner, A., Seidel, M., Vohland, M., 2021.
1436 VNIR and MIR spectroscopy of PLFA-derived soil microbial properties and
1437 associated soil physicochemical characteristics in an experimental plant
1438 diversity gradient. *Soil Biology and Biochemistry* 160, 108319.
- 1439 Jury, W.A., Horton, R., 2004. *Soil physics*. John Wiley & Sons.
- 1440 Karniadakis, G.E., Kevrekidis, I.G., Lu, L., Perdikaris, P., Wang, S., Yang, L., 2021.
1441 Physics-informed machine learning. *Nature Reviews Physics* 3(6), 422-440.

- 1442 Kashinath, K., Mustafa, M., Albert, A., Wu, J.-L., Jiang, C., Esmaeilzadeh, S.,
 1443 Azizzadenesheli, K., Wang, R., Chattopadhyay, A., Singh, A. et al., 2021.
 1444 Physics-informed machine learning: case studies for weather and climate
 1445 modelling. *Philosophical Transactions of the Royal Society A: Mathematical,*
 1446 *Physical and Engineering Sciences* 379(2194), 20200093.
- 1447 Kilane, N.G., 2024. Application of physics-informed neural network approach in soil
 1448 moisture retrieval using GNSS reflectometry, York University, Ontario,
 1449 Canada.
- 1450 Koch, J., Stisen, S., Refsgaard, J.C., Ernstsens, V., Jakobsen, P.R., Højberg, A.L., 2019.
 1451 Modeling depth of the redox interface at high resolution at national scale
 1452 using random forest and residual gaussian simulation. *Water Resources*
 1453 *Research* 55(2), 1451-1469.
- 1454 Kosugi, K.i., 1996. Lognormal distribution model for unsaturated soil hydraulic
 1455 properties. *Water Resources Research* 32(9), 2697-2703.
- 1456 Kubelka, P., Munk, F., 1931. Ein Beitrag zur Optik der Farbanstriche, 12. *Zeitschrift*
 1457 *für Technische Physik*.
- 1458 Laborczi, A., Szatmári, G., Kaposi, A.D., Pásztor, L., 2019. Comparison of soil texture
 1459 maps synthesized from standard depth layers with directly compiled products.
 1460 *Geoderma* 352, 360-372.
- 1461 Lamichhane, S., Kumar, L., Adhikari, K., 2021. Updating the national soil map of
 1462 Nepal through digital soil mapping. *Geoderma* 394, 115041.
- 1463 Lark, R.M., Bishop, T.F.A., Webster, R., 2007. Using expert knowledge with control of
 1464 false discovery rate to select regressors for prediction of soil properties.
 1465 *Geoderma* 138(1), 65-78.
- 1466 Lebeau, M., Konrad, J.-M., 2010. A new capillary and thin film flow model for
 1467 predicting the hydraulic conductivity of unsaturated porous media. *Water*
 1468 *Resources Research* 46(12).
- 1469 Lembrechts, J.J., van den Hoogen, J., Aalto, J., Ashcroft, M.B., De Frenne, P.,
 1470 Kemppinen, J., Kopecký, M., Luoto, M., Maclean, I.M.D., Crowther, T.W. et al.,
 1471 2022. Global maps of soil temperature. *Global Change Biology* 28(9), 3110-
 1472 3144.
- 1473 Li, Q., Zhu, Y., Shanguan, W., Wang, X., Li, L., Yu, F., 2022. An attention-aware LSTM
 1474 model for soil moisture and soil temperature prediction. *Geoderma* 409,
 1475 115651.
- 1476 Li, X., Nieber, J.L., Kumar, V., 2024. Machine learning applications in vadose zone
 1477 hydrology: A review. *Vadose Zone Journal* 23(4), e20361.
- 1478 Libohova, Z., Seybold, C., Adhikari, K., Wills, S., Beaudette, D., Peaslee, S., Lindbo,
 1479 D., Owens, P.R., 2019. The anatomy of uncertainty for soil pH measurements
 1480 and predictions: Implications for modellers and practitioners. *European*
 1481 *Journal of Soil Science* 70(1), 185-199.
- 1482 Liu, W., Baret, F., Xingfa, G., Qingxi, T., Lanfen, Z., Bing, Z., 2002. Relating soil
 1483 surface moisture to reflectance. *Remote Sens Environ* 81(2), 238-246.
- 1484 Lobell, D.B., Asner, G.P., 2002. Moisture Effects on Soil Reflectance. *Soil Science*
 1485 *Society of America Journal* 66(3), 722-727.
- 1486 Luo, Z., Eady, S., Sharma, B., Grant, T., Liu, D.L., Cowie, A., Farquharson, R.,
 1487 Simmons, A., Crawford, D., Searle, R. et al., 2019. Mapping future soil carbon
 1488 change and its uncertainty in croplands using simple surrogates of a complex
 1489 farming system model. *Geoderma* 337, 311-321.

- 1490 Ma, Y., Minasny, B., Demattê, J.A.M., McBratney, A.B., 2023. Incorporating soil
1491 knowledge into machine-learning prediction of soil properties from soil
1492 spectra. *European Journal of Soil Science* 74(6), e13438.
- 1493 Ma, Y., Minasny, B., Malone, B.P., Mcbratney, A.B., 2019. Pedology and digital soil
1494 mapping (DSM). *European Journal of Soil Science* 70(2), 216-235.
- 1495 Ma, Y., Minasny, B., McBratney, A., Poggio, L., Fajardo, M., 2021. Predicting soil
1496 properties in 3D: Should depth be a covariate? *Geoderma* 383, 114794.
- 1497 McBratney, A.B., Mendonça Santos, M.L., Minasny, B., 2003. On digital soil mapping.
1498 *Geoderma* 117(1), 3-52.
- 1499 Mendonça Santos, M.L., Guenat, C., Bouzelboudjen, M., Golay, F., 2000. Three-
1500 dimensional GIS cartography applied to the study of the spatial variation of
1501 soil horizons in a Swiss floodplain. *Geoderma* 97(3), 351-366.
- 1502 Michéli, E., Láng, V., Owens, P.R., McBratney, A., Hempel, J., 2016. Testing the
1503 pedometric evaluation of taxonomic units on soil taxonomy — A step in
1504 advancing towards a universal soil classification system. *Geoderma* 264, 340-
1505 349.
- 1506 Miller, B.A., Koszinski, S., Wehrhan, M., Sommer, M., 2015. Impact of multi-scale
1507 predictor selection for modeling soil properties. *Geoderma* 239-240, 97-106.
- 1508 Minasny, B., McBratney, A.B., 2002. The neuro-m method for fitting neural network
1509 parametric pedotransfer functions. *Soil Science Society of America Journal*
1510 66(2), 352-361.
- 1511 Minasny, B., McBratney, A.B., 2007. Incorporating taxonomic distance into spatial
1512 prediction and digital mapping of soil classes. *Geoderma* 142(3), 285-293.
- 1513 Mualem, Y., 1976. A new model for predicting the hydraulic conductivity of
1514 unsaturated porous media. *Water Resources Research* 12(3), 513-522.
- 1515 Norouzi, S., Pesch, C., Arthur, E., Norgaard, T., Møldrup, P., Greve, M.H., Beucher, A.,
1516 Sadeghi, M., Zaresourmanabad, M., Tuller, M. et al., 2024. Physics-informed
1517 neural networks for estimating a continuous form of the soil water retention
1518 curve from basic soil properties. *ESS Open Archive*.
- 1519 Norouzi, S., Sadeghi, M., Tuller, M., Ebrahimian, H., Liaghat, A., Jones, S.B., de Jonge,
1520 L.W., 2023. A novel laboratory method for the retrieval of the soil water
1521 retention curve from shortwave infrared reflectance. *Journal of Hydrology*
1522 626, 130284.
- 1523 Norouzi, S., Sadeghi, M., Tuller, M., Liaghat, A., Jones, S.B., Ebrahimian, H., 2022. A
1524 novel physical-empirical model linking shortwave infrared reflectance and soil
1525 water retention. *Journal of Hydrology* 614, 128653.
- 1526 Ochsner, T.E., Horton, R., Ren, T., 2001. A new perspective on soil thermal
1527 properties. *Soil science society of America Journal* 65(6), 1641-1647.
- 1528 Odgers, N.P., Sun, W., McBratney, A.B., Minasny, B., Clifford, D., 2014.
1529 Disaggregating and harmonising soil map units through resampled
1530 classification trees. *Geoderma* 214-215, 91-100.
- 1531 Or, D., 2020. The tyranny of small scales—On representing soil processes in global
1532 land surface models. *Water Resources Research* 56(6).
- 1533 Orton, T.G., Pringle, M.J., Bishop, T.F.A., Menzies, N.W., Dang, Y.P., 2020. Increment-
1534 averaged kriging for 3-D modelling and mapping soil properties: Combining
1535 machine learning and geostatistical methods. *Geoderma* 361, 114094.
- 1536 Ouzzane, M., Eslami-Nejad, P., Aidoun, Z., Lamarche, L., 2014. Analysis of the
1537 convective heat exchange effect on the undisturbed ground temperature.
1538 *Solar Energy* 108, 340-347.

1539 Padarian, J., McBratney, A.B., Minasny, B., 2020a. Game theory interpretation of
1540 digital soil mapping convolutional neural networks. *SOIL* 6(2), 389-397.

1541 Padarian, J., Minasny, B., McBratney, A.B., 2019a. Transfer learning to localise a
1542 continental soil vis-NIR calibration model. *Geoderma* 340, 279-288.

1543 Padarian, J., Minasny, B., McBratney, A.B., 2019b. Using deep learning for digital soil
1544 mapping. *SOIL* 5(1), 79-89.

1545 Padarian, J., Minasny, B., McBratney, A.B., 2020b. Machine learning and soil
1546 sciences: a review aided by machine learning tools. *SOIL* 6(1), 35-52.

1547 Padarian, J., Minasny, B., McBratney, A.B., 2022a. Assessing the uncertainty of deep
1548 learning soil spectral models using Monte Carlo dropout. *Geoderma* 425,
1549 116063.

1550 Padarian, J., Stockmann, U., Minasny, B., McBratney, A.B., 2022b. Monitoring
1551 changes in global soil organic carbon stocks from space. *Remote Sens*
1552 *Environ* 281, 113260.

1553 Perlman, J., Hijmans, R.J., Horwath, W.R., 2014. A metamodelling approach to
1554 estimate global emissions from agricultural soils. *Global Ecology and*
1555 *Biogeography* 23(8), 912-924.

1556 Peters, A., Durner, W., Iden, S., 2024. The PDI model system for parameterizing soil
1557 hydraulic properties. *Vadose Zone Journal* 23(4), e20338.

1558 Poggio, L., de Sousa, L.M., Batjes, N.H., Heuvelink, G.B.M., Kempen, B., Ribeiro, E.,
1559 Rossiter, D., 2021. SoilGrids 2.0: producing soil information for the globe with
1560 quantified spatial uncertainty. *SOIL* 7(1), 217-240.

1561 Poggio, L., Gimona, A., 2014. National scale 3D modelling of soil organic carbon
1562 stocks with uncertainty propagation — An example from Scotland. *Geoderma*
1563 232-234, 284-299.

1564 Qi, F., Zhu, A.X., 2003. Knowledge discovery from soil maps using inductive learning.
1565 *International Journal of Geographical Information Science* 17(8), 771-795.

1566 Raissi, M., Perdikaris, P., Karniadakis, G.E., 2019. Physics-informed neural networks:
1567 A deep learning framework for solving forward and inverse problems
1568 involving nonlinear partial differential equations. *Journal of Computational*
1569 *Physics* 378, 686-707.

1570 Richards, L.A., 1931. Capillary conduction of liquids through porous mediums.
1571 *Journal of Applied Physics* 1, 318-333.

1572 Richardson, L.F., 1922. *Weather prediction by numerical process*. Cambridge
1573 University Press, Cambridge, United Kingdom.

1574 Roscher, R., Bohn, B., Duarte, M.F., Garcke, J., 2020. Explainable machine learning
1575 for scientific insights and discoveries. *IEEE Access* 8, 42200-42216.

1576 Rosin, N.A., Dematté, J.A.M., Poppiel, R.R., Silvero, N.E.Q., Rodriguez-Albarracin, H.S.,
1577 Rosas, J.T.F., Greschuk, L.T., Bellinaso, H., Minasny, B., Gomez, C. et al., 2023.
1578 Mapping Brazilian soil mineralogy using proximal and remote sensing data.
1579 *Geoderma* 432, 116413.

1580 Rudiyanto, Minasny, B., Chaney, N.W., Maggi, F., Goh Eng Giap, S., Shah, R.M.,
1581 Fiantis, D., Setiawan, B.I., 2021. Pedotransfer functions for estimating soil
1582 hydraulic properties from saturation to dryness. *Geoderma* 403, 115194.

1583 Runje, D., Shankaranarayana, S.M., 2023. Constrained monotonic neural networks,
1584 *Proceedings of the 40th International Conference on Machine Learning*. PMLR,
1585 *Proceedings of Machine Learning Research*, pp. 29338--29353.

1586 Sadeghi, M., Jones, S.B., Philpot, W.D., 2015. A linear physically-based model for
1587 remote sensing of soil moisture using short wave infrared bands. *Remote*
1588 *Sens Environ* 164, 66-76.

1589 Safanelli, J.L., Demattê, J.A.M., Chabrillat, S., Poppiel, R.R., Rizzo, R., Dotto, A.C.,
1590 Silvero, N.E.Q., Mendes, W.d.S., Bonfatti, B.R., Ruiz, L.F.C. et al., 2021.
1591 Leveraging the application of Earth observation data for mapping cropland
1592 soils in Brazil. *Geoderma* 396, 115042.

1593 Seidel, M., Vohland, M., Greenberg, I., Ludwig, B., Ortner, M., Thiele-Bruhn, S.,
1594 Hutengs, C., 2022. Soil moisture effects on predictive VNIR and MIR modeling
1595 of soil organic carbon and clay content. *Geoderma* 427, 116103.

1596 Šimůnek, J., van Genuchten, M.T., Šejna, M., 2016. Recent Developments and
1597 Applications of the HYDRUS Computer Software Packages. *Vadose Zone*
1598 *Journal* 15(7), vzj2016.2004.0033.

1599 Sun, X.L., Minasny, B., Wang, H.L., Zhao, Y.G., Zhang, G.L., Wu, Y.J., 2021.
1600 Spatiotemporal modelling of soil organic matter changes in Jiangsu, China
1601 between 1980 and 2006 using INLA-SPDE. *Geoderma* 384, 114808.

1602 Taghizadeh-Mehrjardi, R., Mahdianpari, M., Mohammadimanesh, F., Behrens, T.,
1603 Toomanian, N., Scholten, T., Schmidt, K., 2020. Multi-task convolutional neural
1604 networks outperformed random forest for mapping soil particle size fractions
1605 in central Iran. *Geoderma* 376, 114552.

1606 Talebi, H., Peeters, L.J.M., Otto, A., Tolosana-Delgado, R., 2022. A truly spatial
1607 random forests algorithm for geoscience data analysis and modelling.
1608 *Mathematical Geosciences* 54(1), 1-22.

1609 Tang, J., Riley, W.J., Manzoni, S., Maggi, F., 2024. Feasibility of formulating ecosystem
1610 biogeochemical models from established physical rules. *Journal of*
1611 *Geophysical Research: Biogeosciences* 129(6), e2023JG007674.

1612 Tartakovsky, A.M., Marrero, C.O., Perdikaris, P., Tartakovsky, G.D., Barajas-Solano, D.,
1613 2020. Physics-informed deep neural networks for learning parameters and
1614 constitutive relationships in subsurface flow problems. *Water Resources*
1615 *Research* 56(5), e2019WR026731.

1616 Tziolas, N., Tsakiridis, N., Ogen, Y., Kalopesa, E., Ben-Dor, E., Theocharis, J., Zalidis,
1617 G., 2020. An integrated methodology using open soil spectral libraries and
1618 Earth Observation data for soil organic carbon estimations in support of soil-
1619 related SDGs. *Remote Sens Environ* 244, 111793.

1620 van Genuchten, M.T., 1980. A closed-form equation for predicting the hydraulic
1621 conductivity of unsaturated soils. *Soil Science Society of America Journal*
1622 44(5), 892-898.

1623 van Zijl, G., van Tol, J., Tinnefeld, M., Le Roux, P., 2019. A hillslope based digital soil
1624 mapping approach, for hydrogeological assessments. *Geoderma* 354,
1625 113888.

1626 Vohland, M., Ludwig, B., Seidel, M., Hutengs, C., 2022. Quantification of soil organic
1627 carbon at regional scale: Benefits of fusing vis-NIR and MIR diffuse reflectance
1628 data are greater for in situ than for laboratory-based modelling approaches.
1629 *Geoderma* 405, 115426.

1630 von Rueden, L., Mayer, S., Beckh, K., Georgiev, B., Giesselbach, S., Heese, R., Kirsch,
1631 B., Pfrommer, J., Pick, A., Ramamurthy, R. et al., 2023. Informed machine
1632 learning - A taxonomy and survey of integrating prior knowledge into learning
1633 systems. *IEEE Transactions on Knowledge and Data Engineering* 35(1), 614-
1634 633.

1635 Wadoux, A.M.J.C., Minasny, B., McBratney, A.B., 2020. Machine learning for digital
1636 soil mapping: Applications, challenges and suggested solutions. *Earth-Science*
1637 *Reviews* 210, 103359.

- 1638 Wadoux, A.M.J.C., Molnar, C., 2022. Beyond prediction: methods for interpreting
1639 complex models of soil variation. *Geoderma* 422, 115953.
- 1640 Wang, J., Filippi, P., Haan, S., Pozza, L., Whelan, B., Bishop, T.F.A., 2024. Gaussian
1641 process regression for three-dimensional soil mapping over multiple spatial
1642 supports. *Geoderma* 446, 116899.
- 1643 Wang, S., Sankaran, S., Wang, H., Perdikaris, P., 2023a. An Expert's Guide to Training
1644 Physics-informed Neural Networks.
- 1645 Wang, Y., Jin, M., Deng, Z., 2018. Alternative model for predicting soil hydraulic
1646 conductivity over the complete moisture range. *Water Resources Research*
1647 54(9), 6860-6876.
- 1648 Wang, Y., Shi, L., Hu, X., Song, W., Wang, L., 2023b. Multiphysics-Informed Neural
1649 Networks for Coupled Soil Hydrothermal Modeling. *Water Resources Research*
1650 59(1), e2022WR031960.
- 1651 Weber, T.K.D., Finkel, M., da Conceição Gonçalves, M., Vereecken, H.,
1652 Diamantopoulos, E., 2020. Pedotransfer function for the Brunswick soil
1653 hydraulic property model and comparison to the van Genuchten-mualem
1654 model. *Water Resources Research* 56(9), e2019WR026820.
- 1655 Weber, T.K.D., Weihermüller, L., Nemes, A., Bechtold, M., Degré, A., Diamantopoulos,
1656 E., Fatichi, S., Filipović, V., Gupta, S., Hohenbrink, T.L. et al., 2024. Hydro-
1657 pedotransfer functions: a roadmap for future development. *Hydrol. Earth*
1658 *Syst. Sci.* 28(14), 3391-3433.
- 1659 Weindorf, D.C., Chakraborty, S., 2024. Balancing machine learning and artificial
1660 intelligence in soil science with human perspective and experience.
1661 *Pedosphere* 34(1), 9-12.
- 1662 Wessolek, G., Bohne, K., Trinks, S., 2023. Validation of soil thermal conductivity
1663 models. *International Journal of Thermophysics* 44(2), 20.
- 1664 Widyastuti, M.T., Minasny, B., Padarian, J., Maggi, F., Aitkenhead, M., Beucher, A.,
1665 Connolly, J., Fiantis, D., Kidd, D., Ma, Y. et al., 2024. PEATGRIDS: Mapping
1666 thickness and carbon stock of global peatlands via digital soil mapping. *Earth*
1667 *Syst. Sci. Data Discuss.* 2024, 1-29.
- 1668 Willard, J.D., Jia, X., Xu, S., Steinbach, M.S., Kumar, V., 2020. Integrating physics-
1669 based modeling with machine learning: A survey. *ArXiv abs/2003.04919*.
- 1670 Wu, F., Tan, K., Wang, X., Ding, J., Liu, Z., 2023. A novel semi-empirical soil multi-
1671 factor radiative transfer model for soil organic matter estimation based on
1672 hyperspectral imagery. *Geoderma* 437, 116605.
- 1673 Xie, E., Zhang, X., Lu, F., Peng, Y., Chen, J., Zhao, Y., 2022. Integration of a process-
1674 based model into the digital soil mapping improves the space-time soil
1675 organic carbon modelling in intensively human-impacted area. *Geoderma*
1676 409, 115599.
- 1677 Xie, X., Yan, H., Lu, Y., Zeng, L., 2024. Simulating field soil temperature variations
1678 with physics-informed neural networks. *Soil and Tillage Research* 244,
1679 106236.
- 1680 Yang, R.M., Liu, L.A., Zhang, X., He, R.X., Zhu, C.M., Zhang, Z.Q., Li, J.G., 2022. The
1681 effectiveness of digital soil mapping with temporal variables in modeling soil
1682 organic carbon changes. *Geoderma* 405, 115407.
- 1683 Yang, W.-H., Clifford, D., Minasny, B., 2015. Mapping soil water retention curves via
1684 spatial Bayesian hierarchical models. *Journal of Hydrology* 524, 768-779.
- 1685 Zaman, B., McKee, M., Neale, C.M.U., 2012. Fusion of remotely sensed data for soil
1686 moisture estimation using relevance vector and support vector machines.
1687 *International Journal of Remote Sensing* 33(20), 6516-6552.

1688 Zhang, L., Heuvelink, G.B.M., Mulder, V.L., Chen, S., Deng, X., Yang, L., 2024. Using
1689 process-oriented model output to enhance machine learning-based soil
1690 organic carbon prediction in space and time. *Science of The Total*
1691 *Environment* 922, 170778.

1692 Zhang, N., Wang, Z., 2017. Review of soil thermal conductivity and predictive
1693 models. *International Journal of Thermal Sciences* 117, 172-183.

1694 Zhang, X., Xie, E., Chen, J., Peng, Y., Yan, G., Zhao, Y., 2023. Modelling the
1695 spatiotemporal dynamics of cropland soil organic carbon by integrating
1696 process-based models differing in structures with machine learning. *Journal of*
1697 *Soils and Sediments* 23(7), 2816-2831.

1698 Zhang, Y., Schaap, M.G., 2017. Weighted recalibration of the Rosetta pedotransfer
1699 model with improved estimates of hydraulic parameter distributions and
1700 summary statistics (Rosetta3). *Journal of Hydrology* 547, 39-53.

1701