



33 collated about 20,000 soil profile data and stored them in a central database. The data were cleaned
34 and harmonized using the latest soil profile data template and prepared 14,681 profile data for
35 modelling. Random Forest was used to develop a continuous quantitative digital map of 18 WRB
36 reference soil groups at 250 m resolution by integrating environmental variables-covariates
37 representing major Ethiopian soil-forming factors. The validated map will have tremendous
38 significance in soil management and other land-based development planning, given its improved
39 spatial nature and quantitative digital representation.

40 **Keywords:** soil profiles, environmental covariates, modelling, expert validation, Reference Soil
41 Group

42 **1 Introduction**

43 Soils are important resources that support the development and production of various economic,
44 social, and ecosystem services, and are useful in climate change mitigation and adaptation (Baveye
45 et al., 2016). Data on soils' physical and chemical characteristics and spatial distribution are needed
46 to define and plan their functions over time and space, which is an important step toward the
47 sustainable use and management of soils (Elias, 2016; Hengl et al., 2017).

48 In Ethiopia, soil surveys and mapping have been conducted at various scales with varying scope,
49 approach, methodology, quality, and level of detail (Abayneh, 2001; Abayneh and Berhanu, 2007;
50 Berhanu, 1994; Elias, 2016; Zewdie, 2013). The most recent country-wide digital soil mapping
51 efforts focused primarily on soil characteristics (Ali et al., 2020; Iticha and Chalsissa, 2019; Tamene
52 et al., 2017), although soil class maps are equally important for allocating a particular soil unit for
53 specific use (Leenaars et al., 2020a; Wadoux et al., 2020). Many notable attempts have been made to
54 improve digital soil information system (Hengl et al., 2021, 2017; 2015; Poggio et al., 2020).
55 However, such initiatives were based on limited and unevenly distributed soil profile data (e.g., 1.15
56 soil profiles per 1,000 km² for Ethiopia) which limits the accuracy and applicability of the products.

57 Thousands of soil profile data were collected since the 1960s (Erkossa et al., 2022), but these data
58 were hardly accessible as they were scattered across different institutions and individuals (Ali et al.,
59 2020). Furthermore, country-wide quantitative and grided spatial soil type information is hardly
60 available (Elias, 2016). The Ethiopian Soil Information System (EthioSIS) project attempted to
61 develop a countrywide digital soil map focusing on topsoil characteristics, including plant nutrient



62 content, but overlooked soil resource mapping (Ali et al., 2020; Elias, 2016), despite a strong need
63 for a high-resolution soil resource map (Mulualem et al., 2018).

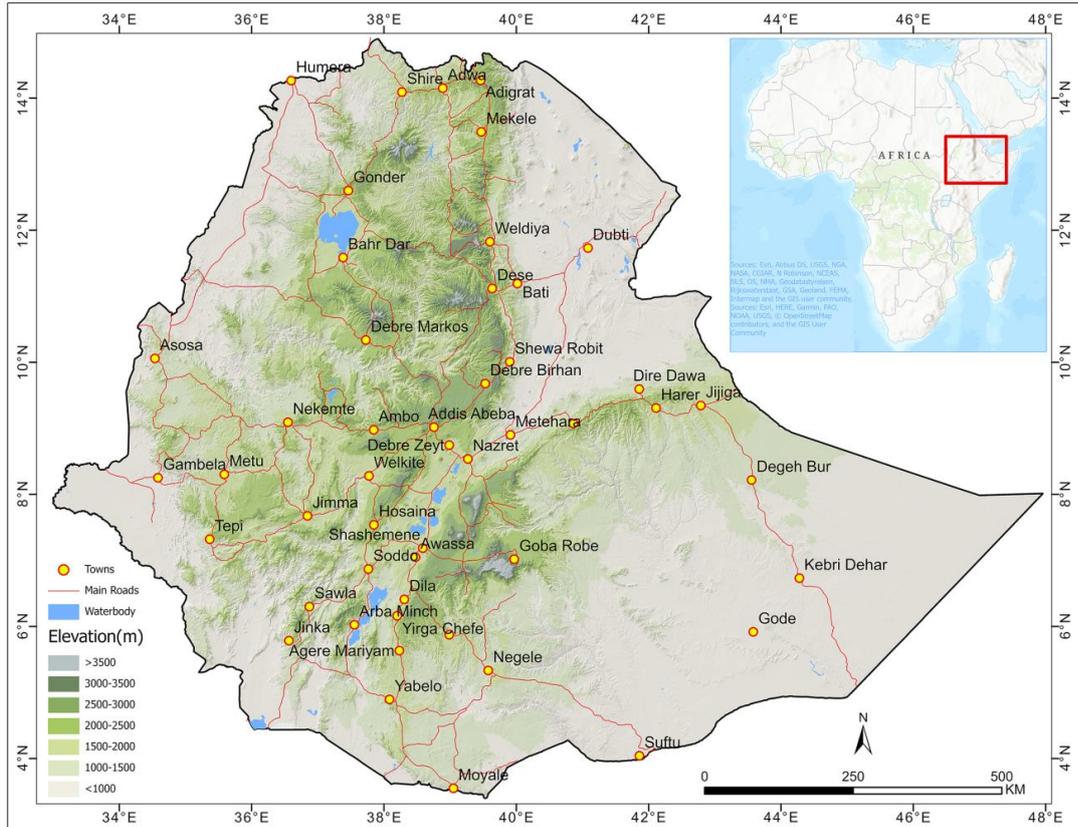
64 Ethiopia has an area of about 1.14 mill. km² consisting of varied environments, making its soils
65 extremely heterogeneous; thus capturing heterogeneity using conventional soil survey and mapping
66 approaches is a resource- and time-consuming endeavour (Hounkpatin et al., 2018). This can be
67 circumvented using available legacy soil profile data accumulated over time coupled with advanced
68 analytical techniques to develop high-resolution digital soil maps (Hounkpatin et al., 2018; Kempen,
69 2012, 2009).

70 The objectives of this study were to (1) develop a national legacy soil profile dataset that can be used
71 as an input for various digital soil mapping exercises, and (2) generate an improved 250 m digital
72 International Union of Soil Science (IUSS) World Reference Base (WRB) Reference Soil Groups
73 (RSGs) map of Ethiopia using the legacy soil profile dataset and advanced machine learning
74 techniques.

75 **2 Methods**

76 **2.1 The study area**

77 The study area covered the entire area of Ethiopia (1.14 mill. km²) located between 3°N and 15° N,
78 and between 33° E and 48° E (Figure 1). The topography of the country is marked by a large
79 altitudinal variation, ranging from 126 meters below sea level at Dalol to 4,620 m at Ras Dashen
80 Mountain in the northwest part of the highlands (Billi, 2015; Enyew and Steeneveld, 2014). The
81 country embraces diverse agroecological zones and farming systems. Ethiopia's wide range of
82 topography, climate, parent material, and land use types created conditions for the formation of
83 different soil types (Abayneh., 2005; Donahue, 1962; Mesfin, 1998; Zewdie, 2013, 1999). More than
84 33% of the country is covered by the central upper and highland complex (Abegaz et al., 2022),
85 which embraces Africa's most prominent mountain system, reaching a maximum altitude of 4,620 m
86 above sea level (Hurni, 1998).



87

88 **Figure 1.** Location map of Ethiopia, overview map © Esri World Topographic Map.

89 **2.2 Legacy soil profile data collation and preparation**

90 In Ethiopia, soil profile data have been generated over decades through various soil survey missions
91 but kept in a variety of formats and quality with limited accessibility. There has been no institution
92 with a national mandate to coordinate the generation, collation, harmonization, and sharing of soil
93 profile data. This has led to the formation of the Coalition of the Willing (CoW) in 2018—a group of
94 individuals and institutions willing to exchange soil and agronomy data to overcome the challenges
95 posed by the lack of data access and sharing mechanism in the country (Tamene et al., 2021).

96 The CoW conducted a national soil and agronomy data ecosystem mapping which revealed that a
97 plethora of legacy soil resource data sets do exist but are scattered across different institutions and

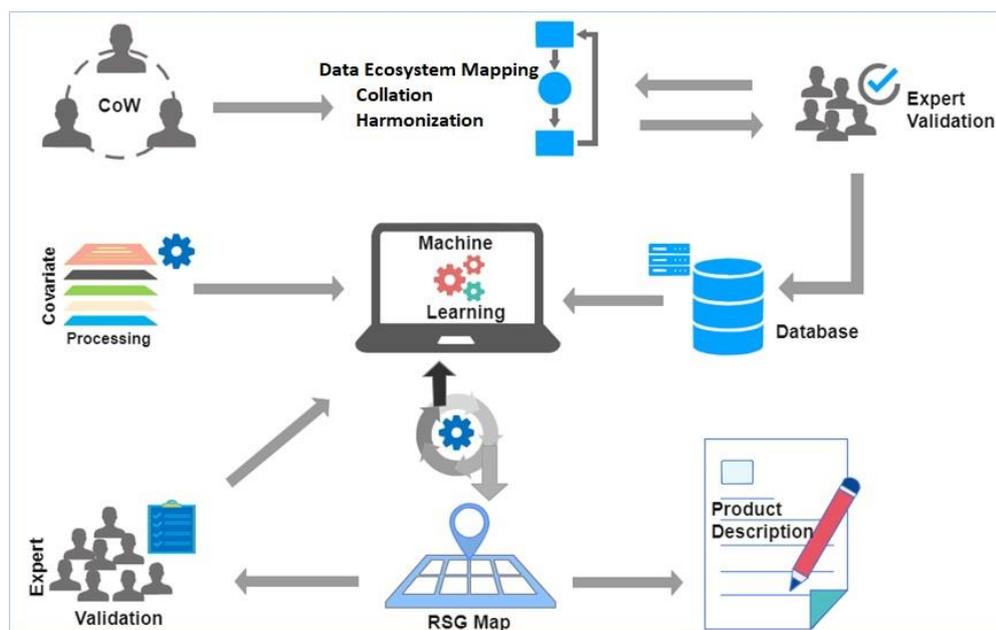


98 individuals (Ali et al., 2020). The assessment also revealed that a sizable proportion of the data
99 holders were willing to share the data in their custody, provided that some regulations are put in
100 place to administer the data. The CoW supported and facilitated data collation campaigns, which
101 involved both formal and informal approaches to data holders.

102 Soil profile data collected from the 1970s to 2021 were acquired from over 88 diverse sources
103 through a data collation campaign (Tamene et al., 2022). Initially, 8000 profile data points were
104 collated and subjected to improved modelling techniques to create a provisional WRB reference soil
105 group map of Ethiopia. This was presented for various partners and data holding institutions to
106 demonstrate the power of data sharing. This created awareness and enabled to mobilise and collate
107 over 20,000 legacy soil profile data. These data were then added to the national data repository.

108 The data had varying levels of completeness in terms of soil field and environmental descriptions
109 and laboratory analysis. This required a rigorous expert-based quality assessment and
110 standardization before compilation into a harmonized format. The expanded version of the Africa
111 Soil Profile (AfSP) database (Leenaars et al., 2014) template was used for standardizing and
112 harmonizing the data. Out of the collated soil profile data, 14,681 georeferenced data points were
113 extracted based on completeness and cleanness for the purposes of modelling. The cleaned soil
114 profile data set contains at least the reference soil group (RSG) nomenclature as outlined in the WRB
115 legend. While the original soil profile records were set in different coordinate systems, all were
116 projected into the adopted standard georeferencing system, namely WGS84, decimal degrees in the
117 QGIS (3.20.2) environment (QGIS Development Team, 2021). To verify their position, soil profile
118 locations were plotted using a standard WGS84 coordinate system to verify that points are matching
119 with the site description, geomorphological settings, and at the very least the source project
120 boundary outline.

121 The accuracy of the data depends on the quality and reliability of the survey data itself which in turn
122 requires expert knowledge and experience in soil description (Leenaars et al., 2020a). In this study,
123 data cleaning, validation, reclassification, and verification were carried out by a team of prominent
124 national pedologists and soil surveyors, including those involved in the generation of some of the
125 soil profile data themselves (Figure 2).



126

127 **Figure 2.** Schematic presentation of data acquisition and workflow.

128 In addition, the Ministry of Agriculture (MoA) soil survey and mapping experts and other volunteers
129 have validated the legacy soil profile observations. This led to the reclassification of the soil types as
130 deemed necessary. Such validation and reclassification involved re-examining the geomorphological
131 setup of the soil profile locations using Google earth as well as reviewing the site and soil
132 description and the corresponding laboratory data and reviewing the proposed soil type. The
133 harmonised data sets in the database were used as input soil profile data for modelling and mapping
134 IUSS WRB reference soil groups.

135 **2.3 Selection and pre-processing of covariates**

136 In order to develop spatially continuous soil class/type maps, data on environmental covariates that
137 represent directly or indirectly the soil-forming factors have to be integrated with soil profile data
138 (Hengl and MacMillan, 2019). Environmental covariates representing soil-forming factors (climate,
139 organisms, relief, parent material, and time) were derived from diverse remote sensing products and
140 thematic maps (Hengl and MacMillan, 2019; McBratney et al., 2003). Selected environmental
141 covariate layers were then used to predict the soil property across the full extent of the prediction



142 area using the soil observation data from the sampling locations (McBratney et al., 2003, Miller et
143 al., 2021).

144 In this study, a set of 27 covariate layers (Appendix B), from 68 potential covariates, were prepared
145 in GeoTiff format with 250 m resolution and Lambert azimuthal equal-area projection with the
146 latitude of origin 8.65 and centre of meridian 39.64 which is the centre point for Ethiopia. This
147 projection was selected since it is effective in minimizing area distortions over land. All layers were
148 masked for buildings and water bodies by the national boundary of Ethiopia and stacked using the
149 stack () function of the raster package in R [version 4.05] (R Core Team, 2020). A 250 m spatial
150 resolution was chosen to accommodate both the spatial resolution of the major co-variate inputs and
151 make it applicable for large-scale analysis.

152 The covariates included terrain variables derived from the 90-meter Shuttle Radar Topography
153 Mission (SRTM) digital elevation model (DEM) (Vågen, 2010), climatic variables from Enhancing
154 National Climate Services (ENACTS) (Dinku et al., 2014), Moderate Resolution Imaging
155 Spectroradiometer (MODIS) imagery raw bands and derived indices (Vågen, 2010), national
156 geology map of Ethiopia (Tefera et al., 1996), and land use/ cover map of Ethiopia (WLRC-AAU,
157 2010) (Table 1).

158 A 4 km climate grid data from the National Meteorological Agency's (NMA) ENACTS initiative
159 was used because it addresses the spatial and temporal gaps and quality problems of other climatic
160 data sources for Ethiopia (Dinku et al., 2014). The long-term mean, minimum, maximum, and
161 standard deviation temperature, and precipitation data for the period between 1983 and 2016 from
162 the ENACT-NMA initiatives (Dinku et al., 2014) were used. In addition, the hydrologically
163 corrected DEM of the Africa soil information service (Vågen, 2010) and DEM derivatives were
164 calculated using SAGA-GIS version 7.3.0 (Conrad et al., 2015) for topography as a soil-forming
165 factor. We used national geological (Tefera et al., 1996) and land use/land cover (WLRC-AAU,
166 2010) thematic maps of Ethiopia to represent parent material and organisms, respectively.

167 The covariate pre-processing, visual inspection for inconsistencies, resampling to a target grid of 250
168 m and compilations were conducted in QGIS [3.20.2] (QGIS Development Team, 2021), SAGA GIS
169 [7.8.2] (Conrad et al., 2015) and R [version 4.05] (R Core Team, 2020) software packages. Once
170 each covariate was adjusted to have an identical spatial resolution, extent and projection, continuous



171 covariates were resampled using the bilinear spline method whereas categorical covariates were
172 resampled using the nearest neighbour method.

173 The near-zero variance, available in the near *ZeroVar* function *caret* package in R (Kuhn, 2008) was
174 used to identify and remove environmental variables that have little or no variance. After expert
175 judgement to determine the type of covariates for modelling RSGs and near-zero variance analysis, a
176 total of 27 environmental variables (24 continuous and 3 categorical) were used for the modelling.

177 **2.4 Modelling and mapping soil types/reference soil groups**

178 **2.4.1 Model tuning and quantitative evaluation**

179 Recent developments in data analytics showed the potential to undertake sophisticated analysis
180 involving large datasets within a relatively short time using models. In digital soil mapping,
181 machine-learning techniques have been extensively used to determine the relationship between soil
182 types and environmental variables (McBratney et al., 2003). Many machine learning models were
183 developed in the past decades for digital soil mapping to spatially predict soil classes based on
184 existing soil data and soil-forming environmental covariates (Heung et al., 2016). Random Forest
185 (RF), a tree-based ensemble method, is one of the most promising machine learning techniques
186 available for digital soil mapping (Breiman, 2001; Heung et al., 2016), which has gained tremendous
187 popularity due to its high overall accuracy and has been widely used in predictive soil mapping
188 (Brungard, 2015; Hengl et al., 2018).

189 Examples of the main strengths of the RF model are its ability to handle numerical and categorical
190 data without any assumption of the probability distribution; and its robustness against nonlinearity
191 and overfitting (Breiman, 2001; Svetnik et al., 2003). In the RF model, data are split into training (80
192 %) and testing (20 %) components for building the model and model testing, respectively (Kuhn,
193 2008).

194 Hyper-parameter optimization and cross-validation on the training dataset have been performed for
195 optimal model application using *Caret* package (Kuhn, 2008). Model tuning was performed with a
196 repeated 10-fold cross-validation procedure and applied multiple combinations of hyper-parameters
197 for the ranger method, which is a fast implementation of RF, particularly suited for high dimensional
198 data (Wright and Ziegler, 2017). Three parameters, i.e., the number of covariates used for the splits
199 (*mtry*), splitting rules (*splitrule*) and minimum node size (*min.node.size*) were optimised. The values



200 of 1,000 number of trees (ntree) with mtry ranged from 10 to 20, min.node.size ranged from 5 – 15
201 with an interval of five and extra trees as splitrule fed for the optimization procedure.

202 The accuracy of the testing dataset was related to the model performance for the new dataset,
203 indicating the capacity of the model to predict at the unsampled location. A confusion matrix was
204 also used to calculate a cross-tabulation of observed and predicted classes with associated statistics
205 i.e., producer's accuracy and user's accuracy. The computational framework was based on open-
206 source software and was performed on a windows server 2016 standard with 250 GB of working
207 memory.

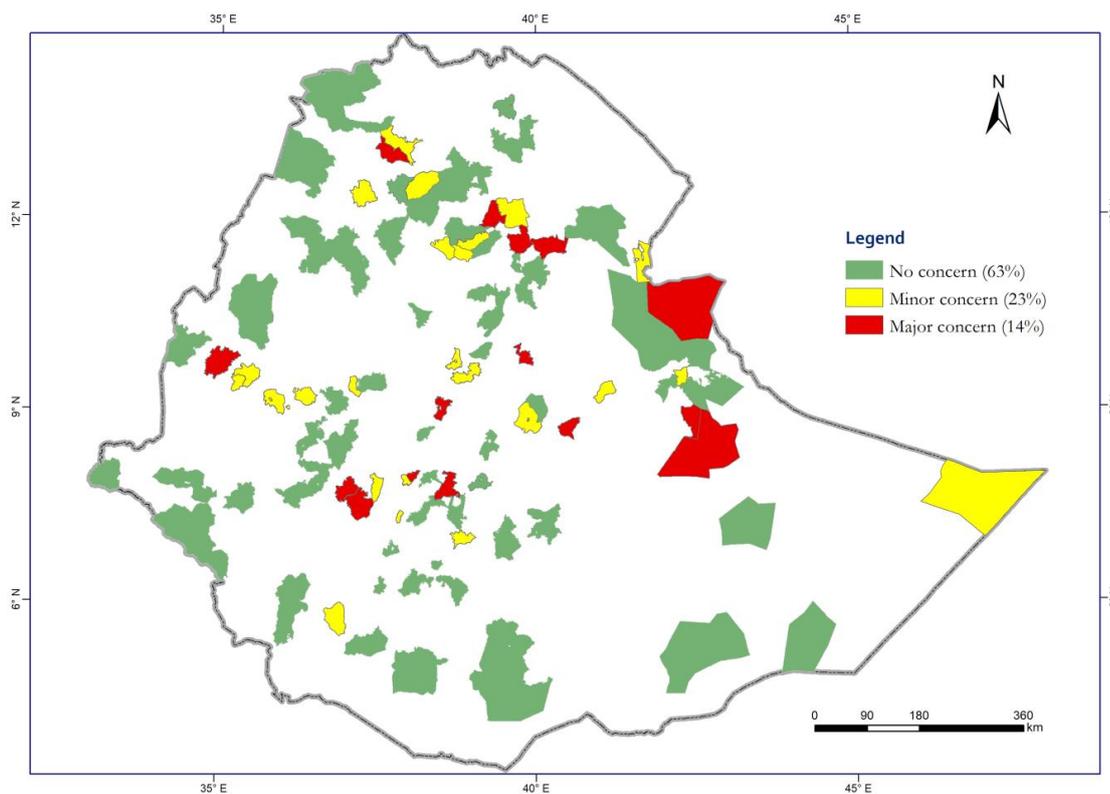
208 **2.4.2 Qualitative evaluation of spatial patterns of the beta-version soil map**

209 Expert knowledge of soil-landscape relations and soil distribution remains important to evaluate the
210 predictive soil mapping results and assess if predicted spatial patterns make sense from a pedological
211 viewpoint (Hengl et al., 2017; Poggio et al., 2020). An important step in model evaluation is,
212 therefore, expert assessment whereby professionals with broad experience in soil survey and
213 mapping can evaluate and improve the quality of the soil resource map. Accordingly, an expert
214 validation workshop was conducted using the first version of the reference soil groups (RSGs) map.
215 About 45 multi-disciplinary scientists including soil surveyors, pedologists, geologists, and
216 geomorphologists were drawn from national and international research, development, and higher
217 learning institutions to review the draft RSG map in plenary. This was followed by breakout sessions
218 where groups of experts evaluated the map based on their experience and knowledge of soil-
219 landscape relations of the country.

220 While the plenary discussion provided an overview of the approaches followed in developing the
221 map, the facilitated group discussion helped to have an in-depth review of the selected polygons of
222 the map assigned to them. Participants were split into five groups (with 8-10 members each) and
223 have chosen up to 60 polygons representing areas with which at least one of the group members has
224 sufficient information, including data sources. Overall, the groups have checked a total of 126
225 polygons (Figure 3) which were fairly distributed across the country. In cases where there is
226 ambiguity, the experts overlaid the soil profile locations on Google earth map to evaluate the
227 description and soil lab results. The group members displayed the polygons one by one in a GIS
228 environment and discussed the predicted dominant and associated soil reference soil groups and



229 labelled them in one of three confirmation categories: 1. confirmed with ‘no concern’, 2. confirmed with
230 with “minor concern”, and 3. confirmed with ‘major concern’. Confirmation with ‘no concern’ was
231 made when all members of a group agreed on both the types and relative coverage of the predicted
232 soils within the polygon. Confirmation with ‘minor concern’ was made when all or some of the
233 team members agreed on the predicted soil types within the polygons but did not agree on the order
234 of abundance or the probability occurrence of one or two soils, while confirmation with ‘major
235 concern’ was made when all members of the team did not agree on the predicted soil type, or when
236 the presence of another soil type, other than the predicted ones is noted.



237
238 **Figure 3.** The spatial distribution of districts validated by stakeholders and feedback categories
239 according to the level of concerns raised.

240 After finalising the evaluation at the group’s level assessment, each group presented the results in the
241 plenary followed by a discussion to get feedback from other participants. Following the plenary



242 discussions, the participants created a group of six senior pedologists to work on the
243 recommendations, including validation of the additional data obtained during the event. Based on
244 these outputs, the model was re-run to produce the current version of the soil map.

245 **3 Results and Discussion**

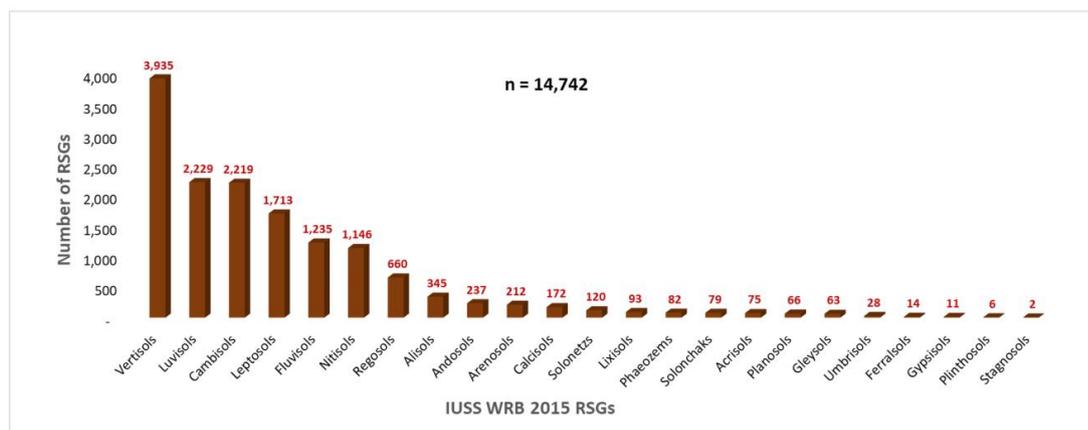
246 **3.1 Soil profile datasets**

247 Using the IUSS WRB, 2015, the preliminary identified 14,742 georeferenced legacy soil profiles
248 were classified/reclassified into twenty-three reference soil groups (RSGs). Nearly 90% of the soil
249 profile points represented Vertisols, followed by Luvisols, Cambisols, Leptosols, Fluvisols, and
250 Nitisols, which were found to be the dominant soil types in Ethiopia (Figure 4). The remaining 10%
251 represented the Regosols, Alisols, Andosols, Arenosols, Calcisols, Solonetz, Lixisols, Phaeozems,
252 Solonchaks, Acrisols, Planosols, Gleysols, Umbrisols, Ferralsols, Gypsisols, Plinthosols, and
253 Stagnosols.

254 The results suggest that about 72 % of the IUSS WRB (2015) RSGs were confirmed to occur in
255 Ethiopia. In this regard, Ethiopia is considered as a soil museum having endowed with a diverse
256 range of soil types owing to the diversities in the pedogenetic factors (Elias, 2016), which is known
257 to have most of the reference soil groups in varying frequencies depending on existing physiographic
258 and agroecological positions (Mishra et al., 2004).

259 One of the challenges with legacy soil data in categorical mapping is that of imbalanced soil
260 samples, in that all classes were not represented equally (Wadoux et al., 2020). For this study, soil
261 profiles with less than 30 observations were objectively excluded from the model after examining
262 the accuracy and spatial distribution of each reference soil group. Five reference soil groups
263 (Umbrisols, Ferralsols, Gypsisols, Plinthosols, and Stagnosols) were excluded from the model and
264 left unmapped in this EthioSoilGrid version 1.0 map.

265



266

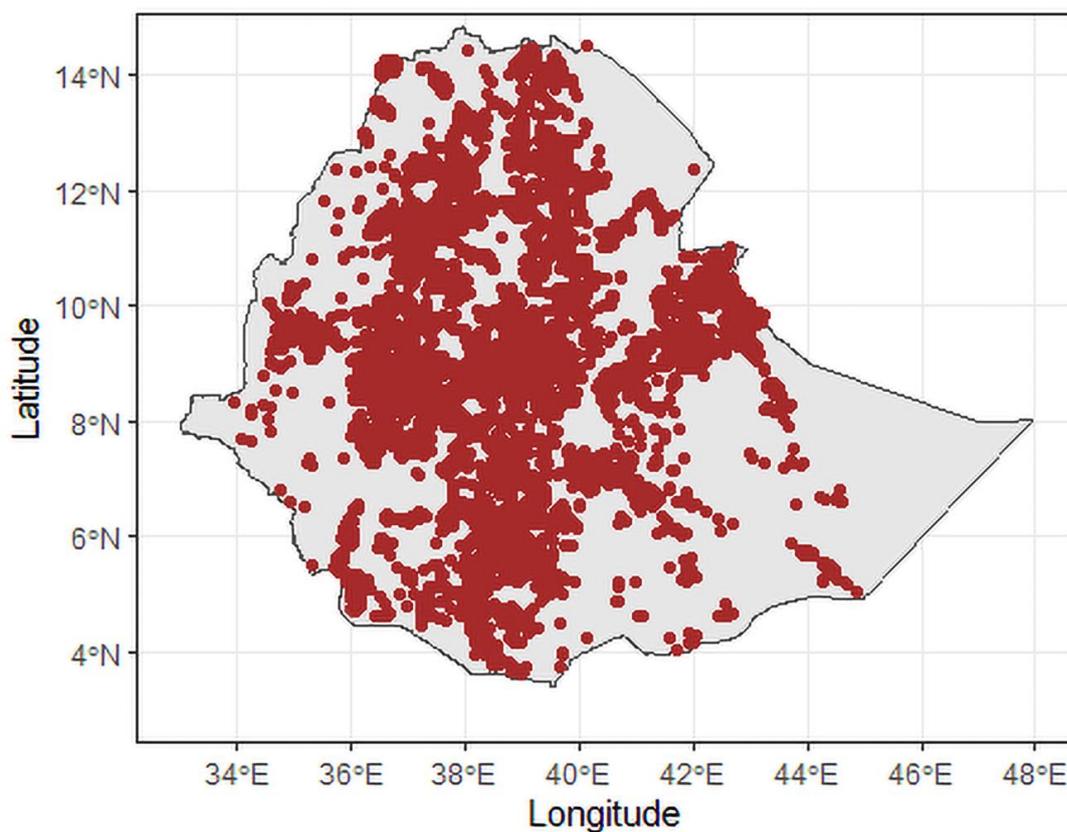
267 **Figure 4.** Number of profile points per WRB reference soil groups.

268 With regards to the total area of Ethiopia and excluding the built-up (urban) and water surface areas,
 269 and, the soil profile spatial distribution (Figure 5) represented an average density of 13.1 soil profiles
 270 per 1,000 km². The actual density of observations varied greatly between different parts of the
 271 country. The variation tends to follow river basins, sub-basins, and agricultural land-use types-based
 272 studies from which legacy soil observations were pulled for the present study. For instance, in 30
 273 intervention districts of the Capacity Building for Scaling up of Evidence-Based Best Practices in
 274 Agricultural Production in Ethiopia (CASCAPE) project, the average profile density was 1 profile
 275 per 11.5 km² (about 87 profiles per 1,000km²) for a total area of about 26,830 km² (Leenars et
 276 al.,2020a). Similar semi-detailed soil mapping missions in 15 districts were conducted through the
 277 Bilateral Ethiopia-Netherlands Effort for Food, Income and Trade (BENEFIT)-REALISE project
 278 which generated about 217 observations per 1,000 km² (Leenars et al., 2020b).

279 A soil type and depth map compilation and updating mission at a 1:250,000 scale by the Water Land
 280 Resource Centre (WLRC) of Addis Ababa University collated and used about 3,949 legacy soil
 281 profiles for the entire country (Ali et al., 2020), about 3.5 profiles per 1,000 km². The existing
 282 accessible compiled legacy soil profile database of Ethiopia prepared by the Africa soil profile
 283 database consisted of 1,712 legacy soil profile observations or 1.5 profiles per 1,000 km² (Batjas et
 284 al., 2020; Leenaars et al., 2014), which indicates that the number of data used in this study is 8.5
 285 times higher than that was used in the former. However, the soil profile distribution across the



286 country was uneven; additional soil survey missions are needed for the eastern lowlands and other
287 less represented areas in the future.



288

289 **Figure 5.** Spatial distribution of collated legacy soil profile data.

290 The soil profiles distribution across the 32 agro-ecological zones (AEZ) of Ethiopia revealed that all,
291 except two—tepid per-humid mid highland (0.13% landmass) and very cold sub-humid sub-afro
292 alpine to afro-alpine (0.03 % landmass)—were represented by soil profiles observations. Furthermore,
293 about 95 % of the profile observations represented 91 % of the AEZs aerial coverage (Appendix A).
294 The distribution of legacy soil profiles varied across AEZs. In general, top-ranked lowland AEZs
295 with roughly 56 % area coverage obtained 23 % of the total profile observations, while top-ranked
296 highland AEZs with 20 % area coverage received 47 % of profile observations. For instance, warm
297 desert, warm moist, hot arid, and warm sub-moist lowlands with area coverage of around 20 %, 15



298 %, 11 %, and 10 %, were represented roughly by 3 %, 11 %, 2 %, and 7 % of the total profiles,
299 respectively. Tepid moist mid highlands (8% area coverage), tepid sub-humid mid highlands (7 %
300 area coverage), and tepid sub-moist mid highlands (5 % area coverage) each were represented by 20
301 %, 15 %, and 12 % of the profiles, respectively.

302 **3.2 Modelling and Mapping**

303 **3.2.1 Variable importance**

304 The reference soil group spatial pattern is primarily influenced by long-term average surface
305 reflectance, flow-based DEM indices, and precipitation. Figure 6 shows variables of importance for
306 determining RSGs spatial prediction. The top-ranked variables were (i) long-term MODIS Near-
307 Infrared (NIR) reflectance; (ii) multiresolution index of valley bottom flatness, (iii) long-term mean
308 day-land surface temperature; (iv) long-term mean soil moisture; (v) standard deviation of long-term
309 precipitation; (vi) long-term mean precipitation; and (vii) topographic wetness index.

310 MODIS long-term mean spectral signatures showed high relative importance. According to Hengl et
311 al (2017), accounting for seasonal vegetation fluctuation and inter-annual variations in surface
312 reflectance, long-term temporal signatures of the soil surface, derived as monthly averages from
313 long-term MODIS imagery were more effective. Furthermore, Hengl and MacMillan (2019)
314 explained that long-term average seasonal signatures of surface reflectance provide a better
315 indication of soil characteristics than only a single snapshot of surface reflectance.

316 The Multi-Resolution Valley Bottom Flatness Index, a DEM-derived topography index, is the
317 second top-ranked covariate driving soil variability across Ethiopia. This hydrological/soil removal
318 and accumulation/deposition index is used to distinguish valley floor and ridgetop landscape
319 positions (Soil Science Division Staff, 2017) highly responsible for multiple soil-forming processes
320 to operate over a particular landscape, resulting in a wide range of soil development. The influence
321 of topography on spatial soil variation is manifested in every landscape of Ethiopia (Belay, 1997;
322 Mesfin, 1998; Zewdie, 2013).

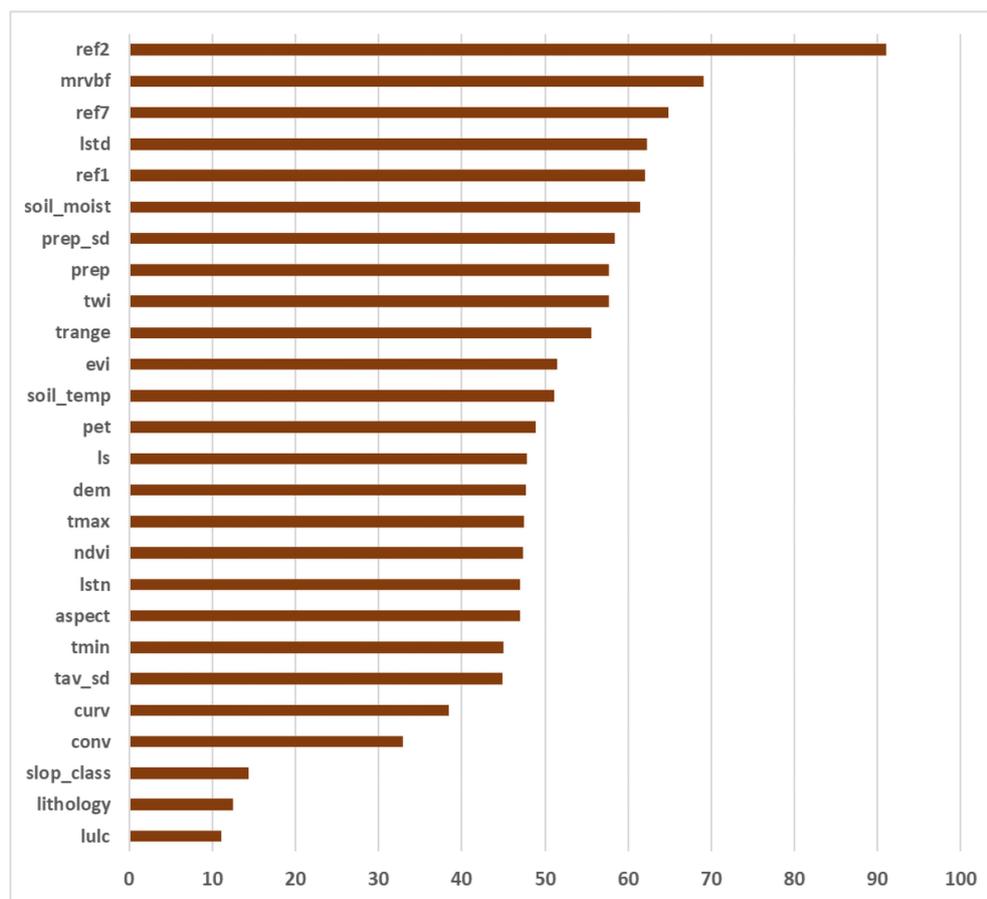
323 Long-term daily mean land surface temperature, mean soil moisture, rainfall standard deviation and
324 mean annual rainfall were among the top-ranked covariates for predicting reference soil groups'
325 spatial variation across the country. In Ethiopia, different soil genesis studies revealed that climate



326 has a significant influence on soil development and properties and is, therefore, responsible for
327 having widely varying soils in the country (Abayneh, 2006, 2005; Fikru, 1988, 1980; Zewdie, 2013).
328 Rainfall variability in Ethiopia is governed by global, regional, and local factors. Ethiopian climate
329 is substantially governed by local factors in which the topography is powerful. It is known as a
330 country of natural contrast; characterised by a complex topography that strongly defines both rainfall
331 and temperature patterns, by modifying the influence of the large-scale ocean-land-atmosphere
332 pattern, thus creating diverse localised climates.

333 Spatially, rainfall in Ethiopia is characterised by a decreasing trend in the direction from west to east,
334 south-north, west-north-east and west-east. The lowlands in the southeast and northeast, covering
335 approximately 55% of the country's land area, are under arid and semi-arid climates. Annual rainfall
336 ranges from less than 300 mm in the south-eastern and north-western lowlands to over 2,000 mm in
337 the southwestern (southern portion of the western highlands). The eastern lowlands get rain twice a
338 year, in April–May and October–November, with two dry periods in between. The total annual
339 precipitation in this regime varies from 500 to 1,000 mm. The driest of all regions is the Denakil
340 Plain, which receives less than 500 mm and sometimes none (Fazzini et al., 2015). Temperatures are
341 also greatly influenced by the rapidly changing altitude in Ethiopia and mean monthly values vary
342 from about 35°C, in the northeast lowland to less than 7.5°C over the north and central highland.

343 Among the most important covariates for predicting reference soil groups in the Ethiopian highlands,
344 (Leenars et al., 2020a), are monthly average soil moisture for January (ranked 3rd), long-term
345 average soil moisture (ranked 4th), and monthly average soil moisture for August (ranked 5th).
346 Similarly, in this study, soil moisture was among the top ten-ranked covariates in modelling and
347 explaining long-distance soil type variability across the country.



348

349 **Figure 6.** Random forest covariate relative importance for modelling RSGs. See Appendix B for
350 abbreviations.

351 In this study, lithology showed a relatively low influence on soil variability. This is against the long-
352 standing fact that Ethiopia is believed to be a land of geologic contrast (Abyneh,2005; Alemayehu
353 et al., 2014; Elias., 2016; Jarvis et al., 2011; Zewdie, 2013) characterised by (i) recent and old
354 volcanic activities; (ii) the highlands consisting of igneous rocks (mainly basalts); (iii) steep-sided
355 valleys characterise by strong colluvial and alluvial deposits; (iv) denudation process exposed
356 metamorphic rocks; and (v) occurrence of various sedimentary rocks like limestone and sandstone in
357 the relatively lower areas. The low influence of lithology may be related to the use of a coarse-scale
358 and less detailed lithology map, which may not sufficiently capture the spatial variability of the
359 parent materials.



360 **3.2.2 Model performance**

361 The parameter optimization process resulted in mtry 20, split rule extra trees and minimum node size
362 5. The overall accuracy of the model was 56.24 % which ranged between 54.43% and 58.1% with a
363 95% confidence interval. The kappa values based on the internal cross-validation and testing dataset
364 showed that the overall model performance produced using 10-fold cross-validation with the
365 repeated fitting was 48%. Considering similar area-based digital soil class mapping efforts, the
366 overall purity (accuracy) was in line with the accuracies that were typically reported for soil class
367 maps developed with random forest model (Leenaars et al., 2020a) and statistical methods (Heung et
368 al., 2016; Holmes et al., 2015). Table 1 shows the confusion matrix at validation/testing points i.e.,
369 20 % of the observation. Further, the matrix indicates the producer's accuracy (class representation
370 of observed versus predicted) and user's accuracy (map purity) were not similar for all RSGs. The
371 map purity is in the order of Lixisols, Calcisols, Alisols, Phaeozems, Vertisols, Andosols,
372 Solonchaks, Fluvisols, Arenosols, Leptosols, Luvisols, Nitisols, and Cambisols. However, Vertisols,
373 Calcisols, and Andosols are the observed classes that are best represented by the map followed by
374 Fluvisols, Alisols, Nitisols, Leptosols, Luvisols and Cambisols.

375 Global Soil Grids at 250 m resolution used machine learning algorithms to map the global WRB
376 reference soil groups with map purity and weighted kappa of 28% and 42%, respectively (Hengl et
377 al., 2017). The Soil Grids 250 m WRB soil groups/classes prediction output-spatial soil patterns
378 were not evaluated based on expert knowledge while in this study we did an extensive back and
379 forth qualitative assessment by a panel of pedologists. The quantitative accuracy in the present study
380 (about 56 %) coupled with an expert-based qualitative evaluation of the predicted maps indicated the
381 development and achievement of a substantially enhanced national product for users of spatial soil
382 resources information. This finding is a step forward and acceptable considering that Soil Grids are
383 not expected to be as accurate as locally produced maps and models that use much more local point
384 data and finer local variables (Mulder et al., 2016). Further, the data and finding in this study can
385 help improve the soil maps of Africa as it partially addresses the concern by Hengl et al. (2017) who
386 recognised that WRB RSGs modelling in the global Soil Grids 250 m is critically uncertain for parts
387 of Africa.

388



389 **Table 1.** Confusion matrix of random forest RSG prediction (at validation/testing observations).

Prediction	Reference																		Total	
	Acrisols	Alisols	Andosols	Arenosols	Calcisols	Cambisols	Fluvisols	Gleysols	Leptosols	Lixisols	Luvvisols	Nitisols	Phaeozems	Planosols	Regosols	Solonchaks	Solonetz	Vertisols		
Acrisols	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0.33	3
Alisols	0	40	0	0	0	0	1	1	0	0	9	4	0	0	2	0	0	2	0.68	59
Andosols	0	0	28	1	1	3	5	0	2	0	2	0	0	0	0	0	1	1	0.64	44
Arenosols	0	0	0	11	0	2	1	0	0	0	5	0	0	0	0	0	0	1	0.55	20
Calcisols	0	0	0	0	21	0	1	0	0	0	2	0	0	0	0	0	0	5	0.72	29
Cambisols	2	3	6	9	1	197	28	2	35	2	47	16	5	1	16	3	3	28	0.49	404
Fluvisols	1	0	3	5	1	34	144	0	9	0	15	7	0	0	1	5	5	17	0.58	247
Gleysols	0	0	0	0	0	0	1	2	0	0	1	0	0	1	0	0	0	0	0.40	5
Leptosols	0	1	4	3	3	47	11	0	176	0	27	7	1	0	32	0	0	24	0.52	336
Lixisols	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	1.00	1
Luvvisols	2	16	3	8	0	34	13	2	33	3	216	30	3	0	25	1	0	41	0.50	430
Nitisols	6	8	0	0	1	23	8	3	18	8	29	132	0	1	8	0	1	21	0.49	267
Phaeozems	0	0	0	0	0	0	0	0	0	0	1	0	2	0	0	0	0	0	0.67	3
Planosols	0	0	0	0	0	0	0	0	0	0	1	1	0	5	1	0	0	1	0.11	9
Regosols	0	0	0	0	0	7	1	0	7	1	8	1	0	0	22	0	0	5	0.42	52
Solonchaks	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	3	1	0	0.60	5
Solonetz	0	0	0	0	1	4	1	0	0	0	0	0	0	0	0	1	6	0	0.46	13
Vertisols	3	1	3	5	5	92	32	2	61	3	81	31	5	5	25	2	6	641	0.64	1,003
Producer Accuracy	0.07	0.58	0.60	0.26	0.62	0.44	0.58	0.17	0.51	0.06	0.49	0.58	0.13	0.38	0.17	0.20	0.25	0.81	0.56	-
Total	15	69	47	42	34	443	247	12	342	18	445	229	16	13	132	15	24	787	-	2,930

390 **3.2.3 Modelling and Mapping: EthioSoilGrids Version 1.0**

391 The study identified eighteen reference soil groups in Ethiopia, mapped at 250 m resolution (Figure
 392 7). The model prediction showed that seven soil reference groups including Cambisols, Leptosols,
 393 Vertisols, Fluvisols, Nitisols, Luvvisols, and Calcisols covered nearly 98% of the total land area of the
 394 country (Figure 8). Five soil reference groups (Solonchaks, Arenosols, Regosols, Andosols, and

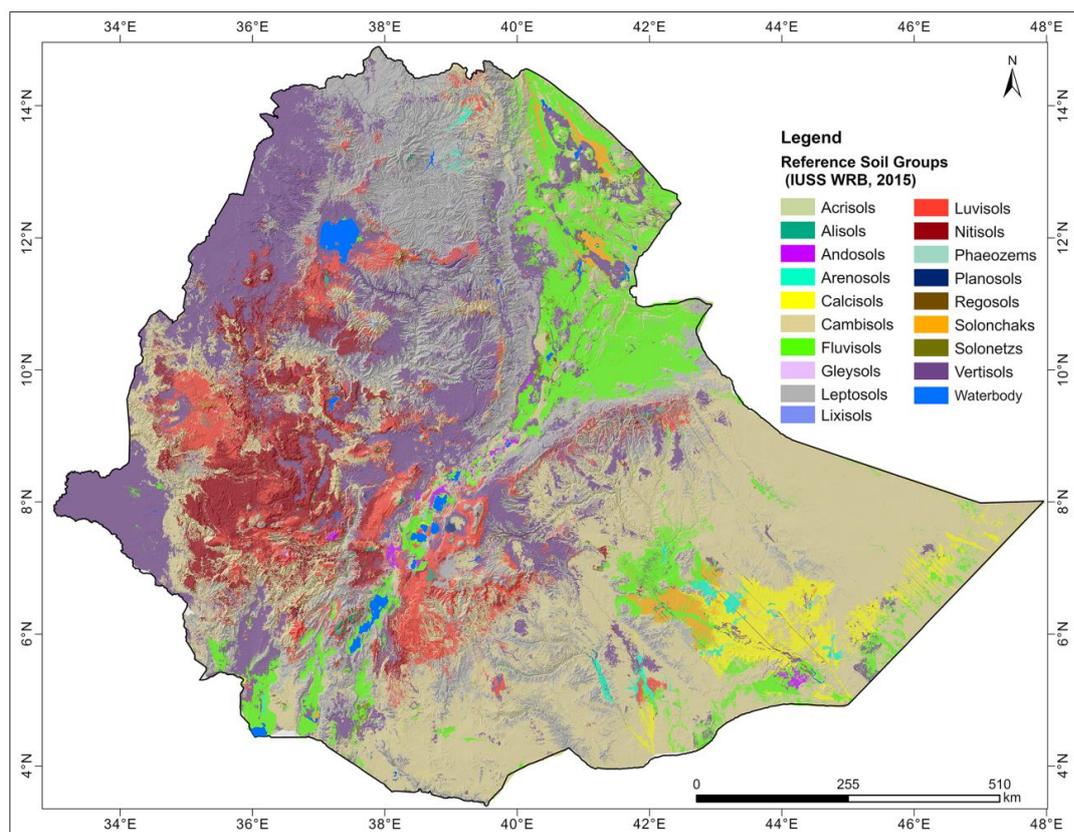


395 Alisols) were estimated to cover about 2% of the land area, while trace coverages of Solonetz (,
396 Planosols, Acrisols, Lixisols, Phaeozems, and Gleysols were also found in some pocket areas.

397 In terms of spatial distribution, Nitisols and Luvisols dominated the northwestern and south-western
398 highlands while the south-eastern lowlands were dominantly covered by Cambisols, Calcisols, and
399 Fluvisols with some Solonchaks. The Vertisols extensively covered the north and south-western
400 lowlands along with the Ethio-Sudan border areas and central highland plateaus. Overall, each RSG
401 position, with other RSGs, along the landscapes/catena/topo-sequence, is in good agreement with the
402 established schematic soil sequence, previous spatial soil information of Ethiopia and with experts'
403 opinions validated across 126 geographic windows of the country.

404 The probability of occurrence of each RSG was mapped (Appendix C) in each modelling spatial
405 window (i.e., the cell size of 250-meter X 250 m). The dominant RSGs were aggregated based on
406 the most probable RSG in each spatial modelling window. There was high correspondence between
407 the top seven ranked prediction probabilities and observed soil types as confirmed visually by
408 overlaying observed classes and prediction probabilities.

409 The overall occurrence and the relative position of each of the RSG along the topo-sequence and its
410 association with other RSGs agree with previous works (Abayneh, 2006; Ali et al., 2010; Abdenna et
411 al., 2018; Asmamaw and Mohammed, 2012; Belay, 2000, 1998, 1997, 1996; Driessen et al., 2001;
412 Elias, 2016; FAO 1984a; Fikre, 2003; Mitku, 1987; Mohammed and Belay, 2008; Mohammed and
413 Solomon, 2012; Mulugeta et al., 2021; Sheleme, 2017; Shimeles et al., 2007; Tolossa, 2015; Zewdie,
414 2013). However, there were cases where the RSGs' position along the topo-sequence and association
415 with other RSGs require further investigation, which was not adequately captured and explained in
416 this study. This might be attributed to the positional accuracy of legacy point observations,
417 modelling approach, and most importantly the level of details and scale/resolution of the
418 environmental variables used in this study.



419

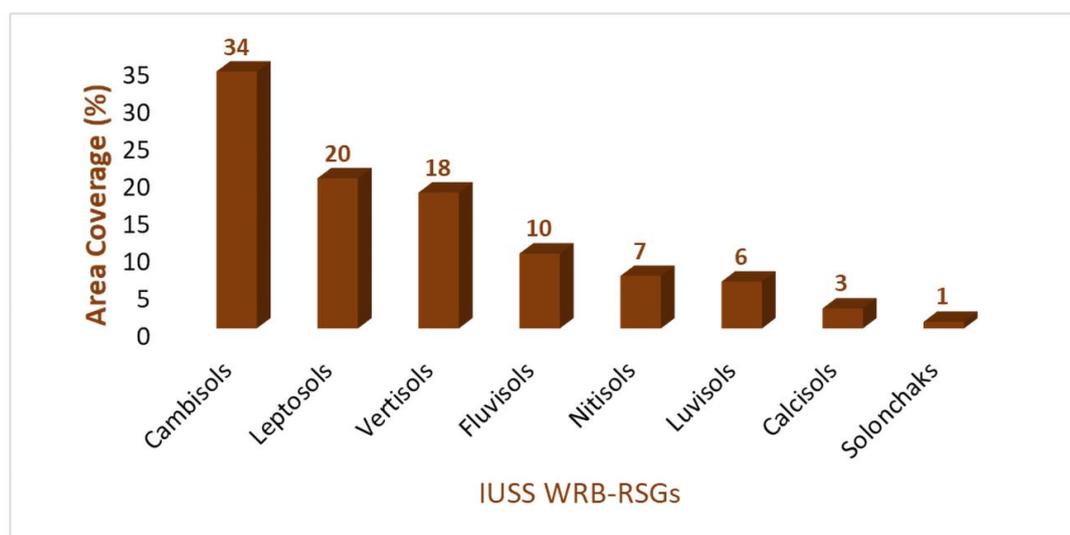
420 **Figure 7.** Major reference soil groups of Ethiopia (EthioSoilGrid V1.0).

421 Considering the third position of Cambisols in the order of frequency occurrence of RSGs per point
 422 observations (following Vertisols and Luvisols), these soils seem to be over-represented on the map
 423 (ranked 1st) apparently at the expense of Vertisols and Luvisols, and to some extent in places of
 424 Leptosols and other RSGs. This might be attributed to the fact that Cambisols create a geographical
 425 continuation with Vertisols and/or Luvisols at the lower slopes and Leptosols/ Regosols at the higher
 426 slopes, suggesting the presence of some bordering soil qualities in respective transitional zones (Ali
 427 et al., 2010; Asmamaw and Mohammed, 2012; Sheleme, 2017; Zewdie, 2013).

428 The proportion of area mapped as Cambisols (34 %) revealed new insights compared with the
 429 information from the most cited spatial soil maps: Cambisols ranked 2nd (21 %), 2nd (16 %), 4th (9
 430 %), and 4th (8 %) as reported by Berhanu (1980), FAO (1984b), FAO (1998), and Soil Grids- Hengl



431 et al (2017), respectively. This might be due to: (i) the number and distribution of profile
432 observations, which is more extensive than the previous ones, (ii) the type and level of details of
433 covariates considered; (iii) variations and rearrangements in the keys for Classification of the RSGs
434 among soil classification versions used in previous studies and misclassification/confusion of
435 Vertisols with Vertic Cambisols, as legacy soil profile data coming from diverse sources.



436 **Figure 8.** The area coverage (in %) for the major WRB RSGs (Note: the remaining 10 RSGs-
437 Arenosols (0.44 %), Regosols (0.35 %), Andosols (0.31%), Alisols (0.16 %), Solonetz (0.04 %),
438 Planosols (0.04 %), Acrisols (0.02 %), Lixisols (0.02%), Phaeozems (0.02 %), and Gleysols (0.01
439 %) were not plotted because of their relatively small area coverage).
440

441 Balanced datasets are ideal to allow decision trees algorithms to produce better classification but for
442 datasets with uneven class size, the generated classification model might be biased towards the
443 majority class (Hounkpatin et al., 2018; Wadoux et al., 2020). This likely scenario requires further
444 investigation for future similar studies and prediction accuracy enhancement.

445 Considering the number and distribution of legacy soil profiles used, the quality monitoring process
446 method was followed to filter dubious soil profiles, and soil classification harmonization protocols
447 were implemented. The study followed a robust modelling framework and generated new insights
448 into the relative area coverage of WRB RSGs in Ethiopia. Further, it provided coherent and up-to-
449 date digital quantitative gridded spatial soil resource information to support the successful



450 implementation of various digital agricultural solutions. The approach used demonstrates the power
451 of data and analytics, and the output is an exemplary use case for similar digital content development
452 efforts in Ethiopia. However, the EthioSoilGrids v1.0 product from this first country-wide RSGs
453 modelling effort requires complementary activities. These include modelling and mapping that
454 should go beyond RSGs and need to include 2nd level classifications. This will be achieved through
455 modelling and mapping a set of principal and supplementary qualifiers along with RSGs which will
456 enable the integration of taxonomy details and requirements with spatial scale protocols, as outlined
457 in IUSS WRB 2015 classification system.

458 **3.3 Expert validation of the soil map**

459 Expert knowledge of soil-landscape relations and soil distribution is important in evaluating the
460 predictive soil mapping results and assessing if predicted spatial patterns make sense from a
461 pedological viewpoint (Hengl et al., 2017). The expert validation workshop participants have
462 commended the initiative and the approach that led to the development of the national soil resource
463 map, including the commitment of the technical experts involved and resources invested in it by
464 partner organizations. Overall, they expressed that the map passed meticulous quality-enhancing
465 processes and that its content and accuracy exceeded their expectations.

466 All three groups have rated the accuracy of the map at 60 +%; of the 126 polygons, they have
467 expressed no concern for 63 %, minor concern for 23 % and a major concern for 14 % of the
468 polygons. While the minor concerns are mostly related to the accuracy of the relative coverage of the
469 predicted dominant soil types, the major concerns may indicate a possible mismatch between the
470 predicted soil type and the experience of some of the group members of the target area such as an
471 important soil type missed out (expected by the experts based on their knowledge of soil coverages
472 and prevailing soil-forming factors in specific areas).

473 After the plenary discussions that followed group presentations, participants have suggested that the
474 final version of the map be released for use after additional desk validation and improvements,
475 especially for the polygons with major concerns. It was recommended to re-run the model after
476 revising the data for the polygons where concerns are reported and use additional data obtained
477 during the event. A small team of senior pedologists was formed to support the core group in



478 revising the data from polygons with reported major concerns. Newly acquired data were cleaned
479 and validated before re-running the model to generate the final version of the map.

480 **4 Conclusions**

481 Coherent and up-to-date country-wide digital soil information is essential to support digital
482 agricultural transformation efforts. This study involved collation, cleaning, harmonization, and
483 validation of the legacy soil profile data sets, involving soil scientists with different backgrounds
484 individually and in groups. To develop the 250 m digital soil resource map, a machine learning
485 modelling approach and expert validation were applied to the harmonised soil database and
486 environmental covariates affecting soil-forming processes. Accordingly, about 20,000 soil profile
487 data have been collated, out of which, about 14,681 were used for the modelling and mapping of
488 eighteen RSGs out of the identified twenty-three RSGs. Although unevenly distributed, the legacy
489 soil profile data used in the modelling covered most of the agro-ecologies of the country. Among the
490 mapped 18 RSGs, the highest number of observed (3,935) profiles represent Vertisols, followed by
491 Luvisols, Cambisols and Leptosols, while Gleysols were represented with the lowest number (63) of
492 profiles. The modelling revealed that MODIS long term reflectance, multiresolution index of valley
493 bottom flatness, land surface temperature, soil moisture, long-term mean annual rainfall, and wetness
494 index of the landscape is the most important covariates for predicting reference soil groups in
495 Ethiopia.

496 Our ten-fold spatial cross-validation result showed an overall accuracy of about 56 % with varying
497 accuracy levels among RSGs. The modelling result revealed that seven major soil reference groups
498 including Cambisols (34 %), Leptosols (20 %), Vertisols (18 %), Fluvisols (10 %) Nitisols (7 %),
499 Luvisols (6 %) and Calcisols (3 %) covered nearly 98 % of the total land area of the country, while
500 minor coverage of other reference soil groups (Solonchaks, Arenosols, Regosols, Andosols, Alisols,
501 Solonetz, Planosols, Acrisols, Lixisols, Phaeozems, and Gleysols) were also detected in some areas.
502 Compared to the existing soil resource map, the coverage of the first three major soil groups has
503 substantially increased which is related to the increased availability of soil profile data covering
504 larger areas of the country, implying that these soils were previously underestimated. Cambisols and
505 Vertisols which together represent nearly half of the total land area are relatively young with
506 inherent fertility, implying the high agricultural potential for the country. However, given their



507 limitations, these and the other soil types require the implementation of suitable land, water, and
508 crop management techniques to sustainably exploit their potential.

509 Given its resolution and quantitative digital representation, the map will have tremendous
510 significance in both agricultural and other land-based development planning while safeguarding the
511 environment. For instance, the accessibility of good quality digital soil data is crucial for developing
512 and using decision support tools (DSTs) such as land use and management decisions. However,
513 effective use of the map requires that the associated WRB second-level classification including
514 principal and supplementary qualifiers and soil atlas providing details of the soil physicochemical
515 properties be accessed together with the map, which the authors and others responsible need to
516 prioritize in their future endeavours.

517



518 **Appendix A: Legacy soil profile data distribution**

519 **Table A1.** Distribution of legacy soil profile data by agroecology zones.

MAJOR_AGRO	AEZ area coverage (%)*	Profiles Observation (%)**
Warm arid lowland plains	19.76	3.40
Warm moist lowlands	15.12	10.74
Hot arid lowland plains	10.79	2.44
Warm sub-moist lowlands	9.63	6.94
Tepid moist mid highlands	8.05	20.21
Warm sub-humid lowlands	7.11	5.69
Tepid sub-humid mid highlands	6.63	15.26
Tepid sub-moist mid highlands	5.17	12.39
Warm semi-arid lowlands	2.75	3.23
Tepid humid mid highlands	2.65	2.48
Warm humid lowlands	2.29	0.45
Cool moist mid highlands	1.74	4.15
Hot sub-humid lowlands	1.67	0.07
Cool sub-moist mid highlands	1.16	3.00
Cool humid mid highlands	0.82	1.01
Warm per-humid lowlands	0.68	0.01



MAJOR_AGRO	AEZ area coverage (%)*	Profiles Observation (%)**
Hot moist lowlands	0.59	3.56
Hot sub-moist lowlands	0.56	0.03
Cool sub-humid mid highlands	0.52	1.38
Tepid arid mid highlands	0.43	0.39
Hot semi-arid lowlands	0.40	2.05
Tepid semi-arid mid highlands	0.19	0.67
Cold moist sub-afro-alpine to afro-alpine	0.07	0.16
Cold sub-moist mid highlands	0.07	0.04
Cold sub-humid sub-afro-alpine to afro-alpine	0.06	0.03
Cold humid sub-afro-alpine to afro-alpine	0.06	0.01
Very cold humid sub-afro-alpine	0.04	0.02
Very cold sub-moist mid highlands	0.02	0.02
Very cold moist sub-afro-alpine to afro-alpine	0.01V	0.03
Hot per-humid lowlands	0.01	0.15
Tepid perhumid mid highland	0.13	0
Very cold sub-humid sub-afro alpine to afro-alpine	0.03	0

520 Note: *= total area of Ethiopia 1.14mln km²; **=total number of profiles 14,681

521



522 **Appendix B: Environmental covariates**

523 **Table B1.** List, description, spatial and temporal extent, and source of covariates used in modelling
524 the reference soil groups.

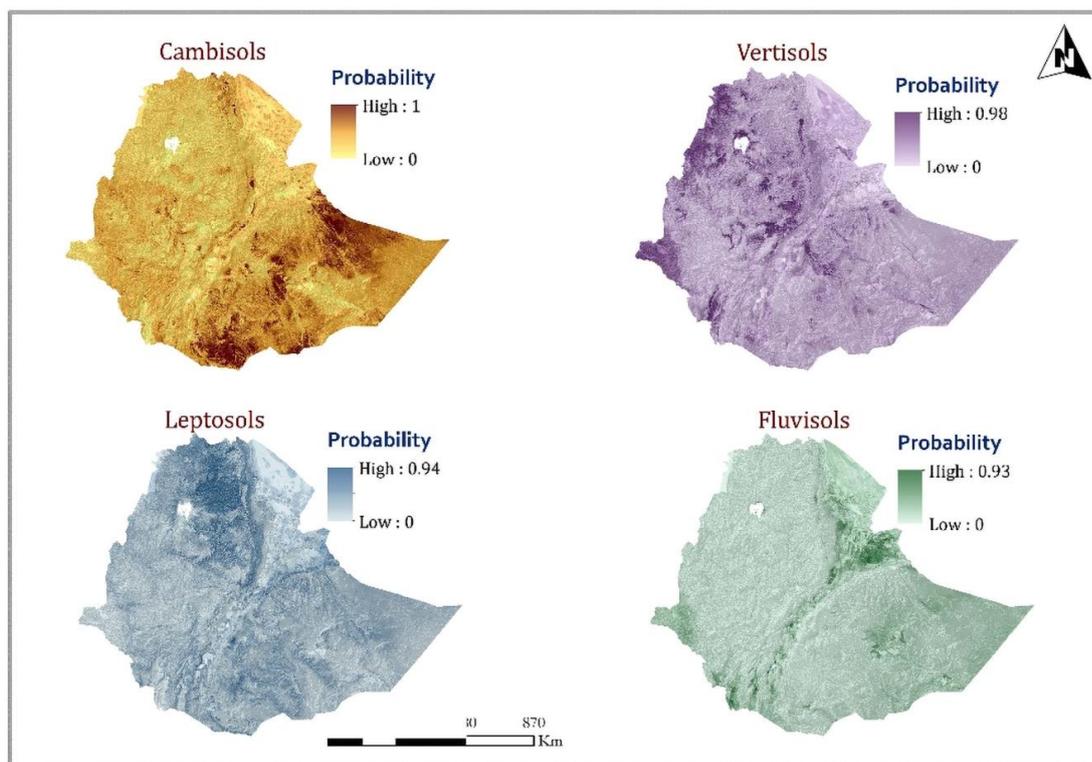
Categories	Covariates	Descriptions	Spatial resolution	Temporal resolution	Source
Climate	prep	Precipitation	4 km	1981 - 2016	ENACTS (Dinku et al.,2014)
	prep_sd	The standard deviation of precipitation	4 km	1981 - 2016	Derived from ENACTS (Dinku et al.,2014)
	tmax	Maximum Temperature	4 km	1983 - 2016	ENACTS (Dinku et al.,2014)
	tmin	Minimum Temperature	4 km	1983 - 2016	ENACTS (Dinku et al.,2014)
	trange	Temperature range	4 km	1983 - 2016	ENACTS (Dinku et al.,2014)
	tav_sd	Standard deviation of average temperature	4 km	1983 - 2016	Derived from ENACTS (Dinku et al.,2014)
	pet	Potential evapotranspiration	4 km	1981 - 2016	Derived from ENACTS (Dinku et al.,2014) using Modified Penman method
	lstn	Land surface temperature- Day (Aqua MODIS- MYD11A2 , time series monthly average)	1000 m	2002-2018	AfSIS ^a
	lstn	Land surface temperature-Night (Aqua MODIS- MYD11A2 , time series monthly average)	1000 m	2002-2018	AfSIS
	soil_moist	Soil Moisture (Derived from one-dimensional soil water balance)	4 km	1981 - 2016	Ethiopian Digital AgroClimate Advisory Platform (EDACaP)
soil_temp	Soil temperature	30 km	1979 - 2019	ERA 5-Reanalysis ECMWF data ^b	
Topography	DEM	Digital elevation model (Elevation)	90 m	-	SRTM- DEM (Vågen, 2010)
	twi	Topographic wetness Index	90 m	-	SAGA GIS-based



Categories	Covariates	Descriptions	Spatial resolution	Temporal resolution	Source
					SRTM-DEM derivative
	aspect	Topographic Aspect	90 m	-	SAGA GIS-based SRTM-DEM derivative
	curv	Topographic Curvature	90 m	-	SAGA GIS-based SRTM-DEM derivative
	conv	Topographic convergence index	90 m	-	SAGA GIS-based SRTM-DEM derivative
	ls	Slope Length and Steepness factor (ls_factor)	90 m	-	SAGA GIS-based SRTM-DEM derivative
	morph	Terrain Morphometry	90 m	-	SAGA GIS-based SRTM-DEM derivative
	mrvbf	Multiresolution index of valley bottom flatness	90 m	-	SAGA GIS-based SRTM-DEM derivative
	slope	Slope class (%)	90 m	-	SAGA GIS-based SRTM-DEM derivative
Vegetation	ndvi	Normalised Difference Vegetation Index (NDVI) (MODIS- MODIS MOD13Q1, time series monthly average)	250 m	2000-2021	AfSIS ^a
	evi	Enhanced Vegetation Index (EVI) (MODIS- MODIS MOD13Q1, time series monthly average)	250 m	2000-2021	AfSIS
	lulc	Land use/ landcover	30 m	2010	Water and Land Resource Centre-Addis Ababa University (WLRC-AAU, 2010)
parent material	lithology	Geology/parent material	1:2,000,000	1996	The Ethiopian Geological Survey (Tefera et al.,1996)
MODIS spectral reflectance	ref1	Red band (MODIS- MODIS MOD13Q1, time series monthly average)	250 m	2000 – 2018	AfSIS ^a
	ref2	Near-Infrared (MODIS- MODIS MOD13Q1, time series monthly average)	250 m	2000 – 2018	AfSIS
	ref7	Mid-Infrared (MODIS- MODIS MOD13Q1, time series monthly average)	250 m	2000 – 2018	AfSIS



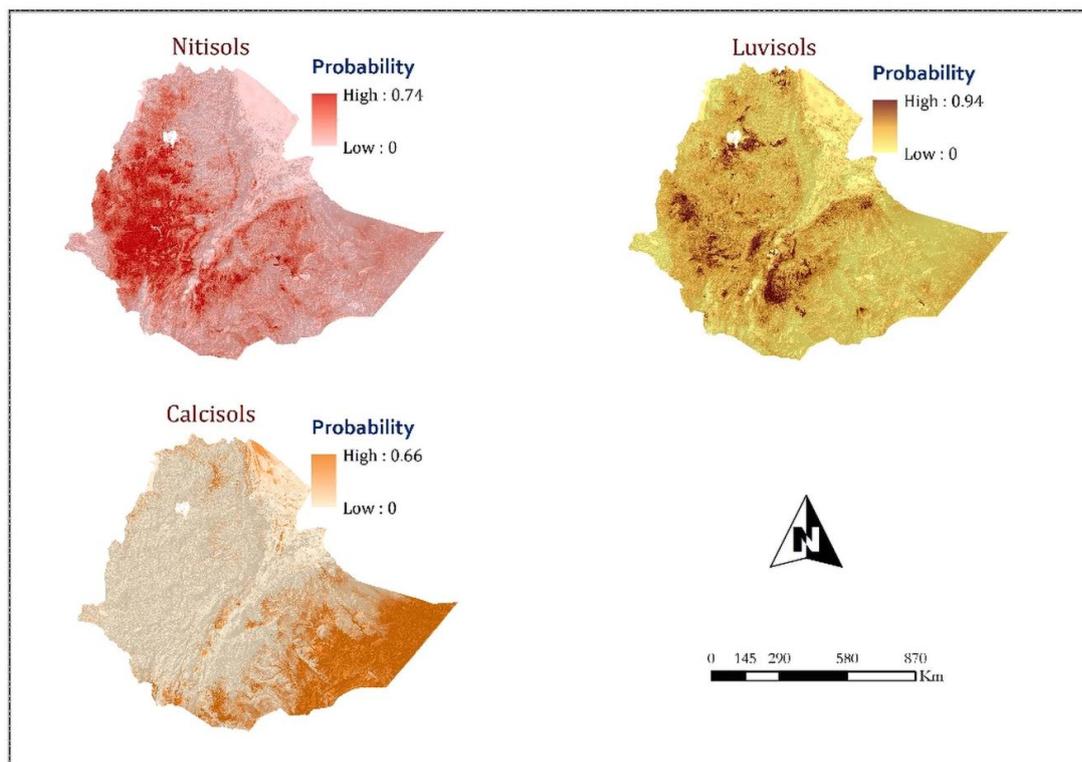
525 **Appendix C: Probability of occurrence of reference soil groups**



526

527 **Figure C1.** Occurrence probability maps of Cambisols, Leptosols, Vertisols, and Fluvisols.

528



529

530

Figure C2. Occurrence probability maps of Nitisols, Luvisols, and Calcisols.

531



532 **Data availability.** Data will be available upon request based on the CoW guideline.

533 **Author contributions.** AA, TE, KG, WA, and LT conceived and designed the study, perform the
534 analysis, and wrote the first draft, with substantial input and feedback from all authors. EM, TM,
535 NH, AY, AM, TA, FW, AL, NT, AA, SG, YA, and BA, contributed to input data preparation, data
536 encoding, and harmonization. Legacy data validation and review of subsequent versions of the paper
537 were performed by MH, WH, AA, DT, GB, MG, SB, MA, AR, YGS, ST, DA, YW, DB, EZ, SC,
538 and EE.

539 **Competing interests.** The authors declare that they have no conflict of interest.

540 **Acknowledgments.** We sincerely appreciate the coalition of the willing (CoW) members who are
541 instrumental in providing, collating, cleaning, standardizing and harmonizing the legacy soil profile
542 data used in generating the soil resource map of Ethiopia at 250 m resolution. The CoW team also
543 deserves credit for inspiring many to share data and develop an integrated national database related
544 to agronomy and soil profile data. The leadership of the Natural Resource Development Sector and
545 Soil Resource Information and Mapping Directorate of the Ministry of Agriculture (MoA) have
546 played a crucial role including assigning experts from the Ministry and other organizations who
547 worked on collating, encoding, harmonizing, and processing the soil survey legacy data are the
548 foundation for the soil resource map. Various institutions, as well as the late and present soil
549 surveyors and pedologists, deserve special recognition for their contributions to the generation and
550 sharing of soil profile data. The senior pedologists and soil surveyors who provided invaluable
551 support to check and harmonize thousands of soil profiles and laboratory results are sincerely
552 appreciated. They worked very hard with positive energy for which we are very grateful. In addition,
553 the same group of experts and additional ones who supported the validation of the preliminary soil
554 resource map deserve credit for their commitment to contributing their experiences. Our sincere
555 appreciation also goes to the continued and persistent support of GIZ-Ethiopia mainly through the
556 project- Supporting Soil Health Interventions in Ethiopia (SSHI), which supported and facilitated the
557 activities of the CoW. The Alliance of Bioversity and CIAT is highly acknowledged for coordinating
558 CoW and its efforts and supporting the implementation of activities that are of high national
559 importance. We would also like to sincerely thank the Excellence in Agronomy (EiA) CGIAR
560 Initiative, which has brought huge contributions to this project in terms of funding and building skill
561 of the various teams. The Water, Land and Ecosystems (WLE) and Climate Change, Agriculture and
562 Food Security (CCAFS) programs of the CGIAR also provided support in various forms. Recently,
563 our work is benefiting from the Accelerating Impacts of CGIAR Climate Research in Africa



564 (AICCRA) project supported by the World Bank in terms of data, analytics, and resources to support
565 data linkage and integration.

566 **Financial support.** This work was supported, in whole or in part, by the Bill & Melinda Gates
567 Foundation [INV-005460]. Under the grant conditions of the Foundation, a Creative Commons
568 Attribution 4.0 Generic License has already been assigned to the Author Accepted Manuscript
569 version that might arise from this submission.

570



571 **References**

- 572 Abayneh, E.: Application of Geographic Information System (GIS) for soil resource study in
573 Ethiopia, in: Proceedings of the National Sensitization Workshop on Agro metrology and
574 GIS, 17-18 December 2001, Addis Ababa, Ethiopia, 162-169, 2001.
- 575 Abayneh, E.: Characteristics, Genesis and Classification of Reddish Soils from Sidamo Region of
576 Ethiopia, PhD Thesis, Universiti Putra Malaysia, 2005.
- 577 Abayneh, E., Zauyah, S., Hanafi, M. M., and Rosenani, A. B.: Genesis and classification of
578 sesquioxidic soils from volcanic rocks in sub-humid tropical highlands of Ethiopia,
579 Geoderma, 136(3-4), 682–695, <https://doi.org/10.1016/j.geoderma.2006.05.006> , 2006.
- 580 Abayneh, E., and Berhanu, D.: Soil Survey in Ethiopia: Past, Present and the Future, in: Proceedings
581 of the 8th Conference of the Ethiopian Society of Soil Science, Soils for sustainable
582 development, 27-28 April, 2006, Addis Ababa, Ethiopia, 2007.
- 583 Abdenna, D., Yli-Halla, M., Mohamed, M., and Wogi, L.: Soil classification of humid Western
584 Ethiopia: A transect study along a toposequence in Didessa watershed, Catena, 163,184-195,
585 <https://doi.org/10.1016/j.catena.2017.12.020>, 2018.
- 586 Abegaz, A., Ashenafi, A., Tamene, L., Abera, W., and Smith, Jo. U.: Modeling long-term attainable
587 soil organic carbon sequestration across the highlands of Ethiopia. Environ. Dev. Sustain.,
588 24, 131–5162, <https://doi.org/10.1007/s10668-021-01653-0>, 2022.
- 589 Abera, W., Tamene, L., Tesfaye, K., Jiménez, D., Dorado, H., Erkossa, T., and Ramirez-Villegas, J.
590 : A data-mining approach for developing site-specific fertilizer response functions across the
591 wheat-growing environments in Ethiopia, *Experimental Agriculture*, 1-1, 2022.
- 592 AfSIS: Africa Soil Information Service project, Covariates for land and climate developed from
593 remotely sensed data, Earth institute, Columbia University, New York,
594 <http://africasoils.net/services/data/remote-sensing/land/>, 2020.
- 595 Alemayehu, R., Van Daele, K., De Paepe, P., Dumon, M., Deckers, J., Asfawossen, A., and Van
596 Ranst, E.: Characterizing weathering intensity and trends of geological materials in the Gilgel
597 Gibe catchment, southwestern Ethiopia, Journal of African Earth Sciences, 99 (2), 568-580,
598 <https://doi.org/10.1016/j.jafrearsci.2014.05.012>, 2014.
- 599 Ali, A., Tamene, L., and Erkossa, T.: Identifying, Cataloguing, and Mapping Soil and Agronomic
600 Data in Ethiopia, CIAT Publication No. 506, International Center for Tropical Agriculture
601 (CIAT), Addis Ababa, Ethiopia, <https://hdl.handle.net/10568/110868>, 2020.



- 602 Ali, A., Abayneh , E., and Sheleme, B.: Characterizing soils of Delbo Wegene watershed, *J. Soil*
603 *Sci. Environ. Manage.*, 1 (8),184-199, 2010.
- 604 Asmamaw, L., and Mohammed, A.: Characteristics and classification of the soils of Gerado
605 catchment, Northeastern Ethiopia, *EJNRS*, 12(1 and 2), 1-22, 2012
606
- 607 Batjes, N., Ribeiro, E., and van Oostrum, Ad.: Standardized soil profile data to support global
608 mapping and modeling (WoSIS snapshot 2019), *Earth Sys. Sc. Data.*, 12, 299-320, 2020.
- 609 Baveye, P.C., Jacques, B., and John, G.: Soil “Ecosystem” Services and Natural Capital: Critical
610 Appraisal of Research on Uncertain Ground, *Front. in Environ. Sci.*, 4:41,
611 <https://www.frontiersin.org/article/10.3389/fenvs.2016.00041>, 2016.
- 612 Belay ,T.: Characteristics and Landscape relationships of Vertisols and Vertic Luvisols of Melbe,
613 Tigray, Ethiopia. *SINET: Ethiopian Journal of Science* 19 (1): 93-115, 1996.
- 614 Belay, T.: Variabilities of Soil Catena on Degraded Hill Slopes of Wtiya Catchment, Wello,
615 Ethiopia, *SINET: J.Sc.*, 20 (2), 151-175, 1997.
- 616 Belay ,T. : Pedogenesis and soil-geomorphic relationships on the Piedmont slopes of Wurgo Valley,
617 South Welo, Ethiopia, *SINET: J.Sc.*, 21(1), 91-111, 1998.
- 618 Belay, T.: Characteristics and classification of soils of Gora Daget forest, South welo highlands,
619 Ethiopia, *SINET: J.Sc.*, 23(1), 35-51, 2000.
- 620 Berhanu, D.: A survey of studies conducted about soil resources appraisal and evaluation for
621 rural development in Ethiopia, IAR, 1980.
- 622 Berhanu, D.: The soils of Ethiopia: Annotated bibliography, Regional Soil Conservation Unit
623 (RSCU), Swedish International Development Authority (SIDA), Tech. handbook no. 9, 1994.
- 624 Billi, P.: Geomorphological landscapes of Ethiopia, in: *Landscapes and Landforms of Ethiopia*,
625 *World Geomorphological Landscapes*, Springer, Dordrecht, 3–32,
626 <https://doi.org/10.1007/978-94-017-8026-1>, 2015.
- 627 Breiman, L.: RandomForests, *Machine Learning*, 45, 5–32,<https://doi.org/10.1023/A:1010933404324> ,
628 2001.
- 629 Brungard, C. W., Boettinger, J. J., Duniway, M. C., Wills, S. A., and Edwards, W. T. C.: Machine
630 learning for predicting soil classes in three semi-arid landscapes, *Geoderma*, 239–240, 68–
631 83, <https://doi.org/10.1016/j.geoderma.2014.09.019>, 2015.
- 632 Brunner, M.: A National Soil Model of Ethiopia: A Geostatistical approach to Create a National Soil
633 Map of Ethiopia on the Basis of an SRTM 90 DEM and SOTWIS Soil Data, A Master’s
634 Thesis, the Univ. of Bern, Switzerland, 2012.



- 635 Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V.,
636 and Böhner, J.: System for Automated Geoscientific Analyses (SAGA) v. 2.1.4, *Geosci.*
637 *Model Dev.*, 8, 1991–2007, <https://doi.org/10.5194/gmd-8-1991-2015>, 2015.
- 638 Dinku, T., Block, P., Sharoff, J., Hailemariam, K., Osgood, D., del Corral, J., Rémi Cousin, R., and
639 Thomson, M. C.: Bridging critical gaps in climate services and applications in Africa. *Earth*
640 *Perspectives*, 1(1), 1-13, <https://doi.org/10.1186/2194-6434-1-15>, 2014.
- 641 Donahue, R. L.: Ethiopia: Taxonomy, cartography and ecology of soils, Michigan State Univ.,
642 African Stud. Center and Inst.Int.Agric.,Comm., Ethiopian Stud., Occasional Papers Series,
643 Monograph 1, 1962.
- 644 Driessen, P. M., Deckers, J., Spaargaren, O., and Nachtergaele, F.: Lecture notes on the major soils
645 of the world, world soil resources reports No. 94, FAO, Rome, www.fao.org/3/a-y1899e,
646 2001.
- 647 Elias, E.: Soils of the Ethiopian Highlands: Geomorphology and Properties, CASCAPE Project,
648 ALTERRA, Wageningen UR, the Netherlands, library.wur.nl/WebQuery/isric/2259099,
649 2016.
- 650 Enyew, B. D., and Steeneveld, G. J.: Analysing the impact of topography on precipitation and
651 flooding on the Ethiopian highlands, *JGeol. Geosci*, 3(2), 2014.
- 652 Erkossa, T., Laekemariam, F., Abera, W., and Tamene, L.: Evolution of soil fertility research and
653 development in Ethiopia: From reconnaissance to data-mining approaches, *Experimental*
654 *Agriculture*, 58, E4. doi:10.1017/S0014479721000235, 2022.
- 655 FAO: Assistance to Land Use-Planning, Ethiopia: Provisional Soil Association Map of Ethiopia,
656 Field document No. 6, The United Nations Development Programme and Food and
657 Agriculture Organization, FAO, Rome, 1984a.
- 658 FAO: Assistance to Land Use-Planning, Ethiopia: Geomorphology and soils, Field Document AG
659 DP/ ETH/78/003, The United Nations Development Programme and FAO, FAO, Rome,
660 1984b.
- 661 FAO: FAO/Unesco Soil Map of the World, revised legend, World Resources Report 60, FAO,
662 Rome, Reprinted, with corrections, as Tec. Pap. 20, ISRIC, Wageningen, 1989,
663 Library.wur.nl/WebQuery/isric/2264662, 1988.
- 664 FAO: Guideline for Soil Description, Fourth edition, FAO, Rome, Italy, 2006.



- 665 FAO: The Soil and Terrain Database for north-eastern Africa, Crop production systems zones of the
666 IGAD sub region, Land and water digital media series 2, FAO,
667 Rome.1998.
- 668 Fazzini, M., Bisci, C., and Billi, P.: The Climate of Ethiopia, in: Landscapes and Landforms of
669 Ethiopia, World Geomorphological Landscapes, edited by: Billi, P., Springer, Dordrecht, the
670 Netherlands, 65 – 87, https://doi.org/10.1007/978-94-017-8026-1_3, 2015.
- 671 Fikre, M.: Pedogenesis of major volcanic soils of the southern central Rift Valley region, Ethiopia,
672 MSc. Thesis. University of Saskatchewan, Saskatoon, Canada, 2003.
- 673 Fikru, A.: Need for Soil Survey Studies, in: Proceedings of the first soils science research review
674 workshop, 11-14 February 1987, 1988.
- 675 Fikru, A.: Soil resources of Ethiopia, in: Natural Resources Degradation a Challenge to Ethiopia,
676 First Natural Resources Conservation conference, IAR, 1980.
- 677 Hengl, T., and MacMillan, R. A.: Predictive Soil Mapping with R, OpenGeoHub foundation,
678 Wageningen, the Netherlands, www.soilmapper.org, ISBN: 978-0-359-30635-0, 2019.
- 679 Hengl, T., Heuvelink, G. B. M., Kempen, B., Leenaars, J. G. B., Walsh, M. G., Shepherd, K. D.,
680 Sila, A., MacMillan, R. A., Mendes de Jesus, J., Tamene, L., and Tondoh, J. E.: Mapping soil
681 properties of Africa at 250 m resolution: random forest significantly improve current
682 predictions, PLoS ONE 10 (6), <https://doi.org/10.1371/journal.pone.0125814>, 2015.
- 683 Hengl, T., Mendes de Jesus, J., Heuvelink, G. B., Ruiperez Gonzalez, M., Kilibarda, M., Blagotić,
684 A., Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A.,
685 Vargas, R., MacMillan, R. A., Batjes, N. H., Leenaars, J. G., Ribeiro, E., Wheeler, I., Mantel,
686 S., and Kempen, B.: SoilGrids250m: Global gridded soil information based on machine
687 learning, PloS one, 12(2), e0169748, <https://doi.org/10.1371/journal.pone.0169748>, 2017.
- 688 Hengl, T., Miller, M., Križan, J., Shepherd, K. D., Sila, A., Kilibarda, M., Antonijević, O., Glušica,
689 L., Dobermann, A., Haefele, S. M., McGrath, S. P., Acquah, G. E., Collinson, J., Parente, L.,
690 Sheykhmousa, M., Saito, K., Johnson, J. M., Chamberlin, J., Silatsa, F., Yemefack, M.,
691 Wendt J, M.R.A., and Crouch, J.: African soil properties and nutrients mapped at 30 m
692 spatial resolution using two-scale ensemble machine learning, Scientific reports, 11(1), 6130,
693 <https://doi.org/10.1038/s41598-021-85639-y>, 2021.



- 694 Hengl, T., Nussbaum, M., Wright, M. N., Heuvelink, G. B. M., and Gräler B.: Random forest as a
695 generic framework for predictive modeling of spatial and spatio-temporal variables, *PeerJ*, 6,
696 <https://doi.org/10.7717/peerj.5518>, 2018.
- 697 Heung, B., Hung, C. H., Zhang, J., Knudby, A., Bulmer, C. E. and Schmidt, M. G.: An overview and
698 comparison of machine-learning techniques for classification purposes in digital soil
699 mapping, *Geoderma*, 265, 62-77, 2016.
- 700 Hounkpatin, K. O. L., Schmidt, K., Stumpf, F., Forkuor, G., Behrens, T., Scholten, T., Amelung, W.,
701 and Welp, G.: Predicting reference soil groups using legacy data: A data pruning and
702 Random Forest approach for tropical environment (Dano catchment, Burkina Faso), *Sci Rep*.
703 2018; 8, 9959, <https://doi.org/10.1038/s41598-018-28244-w>, 2018.
- 704 Hurni H.: Agro-ecological Belts of Ethiopia: Explanatory Notes on three maps at a scale of
705 1:1,000,000, *Soil Cons. Res. Pro.*, University of Bern, (Switzerland) in Association with the
706 Ministry of Agriculture, Addis Ababa, 1998.
- 707 Iticha, B., and Chalsissa, T.: Digital soil mapping for site-specific management of soils, *Geoderma*,
708 351 (85-91), *Geoderma*, <https://doi.org/10.1016/j.geoderma.2019.05.026>, 2019.
- 709 IUSS Working Group WRB.: World Reference Base for Soil Resources 2014, update 2015
710 International soil classification system for naming soils and creating legends for soil maps,
711 *World Soil Resources Reports No. 106*, FAO, Rome, 2015.
- 712 Jarvis, A., Reuter, H. I., Nelson, A., and Guevara, E.: Hole-filled SRTM for the globe Version 4,
713 CGIARCSI SRTM 90m Digital Elevation Database v4.1.,
714 <http://www.cgiarcsi.org/data/elevation/item/45-srtm-90m-digital-elevation-database-v41>,
715 2011.
- 716 Kempen, B., Brus, D. J., Heuvelink, G. B. M., and Stoorvogel, J. J.: Updating the 1:50,000 Dutch
717 soil map using legacy soil data: A multinomial logistic regression approach, *Geoderma*,
718 151(34), 311–326, <https://doi.org/10.1016/j.geoderma.2009.04.023>, 2009.
- 719 Kempen, B., Brus, D. J., Stoorvogel, J. J., Heuvelink, G. B. M., and de Vries, F.: Efficiency
720 comparison of conventional and digital soil mapping for updating soil maps, *SSSA J.*, 76 (6),
721 2097–2115, <https://doi.org/10.2136/sssaj2011.0424>, 2012.
- 722 Kottek, M., Grieser, J., Beck, C., Rudolf, B., and Rubel, F.: World map of the Köppen-Geiger
723 climate classification updated, *Meteorologische Zeitschrift*, 15. 259-263. 10.1127/0941-
724 2948/2006/0130, 2006.



- 725 Kuhn, M.: Building predictive Models in R using the caret package, *Jour. of Stat. Soft.*, 28(5), 1 –
726 26, doi:<http://dx.doi.org/10.18637/jss.v028.i05>, 2008.
- 727 Leenaars, J. G. B., van Oostrum, A.J.M., and Ruiperez ,G.M.: Africa Soil Profiles Database, Version
728 1.2. A compilation of georeferenced and standardised legacy soil profile data for Sub-
729 Saharan Africa (with dataset), ISRIC Report 2014/01, Africa Soil Information Service
730 (AfsIS) project and ISRIC – World Soil Information, Wageningen,
731 library.wur.nl/WebQuery/isric/2259472, 2014.
- 732 Leenaars, J. G. B., Eyasu, E., Wösten, H., Ruiperez González, M., Kempen, B., Ashenafi, A., and
733 Brouwer, F.: Major soil-landscape resources of the cascape intervention woredas, Ethiopia:
734 Soil information in support to scaling up of evidence-based best practices in agricultural
735 production (with dataset), CASCAPE working paper series No. OT_CP_2016_1, Cascape.
736 <https://edepot.wur.nl/428596>, 2016.
- 737 Leenaars, J. G. B., Elias, E., Wösten, J. H. M., Ruiperez-González, M., and Kempen, B.: Mapping
738 the major soil-landscape resources of the Ethiopian Highlands using random forest,
739 *Geoderma*, 361, <https://doi.org/10.1016/j.geoderma.2019.114067>, 2020a.
- 740 Leenaars, J. G. B., Ruiperez, M., González, M., Kempen, B., and Mantel, S.: Semi-detailed soil
741 resource survey and mapping of REALISE woredas in Ethiopia, Project report to the
742 BENEFIT-REALISE programme, December, ISRIC-World Soil Information, Wageningen,
743 2020b.
- 744 McBratney, A. B., Santos, M. M., and Minasny, B.: On digital soil mapping, *Geoderma*, 117 (1-2),
745 3-52, 2003.
- 746 Mesfin, A.: Nature and Management of Ethiopian Soils, ILRI, 272, 1998.
- 747 Mishra, B. B., Gebrekidan, H., and Kibret, K.: Soils of Ethiopia: Perception, appraisal and
748 constraints in relation to food security, *JFAE*, 2(3 and 4): 269-279, 2004.
- 749 Mitiku, H.: Genesis, characteristic and classification of the Central Highland soils of Ethiopia, Ph.D.
750 Thesis, State University of Ghent, Belgium, 1987.
- 751 Mohammed, A., and Belay, T.: Characteristics and classification of the soils of the Plateau of Simen
752 Mountains National Park (SMNP), Ethiopia, *SINET: EJSc.*,31 (2), 89-102, 2008.
- 753 Mohammed, A. and Solomon ,T. : Characteristics and fertility quality of the irrigated soils of
754 Sheneka, Ethiopia, *EJNR*,12 (1 and 2), 1-22, 2012.



- 755 Mulder, V. L., Lacoste, M., Richer de Forges, A. C., and Arrouays, D.: GlobalSoilMap France: high-
756 resolution spatial modelling the soils of France up to two meter depth. *Science of the Total*
757 *Environment* 573, 1352-1369, 2016.
- 758 Mulualem, A., Gobezie, T.B., Kasahun, B., and Demese, M.: Recent Developments in Soil Fertility
759 Mapping and Fertilizer Advisory Services in Ethiopia, A Position Paper,
760 <https://www.researchgate.net/publication/327764748/>, 2018.
- 761 Mulugeta, T., Seid, A., Kefyalew, T., Mulugeta, F., and Tadla , G.: Characterization and
762 Classification of Soils of Askate Subwatershed, Northeastern Ethiopia, *Agri., For. and*
763 *Fisheries*, 10 (3) , 112-122, doi: 10.11648/j.aff.20211003.13, 2021.
- 764 Poggio, L., de Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E., and
765 Rossiter, D.: Soil Grids 2.0: producing soil information for the globe with quantified spatial
766 uncertainty, 2020.
- 767 R Core Team R: A Language and Environment for Statistical Computing, R Foundation for
768 Statistical Computing, Vienna, 2020.
- 769
- 770 Sheleme, B.: Topographic positions and land use impacted soil properties along Humbo Larena-Ofa
771 Sere toposequence, Southern Ethiopia, *JSSEM*, 8(8),135-147,
772 <https://doi.org/10.5897/JSSEM2017.0643>, 2017.
- 773 Shimeles, D., Mohamed, A., and Abayneh, E.: Characteristics and classification of the soils of
774 Tenocha Wenchacher Micro catchment, South west Shewa, Ethiopia. *EJNRS*, 9 (1), 37- 62,
775 2007.
- 776 Soil Science Division Staff: Soil survey manual, edited by: Ditzler, C., Scheffe, K., and Monger,
777 H.C., USDA Handbook 18, Government Printing Office, Washington, D.C., 2017.
- 778 Svetnik, V., Liaw, A., Tong, C., Culberson, J.C., Sheridan, R.P., and Feuston, B.P.: Random forest:
779 a classification and regression tool for compound classification and QSAR modeling, *J. of*
780 *Che. Info. and Com. Sc.*, 43, 1947–1958, doi: 10.1021/ci034160g, 2003.
- 781 Tamene, L. D., Amede, T., Kihara, J., Tibebe, D., and Schulz, S.: A review of soil fertility
782 management and crop response to fertilizer application in Ethiopia: towards
783 development of site- and context-specific fertilizer recommendation, CIAT
784 Publication No. 443, International Center for Tropical Agriculture (CIAT), Addis
785 Ababa, Ethiopia, <hdl.handle.net/10568/82996>, 2017.



- 786 Tamene, L., Erkossa, T., Tafesse, T., Abera, W., and Schultz, S.: A coalition of the willing
787 powering data-driven solutions for Ethiopian agriculture, CIAT Publication No. 518, CIAT,
788 Addis Ababa, Ethiopia, 2021.
- 789 Tefera, M., Chernet, T., and Workneh, H.: Geological Map of Ethiopia, Addis Ababa, Ethiopia:
790 Federal Democratic Republic of Ethiopia, Ministry of Mines and Energy, Ethiopian Institute
791 of Geological Surveys, 1999.
- 792 Tolossa, A.R.: Vertic Planosols in the Highlands of South-Western Ethiopia: Genesis,
793 Characteristics and Use, Ghent University, Faculty of Sciences, 2015.
- 794 Vågen, T.G.: Africa Soil Information Service: Hydrologically Corrected/Adjusted SRTM DEM
795 (AfrHySRTM), International Center for Tropical Agriculture –Tropical Soil Biology and
796 Fertility Institute (CIAT-TSBF), World Agroforestry Centre (ICRAF), Center for
797 International Earth Science Information Network (CIESIN), Columbia University,
798 <https://cmr.earthdata.nasa.gov/search/concepts/C1214155420-SCIOPS>,2010.
- 799 Wadoux, A.M.J.C., Minasny, B., and McBratney, A.B.: Machine learning for digital soil mapping:
800 Applications, challenges and suggested solutions, *Earth Sci. Rev.*, 210, 103359, 2020.
- 801 Wright, M. N., and Ziegler, A.: Ranger: A fast implementation of random forests for high
802 dimensional data in C++ and R, *JSS*, 77(1), <https://doi.org/10.18637/jss.v077.i01>, 2017.
- 803 Water and Land Resource Center-Addis Ababa University (WLRC-AAU): Land use/land cover map
804 of Ethiopia, Addis Ababa, 2010.
- 805 Zewdie, E.: Properties of major Agricultural Soils of Ethiopia, Lambert Academic Publishing, 2013.
- 806 Zwedie, E.: Selected physical, chemical, and mineralogical characteristics of major soils occurring in
807 Chercher highlands, Eastern Ethiopia, *EJNRS*, 1(2), 173 – 185, 1999.