



HAL
open science

Model averaging for mapping topsoil organic carbon in France

Songchao Chen, Vera Leatitia Mulder, Gerard B.M. Heuvelink, Laura Poggio, Manon Caubet, Mercedes Roman Dobarco, Christian Walter, Dominique Arrouays

► **To cite this version:**

Songchao Chen, Vera Leatitia Mulder, Gerard B.M. Heuvelink, Laura Poggio, Manon Caubet, et al.. Model averaging for mapping topsoil organic carbon in France. *Geoderma*, 2020, 366, pp.114237. 10.1016/j.geoderma.2020.114237 . hal-02473703

HAL Id: hal-02473703

<https://hal.science/hal-02473703v1>

Submitted on 5 Jul 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

1 **Title:** Model averaging for mapping topsoil organic carbon in France

2

3 **Authors:**

4 Songchao Chen ^{a, b}. songchao.chen@inra.fr

5 Vera Leatitia Mulder ^c. titia.mulder@wur.nl

6 Gerard B.M. Heuvelink ^{c, d}. gerard.heuvelink@wur.nl

7 Laura Poggio ^d. laura.poggio@isric.org

8 Manon Caubet ^a. manon.caubet@inra.fr

9 Mercedes Román Dobarco ^a. mercedes.roman.dobarco@gmail.com

10 Christian Walter ^b. christian.walter@agrocampus-ouest.fr

11 Dominique Arrouays ^a. dominique.arrouays@inra.fr

12 **Affiliations:**

13 ^a INRAE, Unité InfoSol, 45075 Orléans, France

14 ^b SAS, INRAE, Agrocampus Ouest, 35042 Rennes, France

15 ^c Soil Geography and Landscape Group, Wageningen University, PO Box 47 6700

16 AA Wageningen, The Netherlands

17 ^d ISRIC–World Soil Information, PO Box 353 6700 AJ Wageningen, The Netherlands

18 **Corresponding author:**

19 Songchao Chen: songchao.chen@inra.fr

20 Postal address: INRAE, Unité InfoSol, 2163 Avenue de la Pomme de Pin, CS 40001

21 Ardon, 45075 Orléans, France

22 Telephone: +33(0)602142667

23 **Abstract:**

24 The soil organic carbon (SOC) pool is the largest terrestrial carbon (C) pool and is
25 two to three times larger than the C stored in vegetation and the atmosphere. SOC is
26 a crucial component within the C cycle, and an accurate baseline of SOC is required,
27 especially for biogeochemical and earth system modelling. This baseline will allow
28 better monitoring of SOC dynamics due to land use change and climate change.
29 However, current estimates of SOC stock and its spatial distribution have large
30 uncertainties. In this study, we test whether we can improve the accuracy of the three
31 existing SOC maps of France obtained at national (IGCS), continental (LUCAS), and
32 global (SoilGrids) scales using statistical model averaging approaches. Soil data from
33 the French Soil Monitoring Network (RMQS) were used to calibrate and evaluate five
34 model averaging approaches, i.e., Granger-Ramanathan, Bias-corrected Variance
35 Weighted (BC-VW), Bayesian Modelling Averaging, Cubist and Residual-based
36 Cubist. Cross-validation showed that with a calibration size larger than 100
37 observations, the five model averaging approaches performed better than individual
38 SOC maps. The BC-VW approach performed best and is recommended for model
39 averaging. Our results show that 200 calibration observations were an acceptable
40 calibration strategy for model averaging in France, showing that a fairly small number
41 of spatially stratified observations (sampling density of 1 sample per 2,500 km²)
42 provides sufficient calibration data. We also tested the use of model averaging in
43 data-poor situations by reproducing national SOC maps using various sized subsets
44 of the IGCS dataset for model calibration. The results show that model averaging
45 always performs better than the national SOC map. However, the Modelling
46 Efficiency dropped substantially when the national SOC map was excluded in model
47 averaging. This indicates the necessity of including a national SOC map for model

48 averaging, even if produced with a small dataset (i.e. 200 samples). This study
49 provides a reference for data-poor countries to improve national SOC maps using
50 existing continental and global SOC maps.

51

52 **Keywords:** Soil organic carbon; Digital soil mapping; Bias-corrected Variance
53 Weighted; Sample size requirement; Data-poor countries.

54 1. Introduction

55 Soils are crucial for maintaining ecosystem services such as food production,
56 water regulation, erosion control, biodiversity, and climate regulation (Sanchez et al.,
57 2009; Koch et al., 2013; Adhikari and Hartemink, 2016; Rumpel et al., 2018). To meet
58 the increasing demand for up-to-date and fine-resolution soil information, Digital Soil
59 Mapping (DSM, McBratney et al., 2003) has been widely adopted and is being rapidly
60 developed across different spatial scales since the past decade (e.g., Grunwald et
61 al., 2011; Poggio and Gimona. 2014; Viscarra Rossel et al., 2014; Hengl et al., 2015;
62 Ballabio et al., 2016; Padarian et al., 2017; Sanderman et al., 2018; Chen et al.,
63 2019;). At the global scale, different initiatives aim to deliver fine-resolution gridded
64 soil information. The main examples are the recent Global Soil Partnership GSOC
65 map (<http://54.229.242.119/GSOCmap/>), the *GlobalSoilMap* initiative (Sanchez et al.,
66 2009; Arrouays et al., 2014a), and SoilGrids products (Hengl et al., 2017). SoilGrids
67 adopts a “top-down” approach and produces soil property maps for the entire globe,
68 which are freely distributed and available online (<https://soilgrids.org/>). *GlobalSoilMap*
69 uses a “bottom-up approach” where each country produces soil property maps using
70 its own national soil data and defined specifications (e.g., 3 arc second resolution, six
71 standard depth intervals, quantified prediction uncertainty, Arrouays et al., 2014b).
72 Then, these country-level soil maps are merged into a global map. There are also
73 several initiatives producing soil property maps at the continental scale, such as
74 LUCAS (Tóth et al., 2013) for Europe and AfSIS (Hengl et al., 2015) for Africa. As a
75 result, there are often multiple maps available for a given soil property in a given area
76 produced using various soil databases, environmental covariates, and DSM methods.
77 Users may have multiple maps of the same property with different predictions and
78 different map accuracy which may lead to confusion regarding which map should be

79 used or whether the maps could or should be combined. It is possible to select the
80 most suitable soil property map for a specific region, when the map accuracy can be
81 evaluated using an independent validation dataset. When deciding to combine maps,
82 the hypothesis is that the information provided by the maps is complementary and
83 that a more accurate map may be obtained by merging the input maps using model
84 averaging approaches (Caubet et al., 2019). The model averaging option needs an
85 independent validation dataset and independent calibration data to train the model
86 averaging algorithm. Previous studies showed the potential of model averaging in
87 improving the accuracy of soil property maps of pH, soil texture, and available water
88 capacity (Malone et al., 2014; Padarian et al., 2014; Clifford and Guo, 2015; Román
89 Dobarco et al., 2017; Caubet et al., 2019).

90 The choice between selecting a single map and combining multiple maps is
91 not trivial, and many countries need to make this choice because of the increasing
92 number of different prediction maps of the same soil property. It is particularly
93 relevant to data-poor countries that may have very few or even no data to derive
94 reliable country-based maps, and that could benefit from collecting a limited number
95 of calibration samples to merge the national map with other existing products using
96 model averaging.

97 The objectives of this study are to 1) evaluate the added value of applying
98 model averaging in a data-rich country (e.g. France); 2) determine the most suitable
99 model averaging approach for improving the topsoil (0-20 cm) SOC map of mainland
100 France using three different SOC maps; 3) evaluate how well the model averaging
101 approaches perform for different calibration sizes and optimize the calibration size
102 required in model averaging; and 4) explore the potential of applying model
103 averaging in data-poor situations.

104

105 **2. Data**

106 In this study, we used three SOC maps generated and harmonized from
107 national, continental, and global DSM products and two national soil datasets in
108 France.

109

110 2.1. French national soil organic carbon maps

111 Numerous maps have been generated for France following the *GlobalSoilMap*
112 specifications. The most recent product (Mulder et al., 2016a) used all available point
113 data for France, both from the French Soil Mapping and Inventory Program
114 (Inventaire, Gestion et Conservation des Sols, IGCS) and an systematic grid aiming
115 at monitoring French soil properties (RMQS). More details about these two datasets
116 can be found in the study of Mulder et al. (2016a). For this study, we used the same
117 *GlobalSoilMap* approach as Mulder et al. (2016a), but we set aside the RMQS grid to
118 be used as an independent dataset for calibrating the model averaging algorithms
119 and evaluating map accuracy (see Sections 2.3 and 2.4). A total of 30,381 soil
120 profiles from the IGCS dataset were used to generate SOC maps at the first three
121 *GlobalSoilMap* depth intervals (0-5, 5-15, 15-30 cm). The IGCS dataset is a
122 compilation of soil profiles from many programs that mostly focused on agricultural
123 soils. As a result, the soil profile density is high in some regions (Fig. 1), whereas it is
124 low in other regions; some land uses are over- or under-represented in the calibration
125 dataset. SOC contents at the *GlobalSoilMap* depth intervals were obtained by
126 applying equal area quadratic splines (Bishop et al., 1999; Malone et al., 2009) to soil
127 profile data, as outlined in Mulder et al. (2016b). Spatially exhaustive covariates,
128 including climate zones and meteorological data, vegetation, topography, geology,

129 soils, and land management, were resampled to 90 m resolution. Details about these
130 environmental covariates are given in Mulder et al. (2016a). In this study, the national
131 SOC map (named IGCS SOC map hereafter) for the topsoil (0-20 cm) was calculated
132 from SOC maps of 0-5, 5-15, and 15-30 cm by a weighted averaging approach,
133 where the weights are proportional to the layer thickness (Fig. 2a).

134

135 2.2. Continental and global scale soil organic carbon maps

136 In addition to the aforementioned national SOC map, we also obtained SOC
137 maps for France from continental (LUCAS) and global (SoilGrids) soil map products.

138 The LUCAS SOC map (Fig. 2b) contains SOC predictions for the topsoil (0-20
139 cm) at 1 km resolution for Europe (Aksoy et al., 2016). A total of 23,835 soil samples
140 were used for model calibration. These soil samples were collected from LUCAS
141 (19,860 samples), BioSoil (3,379 plots from forest soil), and SoilTrEC (387 samples
142 from local soil data from six different critical zone observatories in Europe). From
143 these datasets, about 3,500 sites were located in France. A regression kriging model
144 was fitted to generate a SOC map using observed SOC content and 15
145 environmental covariates.

146 The SoilGrids SOC map (<https://soilgrids.org>, v0.5.3, Fig. 2c) was extracted
147 from the study of Hengl et al. (2017), in which SOC was mapped at seven standard
148 depths (0, 5, 15, 30, 60, 100, and 200 cm) at a resolution of 250 m for the globe.
149 These SOC maps were based on about 150,000 soil profiles along with 158 remote
150 sensing-based soil covariates. Maps were produced by fitting an ensemble prediction
151 from random forest and gradient boosting trees. From the 150,000 soil profiles,
152 nearly 3,000 were located in mainland France, mainly originating from the LUCAS
153 database. For this work, the topsoil SOC map was calculated from SoilGrids SOC

154 maps at 0, 5, 15, and 30 cm depth using trapezoidal numerical integration (Hengl et
155 al., 2017).

156 The LUCAS and SoilGrids SOC maps were resampled to 90 m using bilinear
157 interpolation and reprojected to the Lambert 93 coordinate system to match these
158 with the national SOC map.

159

160 2.3. Independent soil data for model averaging calibration and SOC map validation

161 To evaluate the accuracy of the input and merged maps, an independent
162 validation dataset and an independent dataset for calibration of the model averaging
163 algorithm were needed. These datasets were derived from the RMQS French
164 systematic grid, which covers different soil, climate, relief, and land cover conditions
165 (Fig. 1). The RMQS dataset is a 16 km × 16 km square grid where sampling sites are
166 at the centre of each grid cell, covering mainland France (Jolivet et al., 2006). For
167 each site, 25 individual core samples were collected by a hand auger and mixed into
168 a composite sample, both for 0–30 cm and 30–50 cm depth intervals. For more
169 detailed information about the soil sampling design and laboratory analyses, refer to
170 Martin et al. (2009). Because there were no SOC measurements for a depth of 0-20
171 cm for the RMQS sites, we calculated these values depending on land use: 1) for
172 most agricultural soils, SOC concentration decreases at a small rate with depth in the
173 topsoil because of ploughing; thus, SOC content at 0-20 cm is close to that of 0-30
174 cm (Arrouays et al., 2001). We therefore used SOC at 0-30 cm to represent the SOC
175 at 0-20 cm for RMQS sites under agricultural soils; 2) for natural soils (grassland and
176 forest), SOC usually decreases with depth in the topsoil. Therefore, we first
177 calculated SOC at 0-20 cm and at 0-30 cm by equal area quadratic splines using
178 5785 grassland and forest soil profiles from the IGCS dataset. We then fitted a linear

179 model between SOC at 0-20 cm and SOC at 0-30 cm ($SOC_{0-20\text{ cm}} = 1.04 \times SOC_{0-30\text{ cm}} + 0.26$, $R^2 = 0.986$). We used this model to derive SOC at 0-20 cm from SOC at 0-
180 +0.26, $R^2 = 0.986$). We used this model to derive SOC at 0-20 cm from SOC at 0-
181 30 cm for all RMQS sites under natural soils.

182

183 **3. Methods**

184 3.1. Generic framework for model averaging

185 Fig. 3 shows the generic framework for model averaging, which includes four
186 steps. We first explain the procedure used for selecting the calibration and validation
187 subsets from the RMQS dataset. To obtain spatially representative calibration and
188 validation datasets, equal-size clustering (iterative nearest neighbour approach,
189 Monlong, 2018) was applied to the RMQS sites (Step 1), which resulted in spatially
190 compact clusters. This was done for five cluster sample sizes (4, 10, 20, 50, and
191 100). Note that the cluster sample size is only approximately the same for all clusters
192 because the total number of observations (i.e., 1996) is not always a multiple of the
193 cluster sample size. Fig. 4 shows the spatial distribution of the clusters. In Step 2, a
194 k -fold cross-validation framework ($k = 4, 10, 20, 50, 100$) was used to separate a
195 calibration set by randomly allocating one observation per cluster to each fold. Thus,
196 the sample size of each fold was approximately 500, 200, 100, 40, and 20, for $k=4$,
197 10, 20, 50 and 100, respectively. In each of the k times, one of the folds was used to
198 calibrate the model averaging approaches (Step 3), whereas the remaining $k-1$ folds
199 were used for model validation (Step 4, as explained in Section 3.2). By performing
200 this analysis for different values of k , we could also evaluate the performance of the
201 model averaging approaches for different calibration sizes (i.e. 500, 200, 100, 40,
202 and 20). Note that the cross-validation procedure used here has some similarities
203 with spatial cross-validation (Roberts et al., 2017).

204

205 3.2. Model averaging approaches

206 Five model averaging approaches were compared in this study. They are
207 Granger-Ramanathan (Granger and Ramanathan, 1984), Variance Weighted (Bates
208 and Granger, 1969; Heuvelink and Bierkens, 1992), Bayesian model averaging
209 (Hoeting et al., 1999), Piecewise linear decision tree (Quinlan, 1992), and Residual-
210 based piecewise linear decision tree.

211

212 3.2.1. Granger-Ramanathan

213 The Granger-Ramanathan (GR) approach was proposed by Granger and
214 Ramanathan (1984). It assumes that a combination of different model predictions can
215 be approached using a traditional Ordinary Least Square (OLS) method. In our case,
216 a linear regression model was fitted between the measured SOC contents of the
217 calibration set and the SOC predictions of the three SOC maps. The outcome SOC_{GR}
218 from the GR approach can be calculated as

$$219 \quad SOC_{GR} = \sum_{i=1}^p (\alpha_i \cdot SOC_i) + \beta \quad (1)$$

220 where α_i and SOC_i are the regression coefficient and SOC prediction of the i -th SOC
221 map ($p=3$ in this study), and β is the intercept. The α and β coefficients are solved by
222 the OLS method, and the sum of the α_i is not necessarily equal to 1.

223

224 3.2.2. Variance Weighted

225 We used the Bias-corrected Variance Weighted (BC-VW) approach from Ge et
226 al. (2014), which is based on the error variance-covariance matrix that is estimated
227 by comparing model predictions with observations. Thus, the outcome SOC_{BC-VW} is
228 calculated as

229
$$SOC_{BC-VW} = \sum_{i=1}^p \alpha_i \cdot (SOC_i - \beta_i) \quad (2)$$

230 where α_i and SOC_i are the weight and SOC prediction of SOC map i , respectively,
 231 and β_i is the bias correction coefficient for SOC map i . The latter is calculated as

232
$$\beta_i = \frac{1}{m} \sum_{k=1}^m (SOC_{i,k} - SOC_{obs,k}) \quad (3)$$

233 where m is the number of calibration observations, and $SOC_{i,k}$ and $SOC_{obs,k}$ are the
 234 SOC prediction of SOC map i and the SOC observation at the k -th calibration site,
 235 respectively.

236 As described in Ge et al. (2014), the vector $\alpha = [\alpha_1 \cdots \alpha_p]^T$ is calculated by
 237 minimizing the error variance of the model predictions:

238
$$\alpha^T = (\mathbf{1}^T \mathbf{V}^{-1} \mathbf{1})^{-1} \mathbf{1}^T \mathbf{V}^{-1} \quad (4)$$

239 where $\mathbf{1}$ is the p -dimensional identity matrix (recall that $p=3$ in this study), and \mathbf{V} is
 240 the p -dimensional variance-covariance matrix of the prediction error. The elements of
 241 \mathbf{V} are determined as

242
$$\hat{v}_{ij} = \frac{1}{m} \sum_{k=1}^m (SOC_{i,k} - SOC_{obs,k})(SOC_{j,k} - SOC_{obs,k}) \quad (5)$$

243 where $i, j = 1, \dots, n$ represent SOC maps, and m is the number of calibration
 244 observations. Note that the correlations between SOC map errors are considered in
 245 the BC-VW approach.

246

247 3.2.3. Bayesian Model Averaging

248 The Bayesian Model Averaging (BMA) approach assigns a conditional
 249 probability density function (PDF) to each model prediction (Hoeting et al., 1999). The
 250 BMA posterior distribution of the final output (SOC_{BMA}) can be expressed as (Raftery
 251 et al., 2005):

252
$$p(SOC_{BMA}|SOC_{obs}) = \sum_{i=1}^p p(SOC_{BMA}|SOC_{obs}, SOC_i)p(SOC_i|SOC_{obs}) \quad (6)$$

253 where SOC_{obs} are the SOC observations, p is the number of SOC maps (in this study
 254 $p=3$), and SOC_i denote the values of SOC extracted from the SOC map i at the
 255 locations of observations. Therefore, the BMA posterior distribution of SOC_{BMA} is a
 256 weighted average of the posterior distributions of SOC_{BMA} under each of the SOC
 257 maps, weighted by their posterior model probabilities.

258 The posterior model probability of SOC_i is expressed as (Raftery et al., 2005)

259
$$p(SOC_i|SOC_{obs}) = \frac{p(SOC_{obs}|SOC_i)p(SOC_i)}{\sum_{l=1}^p p(SOC_{obs}|SOC_l)p(SOC_l)} \quad (7)$$

260 where $p(SOC_{obs}|SOC_i)$ is the integrated likelihood of SOC_i , and it can be calculated
 261 by BIC approximation (more details can be found in Raftery et al., 2005).

262 We used the R package “BMA” (Raftery et al., 2005) to apply BMA in our case
 263 study.

264

265 3.2.4. Piecewise linear decision tree

266 The Piecewise linear decision tree approach (Cubist) is based on the M5
 267 algorithm (Quinlan, 1992). It partitions the dataset into several subsets within which
 268 inputs (independent variables) are similar. In a given subset, the standard deviation
 269 of the target values is treated as a measure of error and is used as a node splitting
 270 criterion. Every potential split is evaluated by the reduction in standard deviation.
 271 After evaluating all possible splits, Cubist chooses the one split that maximizes the
 272 reduction in error. Then, pruning and smoothing processes are performed to get the
 273 final model. More details are given in Quinlan (1992).

274 In the final Cubist model, partitions are defined by a list of rules, which are
 275 arranged in a hierarchy. Each rule has the following form:

276 **if** [condition] **then** [linear regression model]

277 **else** [apply next rule].

278 A rule indicates that whenever a case satisfies the condition of one rule, the
279 corresponding linear regression model is used to predict the output. In this study, we
280 used the R package “Cubist” (Kuhn et al., 2012).

281 3.2.5. Residual-based piecewise linear decision tree

282 The framework of Residual-based piecewise linear decision tree (Residual-
283 based Cubist, revised from Tao et al., 2018) is as follows: 1) calculate the arithmetic
284 mean SOC value (SOC_{mean}) extracted from IGCS (SOC_{IGCS}), LUCAS (SOC_{LUCAS}),
285 and SoilGrids ($SOC_{\text{SoilGrids}}$) SOC maps at locations of soil observations; 2) calculate
286 the residuals (RES_{IGCS} , RES_{LUCAS} , and $RES_{\text{SoilGrids}}$) between SOC_{mean} and
287 $SOC_{\text{IGCS}}/SOC_{\text{LUCAS}}/SOC_{\text{SoilGrids}}$, which are used as predictors in the Cubist model; 3)
288 calculate the residuals (RES_{obs}) between SOC_{mean} and SOC observations (SOC_{obs}),
289 which are used as the target variable in the Cubist model ; and 4) once the Cubist
290 model is fitted, calculate the final SOC predictions of the Residual-based Cubist by
291 summing up the RES_{obs} (derived from Cubist) and SOC_{mean} .

292

293 3.3. Evaluation of three SOC maps and five model averaging approaches using
294 different calibration sizes

295 The performance of three individual soil SOC maps was assessed using all
296 RMQS data. Based on a k -fold cross-validation framework explained in Section 3.1,
297 we evaluated the five model averaging approaches using different calibration sample
298 sizes (from 20 to 500). Three indicators, the Modelling Efficiency (ME), the Root
299 Mean Square Error (RMSE), and Bias, were used to evaluate prediction accuracy.

300
$$ME = 1 - \frac{\sum_{i=1}^n (\hat{z}_i - z_i)^2}{\sum_{i=1}^n (z_i - \bar{z})^2} \quad (8)$$

301
$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{z}_i - z_i)^2} \quad (9)$$

302
$$\text{Bias} = \frac{1}{n} \sum_{i=1}^n (\hat{z}_i - z_i) \quad (10)$$

303 where n is the size of the cross-validation dataset, z_i and \hat{z}_i are measured and
304 predicted values for the i -th observation in the cross-validation dataset, respectively,
305 and \bar{z} is the mean of the observations in the cross-validation dataset. A negative ME
306 means that the model performs worse than using the average of the observations as
307 a prediction.

308

309 3.4. The effect of national SOC maps on model averaging

310 The IGCS map was generated using the entire IGCS dataset (about 30,000
311 soil profiles), which is very large and hence is an example of a case study in a data-
312 rich country (1 profile per 18 km²). To assess the usefulness of model averaging in
313 data-poor situations, we applied model averaging to a case in which the national
314 SOC map (IGCS) was generated from a much smaller number of soil profiles. To do
315 so, we generated IGCS SOC maps by randomly selecting 10,000, 5,000, 1,000, 800,
316 600, 400, and 200 soil profiles from the whole IGCS dataset. To filter out random
317 sampling effects, we repeated this procedure 100 times for each sample size and
318 reported the average results. These IGCS SOC maps with LUCAS and SoilGrids
319 were finally merged only with the best model averaging approach and using the
320 minimum necessary number of calibration sites as previously estimated. Using the
321 same minimum necessary number of calibration sites, we tested the assumption of
322 SoilGrids and LUCAS providing additional information that is not captured in IGCS
323 SOC map by removing these two SOC maps in model averaging and only using GR

324 approach to calibrate the generated IGCS SOC maps (using 200 to 10,000 soil
325 profiles). We also tested model averaging using only SoilGrids and LUCAS to test the
326 assumption that no national SOC map was available.

327

328 **4. Results**

329 4.1. Summary of IGCS, RMQS, and LUCAS datasets

330 Table 1 summarises SOC statistics of the IGCS, RMQS, and LUCAS (located
331 in France) datasets. About 80% (24,596) of IGCS soil profiles were located in arable
332 soils, and 20% (5,785) were located in forest and permanent grassland soils. In the
333 IGCS soil database, grassland and forest soils (mean SOC of 24.88 g kg⁻¹) had
334 higher SOC values than arable soils (mean SOC of 16.66 g kg⁻¹). Nearly half (985) of
335 the RMQS sampling sites were located in permanent grasslands or forest soils, and
336 the remaining half (1011) were under arable soils. In the RMQS dataset, the mean
337 SOC was 18.19 g kg⁻¹ for arable soils and 35.51 g kg⁻¹ for permanent grassland and
338 forest soils. LUCAS observations had a mean SOC of 26.20 g kg⁻¹ for permanent
339 grassland and arable soils.

340

341 4.2. Evaluation of SOC maps from IGCS, LUCAS, and SoilGrids datasets

342 The IGCS SOC map has the lowest RMSE (18.86 g kg⁻¹) and highest ME
343 (0.25) among the three SOC maps (Fig. 5). The negative Bias (-6.17 g kg⁻¹) indicates
344 that SOC is underestimated in the IGCS SOC map. When the performance of the
345 IGCS SOC map for arable and forest/grassland soils was separately evaluated,
346 arable soils (ME of 0.19 and RMSE of 10.02 g kg⁻¹) were found to have higher
347 accuracy than forest/grassland soils (ME of 0.09 and RMSE of 24.85 g kg⁻¹). SOC
348 maps of LUCAS and SoilGrids have a much higher RMSE of 30.62 and 32.75 g kg⁻¹,

349 and a negative ME of -1.18 and -1.27, respectively. Positive Bias of LUCAS (6.73 g
350 kg⁻¹) and SoilGrids (21.81 g kg⁻¹) showed that these two maps overestimated SOC.
351 The overestimation was larger in SoilGrids than in the LUCAS SOC map.

352

353 4.3. Comparison of five model averaging approaches using different calibration sizes

354 The BC-VW approach performed best among the five model averaging
355 approaches across different calibration sizes, with the lowest RMSE (16.77-18.71 g
356 kg⁻¹) and highest ME (0.23-0.38) (Fig. 6). The GR and BMA ranked second and third
357 when the calibration size was large (100, 200 or 500), with an ME between 0.33 and
358 0.38. The performance of GR substantially decreased when using a calibration
359 sample size of 40 and 20, whereas BMA was more stable (and ranked third) when
360 using a small calibration sample size. Cubist performed worst in the case of a large
361 calibration sample size (100, 200, or 500) but ranked second when the calibration
362 sample size was small (20 or 40). Residual-based Cubist did not perform well across
363 the different calibration sample sizes. It should be noted that BC-VW, GR, and BMA
364 had a Bias close to 0 under different calibration sample sizes, while Cubist and
365 Residual-based Cubist had a large negative ME.

366 All model averaging approaches showed better performance metrics than
367 using the individual LUCAS and SoilGrids SOC maps for all calibration sample sizes.
368 Improvement on the IGCS SOC map only occurred when the calibration sample size
369 was large (100, 200, or 500), while the model averaging approaches performed
370 worse than the IGCS SOC map when the calibration sample size was 20 or 40.

371 In general, the model performance of the five model averaging approaches
372 declined when the calibration size decreased (Fig. 6). Being the best model
373 averaging approach, BC-VW had better performance than the IGCS SOC map when

374 calibration samples were 500, 200, and 100, and it was still slightly better when only
375 40 calibration samples were used. However, 20 calibration samples were not
376 sufficient to improve SOC maps using any of the five model averaging approaches.
377 GR and BMA could improve SOC predictions when calibration sample sizes were
378 500, 200, and 100. However, Cubist and Residual-based Cubist only performed
379 better than the IGCS SOC map when using a calibration sample size of 200 or more.

380 As shown in Fig. 6, only slight differences (ME of 0.37-0.38, and RMSE of
381 16.77-16.90 g kg⁻¹) were observed between 500 and 200 calibration sample sizes
382 when using BC-VW, which was the best model averaging approach. Nevertheless,
383 the model performance of BC-VW showed a steady decline when the calibration
384 sample size decreased from 200 to 20.

385

386 4.4. SOC maps using five model averaging approaches

387 Fig. 7 shows SOC maps obtained from the five model averaging approaches
388 using all RMQS data for calibration. The general spatial patterns of these five SOC
389 maps were quite close, which is consistent with their similar model performance (in
390 the case of a 500 calibration sample size) in Fig. 6. In comparison with the IGCS
391 SOC map (Fig. 2a), these five SOC maps have higher SOC in mountainous regions
392 (e.g., the Alps, the Central Massif, the Pyrenees), forests, and grasslands (e.g., the
393 Landes of Gascony, western Brittany). As shown in Fig. 7f to Fig. 7o, SOC maps
394 derived from GR, BC-VW, and BMA had slightly higher SOC contents than Cubist
395 and Residual-based Cubist. This is particularly visible in Fig. 7k to Fig. 7o, which
396 zooms in on a square area in the Landes of Gascony forest.

397

398 4.5. Influence of national SOC maps on model averaging performance

399 The performance (ME and RMSE) of the IGCS SOC maps derived from
400 different sample sizes showed a slight decline when the number of soil profiles used
401 decreased from 10,000 to 800 (Fig. 8). A stronger decline in performance was
402 observed when the number of soil profiles decreased further from 800 to 200, with
403 ME values dropping from 0.23 to 0.16 and RMSE increasing from 19.11g kg⁻¹ to
404 19.89 g kg⁻¹. The performance of the BC-VW approach on the three SOC maps and
405 the GR approach only on IGCS SOC map showed similar declining trends as the
406 IGCS SOC maps. However, the BC-VW maps always performed better than the
407 IGCS maps (Δ ME > 0.1 and Δ RMSE < -2 g kg⁻¹) and GR maps (Δ ME > 0.04 and
408 Δ RMSE < -1 g kg⁻¹). When using only LUCAS and SoilGrids for model averaging,
409 BC-VW performed much worse than all other SOC maps produced using IGCS,
410 LUCAS, and SoilGrids in model averaging, with a ME of -0.24 and a large RMSE of
411 23.65 g kg⁻¹.

412

413 **5. Discussion**

414 5.1. Performance evaluation of SOC maps from IGCS, LUCAS, and SoilGrids

415 The IGCS SOC map had the best performance indicators among the three
416 source SOC products. However, it showed a slight overall underestimation and a
417 clear tendency to underestimate large SOC values. This may be because the
418 calibration data for generating the IGCS SOC map are dominated by cultivated soils
419 (80% of IGCS dataset), which typically have low SOC values because of
420 management practices (Table 1). As natural soils occupy 45% of the total area of
421 mainland France (Chen et al., 2018), high SOC values are under-represented in the
422 dataset for producing the IGCS SOC map. It consequently resulted in

423 underestimating the effect of some controlling factors driving high SOC values (e.g.,
424 forest or grassland land uses, high elevations). Although the effects of land use and
425 elevation are still clearly visible (Fig. 2a), the spatial patterns of the resulting map are
426 too smooth, as was already described by Mulder et al. (2016a; 2016b). In the French
427 *GlobalSoilMap* product, Mulder et al. (2016a) produced national SOC maps at the
428 first three depth intervals (0-5, 5-15, and 15-30 cm) using both IGCS and RMQS
429 data. The ME evaluated using 10-fold cross-validation ranged from 0.26 to 0.36 for
430 the first three depth intervals. This shows that including RMQS data into national
431 SOC modelling improves model performance. Nevertheless, SOC was still slightly
432 underestimated because the IGCS dataset is almost 15 times larger than the RMQS
433 dataset and IGCS data generally have low SOC content (Table 1).

434 The predictive performance of the LUCAS map and SoilGrids map was much
435 worse than that of the IGCS map, as illustrated in Fig. 2. They both have a tendency
436 to overestimate SOC, either slightly (LUCAS) or largely (SoilGrids). The LUCAS map
437 also exhibited more contrasted and irregular patterns than the IGCS map. Moreover,
438 the LUCAS map showed some areas with artificially rounded boundaries (mainly in
439 southwest France), suggesting a bias linked to the environmental covariates,
440 predictive model, and/or interpolation method used. The SoilGrids map clearly
441 overestimated SOC for the large majority of situations (Fig. 5). It also clearly missed
442 the effect of some land use types on decreasing SOC (e.g., intensively cultivated
443 plains in northern and southwestern parts of France, vineyards in southern France).
444 This suggests that the covariates used for global modelling could not capture these
445 effects; e.g., land use/land cover classes used as covariates for SoilGrids were
446 limited to cultivated land, forests, grasslands, shrublands, wetlands, tundra, artificial
447 surfaces, and bare land cover.

448 Homogenising data to a common depth of 0-20 cm may have induced some
449 additional uncertainty (Laborczi et al., 2018). We also acknowledge that resampling
450 SoilGrids and LUCAS to 90 m resolution may have added a source of discretionality
451 and potential uncertainty.

452

453 5.2. Potential and limitations of model averaging approaches

454 Our results demonstrate the ability of model averaging approaches to improve
455 national SOC maps (Fig. 5 and Fig. 6). The improvement strongly depends on the
456 calibration sample size used for model averaging. It is encouraging that 200 spatially
457 stratified samples (1 sample per 2,500 km²) were enough for producing a sufficiently
458 accurate national SOC map (ME of 0.37 for BC-VW approach) when applying model
459 averaging in France. Note also that the performance of this SOC map is comparable
460 to that of the *GlobalSoilMap* SOC map using IGCS and RMQS datasets (Mulder et
461 al., 2016a).

462 We note that we did not map the uncertainty of SOC predictions when
463 applying model averaging. Prediction uncertainty should be considered in future
464 studies because it is crucial for assessing model quality and robustness. It is also a
465 strongly recommended product outcome, as indicated in the *GlobalSoilMap*
466 specifications (e.g., Arrouays et al., 2014a; Heuvelink, 2014). We could use the
467 method proposed by Ge et al. (2014) to estimate uncertainty when using BC-VW for
468 merging multiple SOC maps.

469 In addition to deriving SOC predictions using model averaging, it would be
470 beneficial to also explicitly quantify the uncertainties associated with these
471 predictions. This can be done using uncertainty propagation techniques such as the
472 Taylor series method and Monte Carlo simulation (Heuvelink, 2018; Román Dobarco

473 et al., 2019) provided that the uncertainties of the input maps and their correlations
474 are quantified. This may be a useful extension of the work presented here. If it is
475 done, it would be useful to also evaluate the validity of the uncertainty maps by
476 computing statistics of the standardised squared prediction error (Lark, 2000) and
477 accuracy plots (Goovaerts, 2001; Wadoux et al., 2018).

478

479 5.3. Comparison with previous model averaging studies

480 Our results suggest that map performance improves when using model
481 averaging approaches and that the BC-VW method is the best approach for SOC
482 mapping in mainland France. Previous studies also showed that model averaging
483 improves map predictions, but different approaches tend to have similar performance
484 (e.g., Malone et al., 2014; Román Dobarco et al., 2017; Caubet et al., 2019). Caubet
485 et al. (2019) applied two model averaging approaches (GR and BC-VW) to improve
486 soil texture maps (clay and sand) and showed that both model averaging approaches
487 improved the accuracy and that GR outperformed BC-VW. Similar results were found
488 by Román Dobarco et al. (2017) for mapping soil texture, and Malone et al. (2014) on
489 pH mapping. Indeed, the best-performing algorithm for model averaging may vary
490 between study areas and for different soil properties, and thus optimization of model
491 averaging methods is case-specific.

492 Caubet et al. (2019) also mentioned the potential use of non-linear models for
493 improving model averaging. However, in our study, non-linear models like Cubist and
494 Residual-based Cubist did not perform better than a linear model like GR. Perhaps
495 this is because three SOC products are not sufficient for calibrating a regression tree
496 or machine learning approach, and that other additional covariates (e.g., elevation,
497 land use, and climatic variables) may be helpful to improve model performance.

498 Especially, the example of the Landes of Gascony (see Fig. 7k to Fig. 7o) shows that
499 the model does not capture the effect of forest land use well in many areas when
500 using a rule-based model such as Cubist.

501 Caubet et al. (2019) found that around 200 to 300 calibration samples were
502 sufficient for model averaging of soil texture over mainland France. This result is
503 consistent with our finding that 200 calibration samples (1 sample per 2,500 km² for a
504 total area of 550,000 km² and a country having a high pedodiversity (Minasny et al.,
505 2010)) selected from equal-size clustering are enough to improve existing SOC maps
506 using model averaging. In our case, it is promising that adding rather few samples
507 improves the SOC maps considerably. This suggests that adding some soil
508 observations uniformly spread over the geographic space helps to correct the bias of
509 the original maps.

510

511 5.4. Contribution of model averaging approaches to data-poor countries

512 We tested model averaging on a situation that may be considered “rich”
513 concerning the amount of available data (Arrouays et al., 2017). In this study, we
514 used 30,000 samples for national SOC mapping, which is 1 sample per 18 km².
515 Although France has numerous point soil data, these data are rather clustered and
516 irregularly cover the territory. They also over-represent some agro-pedoclimatic
517 conditions (e.g., low elevations and intensively cultivated areas). These conditions
518 (irregularity and non-representativeness of samples) are likely to be similar in most
519 data-rich countries that use legacy data for DSM.

520 The fact that the number of samples needed to calibrate the averaging model
521 is rather low is encouraging, i.e. 200 samples for mainland France. This is cost-

522 effective given the limited effort required to gather a fairly small number of soil
523 samples to improve national soil maps.

524 The results shown in Fig. 8 indicate that model averaging always has a
525 substantial added value in terms of model performance compared to using the IGCS
526 SOC map alone. Surprisingly, the added value of model averaging does not depend
527 on the sample size (200 to 10,000 samples) used for producing the national map.
528 This might be due to the fact that our calibration sample size for model averaging (i.e.
529 200 spatially stratified observations) is large enough to capture the main variations of
530 SOC in mainland France. The results shown in Figure 8 also show that removing the
531 LUCAS and SoilGrids SOC maps (GR) decreases the map accuracy in model
532 averaging (BC-VW) which implies that these two SOC maps are complementary to
533 the IGCS SOC map for model averaging. Moreover, the added value of model
534 averaging is larger than that of only increasing the number of profiles used for
535 producing the IGCS SOC map. For example, using 200 samples for model averaging
536 calibration results in an ME increase of 0.12, whereas the ME only increases by 0.07
537 when the number of profiles used for producing the IGCS SOC map increases from
538 200 to 10,000. This indicates that adding a relatively small regular grid of soil
539 samples to merge several maps might be more efficient than expanding the database
540 with a large number of soil samples for which the sample locations are not controlled.
541 In many countries, soil mapping activities are frequently guided by local needs and
542 interests. This explains why national soil datasets are often clustered and why adding
543 more legacy data may sometimes lead to increasing sources of bias (e.g., Poggio et
544 al., 2019). Overall, our study advocates merging predictions in both data-rich and
545 data-poor situations and demonstrates that the added value of merging is relatively
546 higher in data-poor situations. However, notably, the performance of BC-VW drops

547 substantially when excluding the IGCS SOC map and when it only uses LUCAS and
548 SoilGrids for model averaging. This indicates the importance of a national SOC map
549 in model averaging, even if this SOC map is produced with a small dataset (i.e. 200
550 samples).

551

552 **6. Conclusion**

553 We tested the ability of five model averaging approaches for improving
554 existing SOC maps by merging national, continental, and global SOC products. All
555 five model averaging approaches could improve the national SOC map when more
556 than 100 soil samples were used for calibration of the model averaging approaches.
557 The BC-VW approach performed better than the other four approaches. Model
558 averaging approaches using a rather small calibration dataset (i.e. 200 observations
559 uniformly spread over mainland France) for calibration proved to be efficient. The
560 national SOC map was very important and drove performance when merging all SOC
561 maps, however SoilGrids and LUCAS SOC maps had added value by capturing
562 relevant patterns additional to the national SOC map. By reducing the number of
563 national soil samples in France for producing the national SOC map, we found that
564 merging maps using model averaging is also applicable to data-poor situations and
565 might thus be attractive to data-poor countries, provided sufficient soil data are
566 available for calibration of the model averaging approach.

567

568 **Acknowledgements**

569 Soil data collection was supported by the French Scientific Group of Interest
570 on soils: the GIS Sol, involving the French Ministry of Ecology, the French Ministry of
571 Agriculture, the French Environment and Energy Management Agency (ADEME), the

572 French Institute for Research and Development (IRD), the French National
573 Geographic and Forest Inventory Institute (IGN) and the French National Institute for
574 Agronomic Research (INRA). This work was partly funded by the Environment &
575 Agronomy Department of INRA, in the framework of its calls for innovative research
576 (grant no. 6282), and partly funded by the project Coordination of International
577 Research Cooperation on soil CARbon Sequestration in Agriculture (CIRCASA) (grant
578 no. 774378) under H2020-EU.3.2.1.1. Vera Laetitia Mulder, Laura Poggio and
579 Dominique Arrouays are members of a Research Consortium supported by LE
580 STUDIUM Loire Valley Institute for Advanced studies. We also thank all colleagues
581 involved in soil sampling and populating the soil database. Songchao Chen received
582 support from the China Scholarship Council (grant no. 201606320211).

583

584 **References**

- 585 Adhikari, K., Hartemink, A.E., 2016. Linking soils to ecosystem services—A global
586 review. *Geoderma* 262, 101–111.
- 587 Aksoy, E., Yigini, Y., Montanarella, L., 2016. Combining soil databases for topsoil
588 organic carbon mapping in Europe. *PloS One* 11, e0152098.
- 589 Arrouays, D., Deslais, W., Badeau, V., 2001. The carbon content of topsoil and its
590 geographical distribution in France. *Soil Use Manage.* 17, 7–11.
- 591 Arrouays, D., Grundy, M.G., Hartemink, A.E., Hempel, J.W., Heuvelink, G.B.M.,
592 Hong, S.Y., Lagacherie, P., Lelyk, G., McBratney, A.B., McKenzie, N.J.,
593 Mendonca-Santos, M.d.L., Minasny, B., Montanarella, L., Odeh, I.O.A., Sanchez,
594 P.A., Thompson, J.A., Zhang, G.-L., 2014a. Chapter Three — GlobalSoilMap:
595 Toward a Fine-Resolution Global Grid of Soil Properties. *Adv. Agron.* 125, 93–
596 134.
- 597 Arrouays, D., Leenaars, J.G.B., Richer-de-Forges, A.C., Adhikari, K., Ballabio, C.,
598 Greve, M., Grundy, M., Guerrero, E., Hempel, J.W., Hengl, T., Heuvelink,
599 G.B.M., et al., 2017. Soil legacy data rescue via GlobalSoilMap and other
600 international and national initiatives. *GeoResJ* 14, 1–19.
- 601 Arrouays, D., McKenzie, N.J., Hempel, J., Richer de Forges, A.C., McBratney, A.B.,
602 2014b. *GlobalSoilMap: Basis of the Global Spatial Soil Information System*. 1st
603 ed. CRC Press Taylor & Francis Group, pp. 9–12.
- 604 Ballabio, C., Panagos, P., Monatanarella, L., 2016. Mapping topsoil physical
605 properties at European scale using the LUCAS database. *Geoderma* 261, 110–
606 123.
- 607 Bates, J.M., Granger, C.W., 1969. The combination of forecasts. *Oper. Res. Soc.* 20,
608 451–468.

609 Bishop, T.F.A., McBratney, A.B., Laslett, G.M., 1999. Modelling soil attribute depth
610 functions with equal-area quadratic smoothing splines. *Geoderma* 91, 27–45.

611 Caubet, M., Román Dobarco, M., Arrouays, D., Minasny, B., Saby, N.P.A., 2019.
612 Merging country, continental and global predictions of soil texture: Lessons from
613 ensemble modelling in France. *Geoderma* 337, 99–110.

614 Chen, S., Arrouays, D., Angers, D.A., Chenu, C., Barré, P., Martin, M.P., Saby,
615 N.P.A., Walter, C., 2019. National estimation of soil organic carbon storage
616 potential for arable soils: A data-driven approach coupled with carbon-landscape
617 zones. *Sci. Total Environ.* 666, 355–367.

618 Chen, S., Martin, M.P., Saby, N.P.A., Walter, C., Angers, D.A., Arrouays, D., 2018.
619 Fine resolution map of top-and subsoil carbon sequestration potential in France.
620 *Sci. Total Environ.* 630, 389–400.

621 Clifford, D., Guo, Y., 2015. Combining two soil property rasters using an adaptive
622 gating approach. *Soil Res.* 53, 907–912.

623 Ge, Y., Avitabile, V., Heuvelink, G.B.M., Wang, J., Herold, M., 2014. Fusion of pan-
624 tropical biomass maps using weighted averaging and regional calibration data.
625 *Int. J. Appl. Earth Obs.* 31, 13–24.

626 Goovaerts, P., 2001. Geostatistical modelling of uncertainty in soil science.
627 *Geoderma* 103, 3–26.

628 Granger, C.W., Ramanathan, R., 1984. Improved methods of combining forecasts. *J.*
629 *Forecasting* 3, 197–204.

630 Grunwald, S., Thompson, J.A., Boettinger, J.L., 2011. Digital soil mapping and
631 modeling at continental scales: Finding solutions for global issues. *Soil Sci. Soc.*
632 *Am. J.* 75, 1201–1213.

633 Hengl, T., de Jesus, J.M., Heuvelink, G.B.M., Gonzalez, M.R., Kilibarda, M., Blagotić,
634 A., Shangguan, W., Wright, M.N., Geng, X., Bauer-Marschallinger, B., Guevara,
635 M.A., Vargas, R., MacMillan, R.A., Batjes, N.H., Leenaars, J.G.B., Ribeiro, E.,
636 Wheeler, I., Mantel, S., Kempen, B., 2017. SoilGrids250m: Global gridded soil
637 information based on machine learning. PLoS One 122, e0169748.

638 Hengl, T., Heuvelink, G.B.M., Kempen, B., Leenaars, J.G., Walsh, M.G., Shepherd,
639 K.D., Sila, A., MacMillan, R.A., de Jesus, J.M., Tamene, L., Tondoh, J.E., 2015.
640 Mapping soil properties of Africa at 250 m resolution: Random forests
641 significantly improve current predictions. PloS One 10, e0125814.

642 Heuvelink, G.B.M., 2014. Uncertainty quantification of GlobalSoilMap products. In
643 GlobalSoilMap: basis of the global spatial soil information system. 1st ed. CRC
644 Press Taylor & Francis Group, pp. 335–340.

645 Heuvelink, G.B.M., 2018. Uncertainty and uncertainty propagation in soil mapping
646 and modelling. Pedometrics, pp.439–461.

647 Heuvelink, G.B.M., Bierkens, M.F.P., 1992. Combining soil maps with interpolations
648 from point observations to predict quantitative soil properties. Geoderma 55, 1–
649 15.

650 Hoeting, J.A., Madigan, D., Raftery, A.E., Volinsky, C.T., 1999. Bayesian model
651 averaging: a tutorial. Stat. Sci. 14, 382–401.

652 Jolivet, C., Arrouays, D., Boulonne, L., Ratié, C., Saby, N.P.A., 2006. Le réseau de
653 mesures de la qualité des sols de France (RMQS). Etat d'avancement et
654 premiers résultats. Etude et Gestion des Sols 13, 149–164.

655 Koch, A., McBratney, A.B., Adams, M., Field, D., Hill, R., Crawford, J., Minasny, B.,
656 Lal, R., Abbott, L., O'Donnell, A., Angers, D., 2013. Soil security: solving the
657 global soil crisis. Glob. Policy 4, 434–441.

658 Kuhn, M., Weston, S., Keefer, C., Coulter, N., 2012. Cubist models for regression. R
659 package Vignette R package version 0.0, 18.

660 Laborczi, A., Szatmári, G., Kaposi, A.D., Pásztor, L., 2018. Comparison of soil texture
661 maps synthesized from standard depth layers with directly compiled products.
662 Geoderma <https://doi.org/10.1016/j.geoderma.2018.01.020>

663 Lark, R.M., 2000. A comparison of some robust estimators of the variogram for use in
664 soil survey. *Eur. J. Soil Sci.* 51, 137–157.

665 Loiseau, T., Chen, S., Mulder, V.L., Román Dobarco, M., Richer-de-Forges, A.C.,
666 Lehmann, S., Bourennane, H., Saby, N.P.A., Martin, M.P., Vaudour, E., Gomez,
667 C., Lagacherie, P., Arrouays, D., 2019. Satellite data integration for soil clay
668 content modelling at a national scale. *Int. J. Appl. Earth Obs.* 82, 101905.

669 Malone, B.P., McBratney, A.B., Minasny, B., Laslett, G.M., 2009. Mapping continuous
670 depth functions of soil carbon storage and available water capacity. *Geoderma*
671 154, 138–152.

672 Malone, B.P., Minasny, B., Odgers, N.P., McBratney, A.B., 2014. Using model
673 averaging to combine soil property rasters from legacy soil maps and from point
674 data. *Geoderma* 232, 34–44.

675 Martin, M.P., Lo Seen, D., Boulonne, L., Jolivet, C., Nair, K.M., Bourgeon, G.,
676 Arrouays, D., 2009. Optimizing pedotransfer functions for estimating soil bulk
677 density using boosted regression trees. *Soil Sci. Soc. Am. J.* 73, 485–493.

678 McBratney, A.B., Santos, M.M., Minasny, B., 2003. On digital soil mapping.
679 *Geoderma* 117, 3–52.

680 Minasny, B., McBratney, A.B., Hartemink, A.E., 2010. Global pedodiversity,
681 taxonomic distance and the World Reference Base. *Geoderma* 155, 132–139.

682 Monlong, J., 2018. Hippocampus, Github repository,
683 [https://github.com/jmonlong/Hippocampus/blob/master/content/post/2018-06-09-](https://github.com/jmonlong/Hippocampus/blob/master/content/post/2018-06-09-ClusterEqualSize.Rmd)
684 [ClusterEqualSize.Rmd](https://github.com/jmonlong/Hippocampus/blob/master/content/post/2018-06-09-ClusterEqualSize.Rmd)

685 Mulder, V.L., Lacoste, M., Richer-de-Forges, A.C., Arrouays, D., 2016a.
686 GlobalSoilMap France: High-resolution spatial modelling the soils of France up
687 to two meter depth. *Sci. Total Environ.* 573, 1352–1369.

688 Mulder, V.L., Lacoste, M., Richer-de-Forges, A.C., Martin, M.P., Arrouays, D., 2016b.
689 National versus global modelling the 3D distribution of soil organic carbon in
690 mainland France. *Geoderma* 263, 16–34.

691 Padarian, J., Minasny, B., McBratney, A.B., 2017. Chile and the Chilean soil grid: a
692 contribution to GlobalSoilMap. *Geoderma Reg.* 9, 17–28.

693 Padarian, J., Minasny, B., McBratney, A.B., Dalgliesh, N., 2014. Predicting and
694 mapping the soil available water capacity of Australian wheatbelt. *Geoderma*
695 *Reg.* 2, 110–118.

696 Poggio, L., Gimona, A., 2014. National scale 3D modelling of soil organic carbon
697 stocks with uncertainty propagation—an example from Scotland. *Geoderma* 232,
698 284–299.

699 Poggio, L., Lassauce, A., Gimona, A., 2019. Modelling the extent of northern peat
700 soil and its uncertainty with Sentinel: Scotland as example of highly cloudy
701 region. *Geoderma* 346, 63–74.

702 Quinlan, J.R., 1992. Learning with continuous classes. In 5th Australian Joint
703 Conference on Artificial Intelligence, 92, 343–348.

704 Raftery, A.E., Gneiting, T., Balabdaoui, F., Polakowski, M., 2005. Using Bayesian
705 model averaging to calibrate forecast ensembles. *Mon. Weather Rev.* 133,
706 1155–1174.

707 Roberts, D.R., Bahn, V., Ciuti, S., Boyce, M.S., Elith, J., Guillerá-Arroita, G.,
708 Hauenstein, S., Lahoz-Monfort, J.J., Schröder, B., Thuiller, W., Warton, D.I.,
709 Wintle, B.A., Hartig, F., Dormann, C.F., 2017. Cross-validation strategies for data
710 with temporal, spatial, hierarchical, or phylogenetic structure. *Ecography*, 40 (8),
711 913–929.

712 Román Dobarco, M., Arrouays, D., Lagacherie, P., Ciampalini, R., Saby, N.P.A.,
713 2017. Prediction of topsoil texture for Region Centre (France) applying model
714 ensemble methods. *Geoderma* 298, 67–77.

715 Román Dobarco, M., Bourennane, H., Arrouays, D., Saby, N.P.A., Cousin, I., Martin,
716 M.P., 2019. Uncertainty assessment of GlobalSoilMap soil available water
717 capacity products: A French case study. *Geoderma* 344, 14–30.

718 Rumpel, C., Amiraslani, F., Koutika, L., Smith, P., Whitehead, D., Wollenberg, D.,
719 2018. Put more carbon in soils to meet Paris climate pledges. *Nature* 564, 32–
720 34.

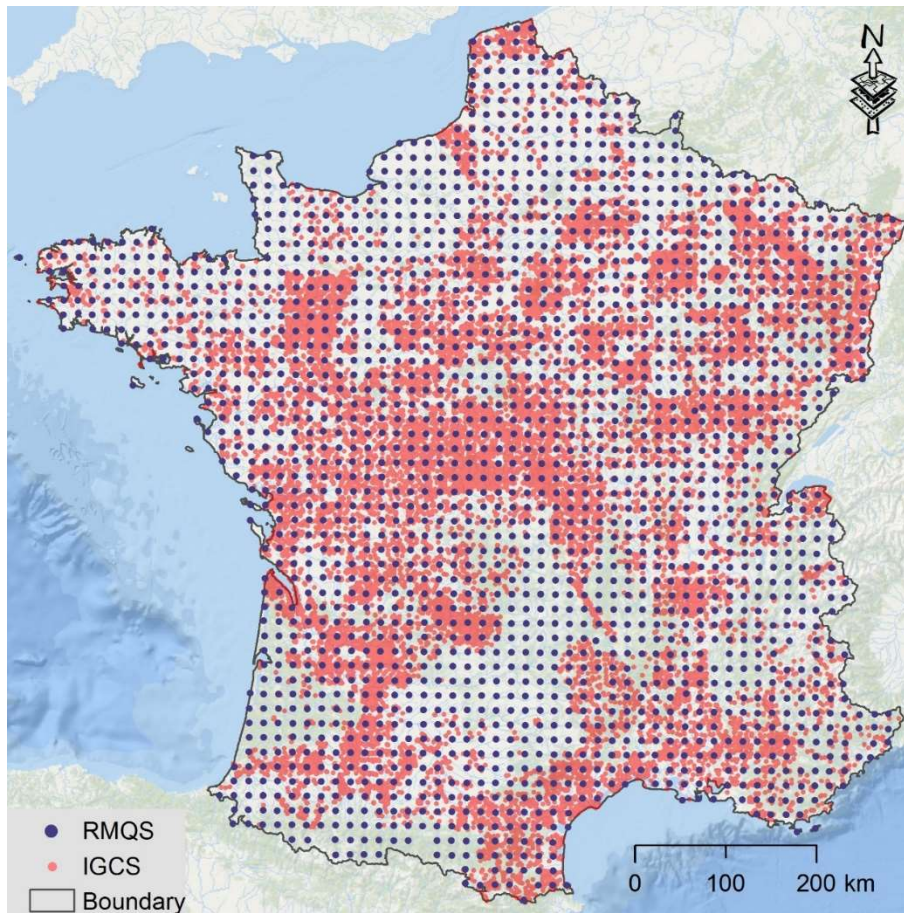
721 Sanchez, P.A., Ahamed, S., Carré, F., Hartemink, A.E., Hempel, J., Huising, J.,
722 Lagacherie, P., McBratney, A.B., McKenzie, N.J., de Lourdes Mendonça-Santos,
723 M., Minasny, B., 2009. Digital soil map of the world. *Science* 325, 680–681.

724 Sanderman, J., Hengl, T., Fiske, G., Solvik, K., Adame, M. F., Benson, L., Bukoski,
725 J.J., Carnell, P., Cifuentes-Jara, M., Donato, D., Duncan, C., Eid, E.M., zu
726 Ermgassen, P., Lewis, C.J.E., Macreadie, P.I., Glass, L., Gress, S., Jardine,
727 S.L., Jones, T.G., Nsombo, E.N., Rahman, M.M., Sanders, C.J., Spalding, M.,
728 Landis, E., 2018. A global map of mangrove forest soil carbon at 30 m spatial
729 resolution. *Environ. Res. Lett.* 13, 055002.

- 730 Tao, Y., Yang, T., Faridzad, M., Jiang, L., He, X., Zhang, X., 2018. Non-stationary
731 bias correction of monthly CMIP5 temperature projections over China using a
732 residual - based bagging tree model. *Int. J. Climatol.* 38, 467–482.
- 733 Tóth, G., Jones, A. Montanarella, L., 2013. The LUCAS topsoil database and derived
734 information on the regional variability of cropland topsoil properties in the
735 European Union. *Environ. Monit. Assess.* 185, 7409–7425.
- 736 Viscarra Rossel, R.A., Webster, R., Bui, E.N., Baldock, J.A., 2014. Baseline map of
737 organic carbon in Australian soil to support national carbon accounting and
738 monitoring under climate change. *Glo. Change Biol.* 20, 2953–2970.
- 739 Wadoux, A.M.C., Brus, D.J., Heuvelink, G.B.M., 2018. Accounting for non-stationary
740 variance in geostatistical mapping of soil properties. *Geoderma* 324, 138–147.

741 **Figures**

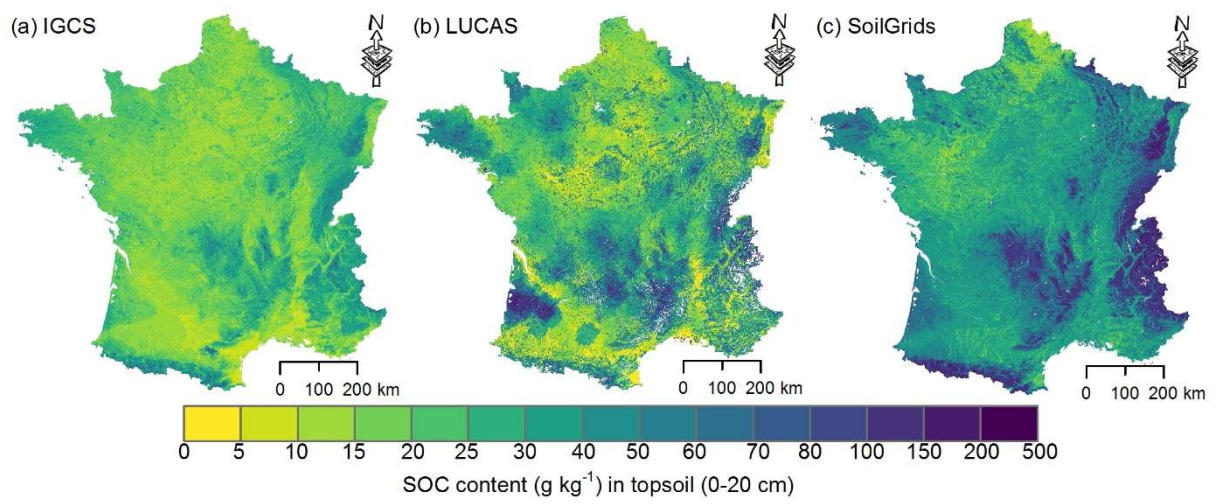
742 Fig. 1 Study area (Mainland France) and soil sampling sites from IGCS and RMQS
743 datasets.



744

745

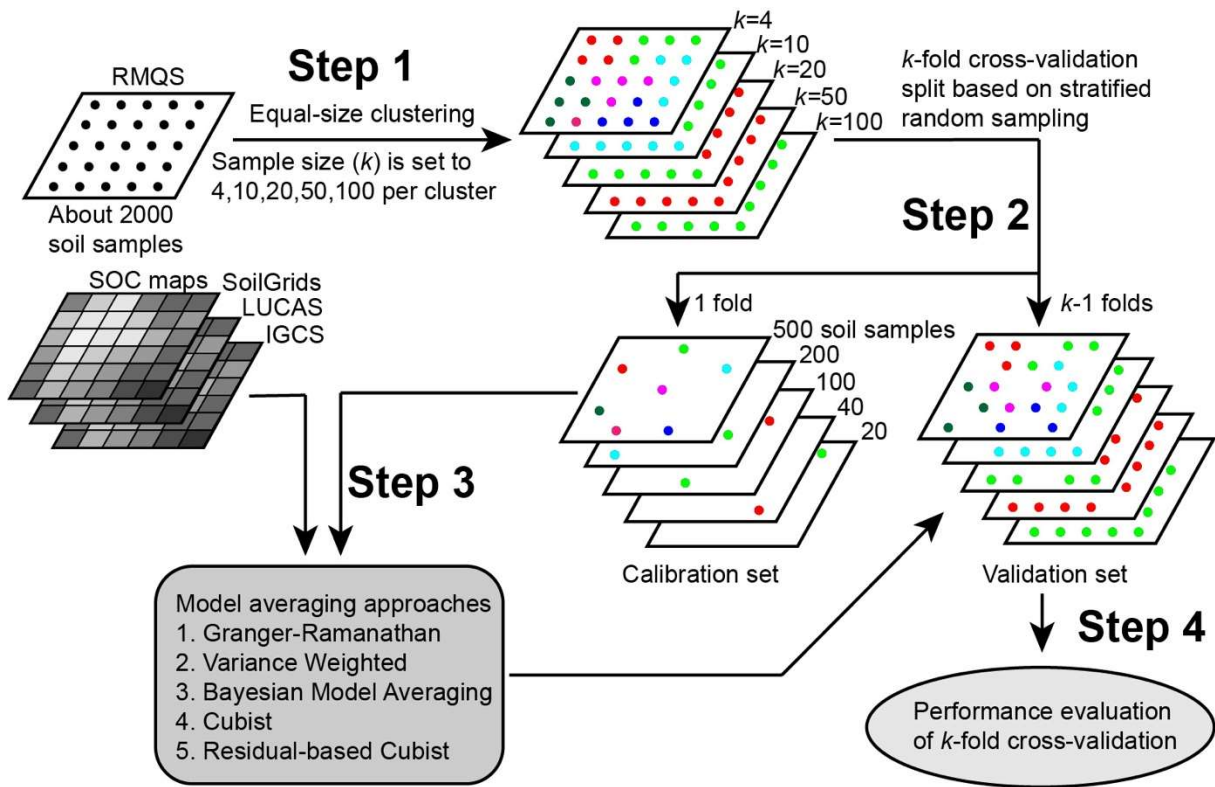
746 Fig. 2 SOC maps of mainland France from IGCS (a), LUCAS (b) and SoilGrids (c).



747

748

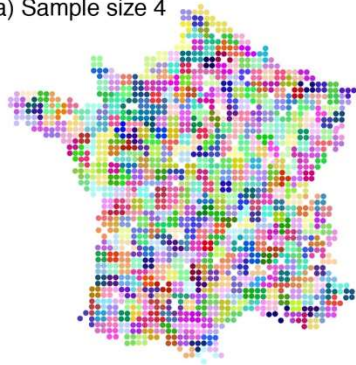
749 Fig. 3 Model averaging workflow.



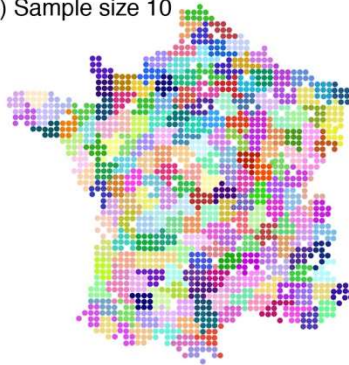
750
751

752 Fig. 4 Spatial cluster distribution of RMQS sites, using equal-size clustering. The
753 cluster sample sizes are 4 (a), 10 (b), 20 (c), 50 (d) and 100 (e).

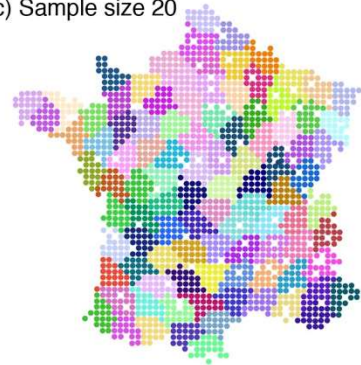
(a) Sample size 4



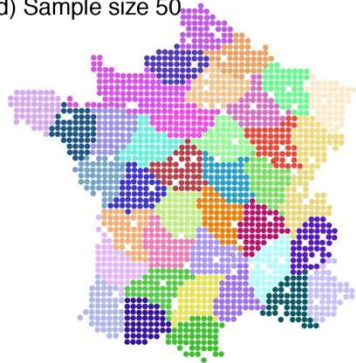
(b) Sample size 10



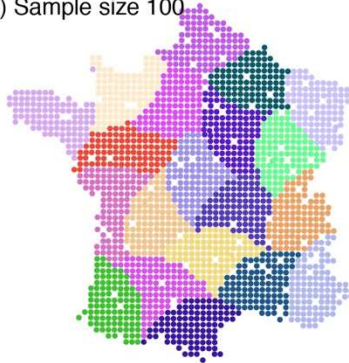
(c) Sample size 20



(d) Sample size 50



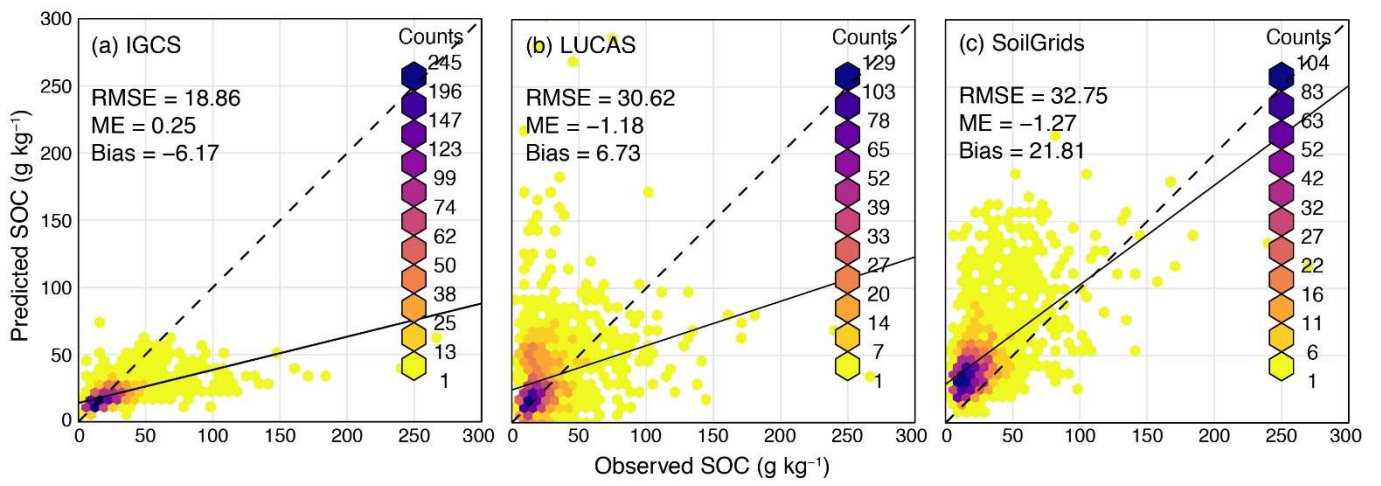
(e) Sample size 100



754

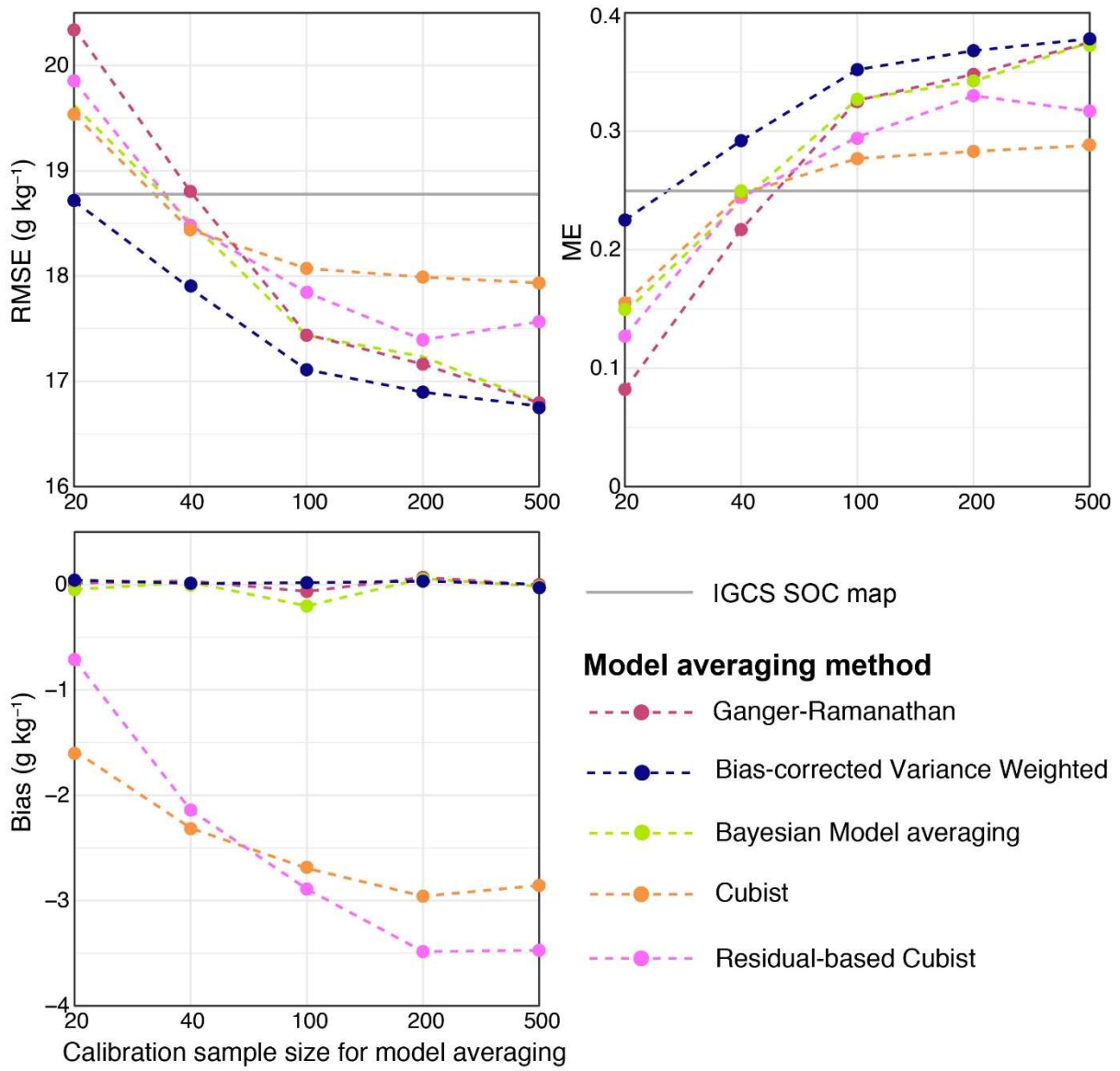
755

756 Fig. 5 Performance of IGCS (a), LUCAS (b) and SoilGrids (c) SOC maps.



757

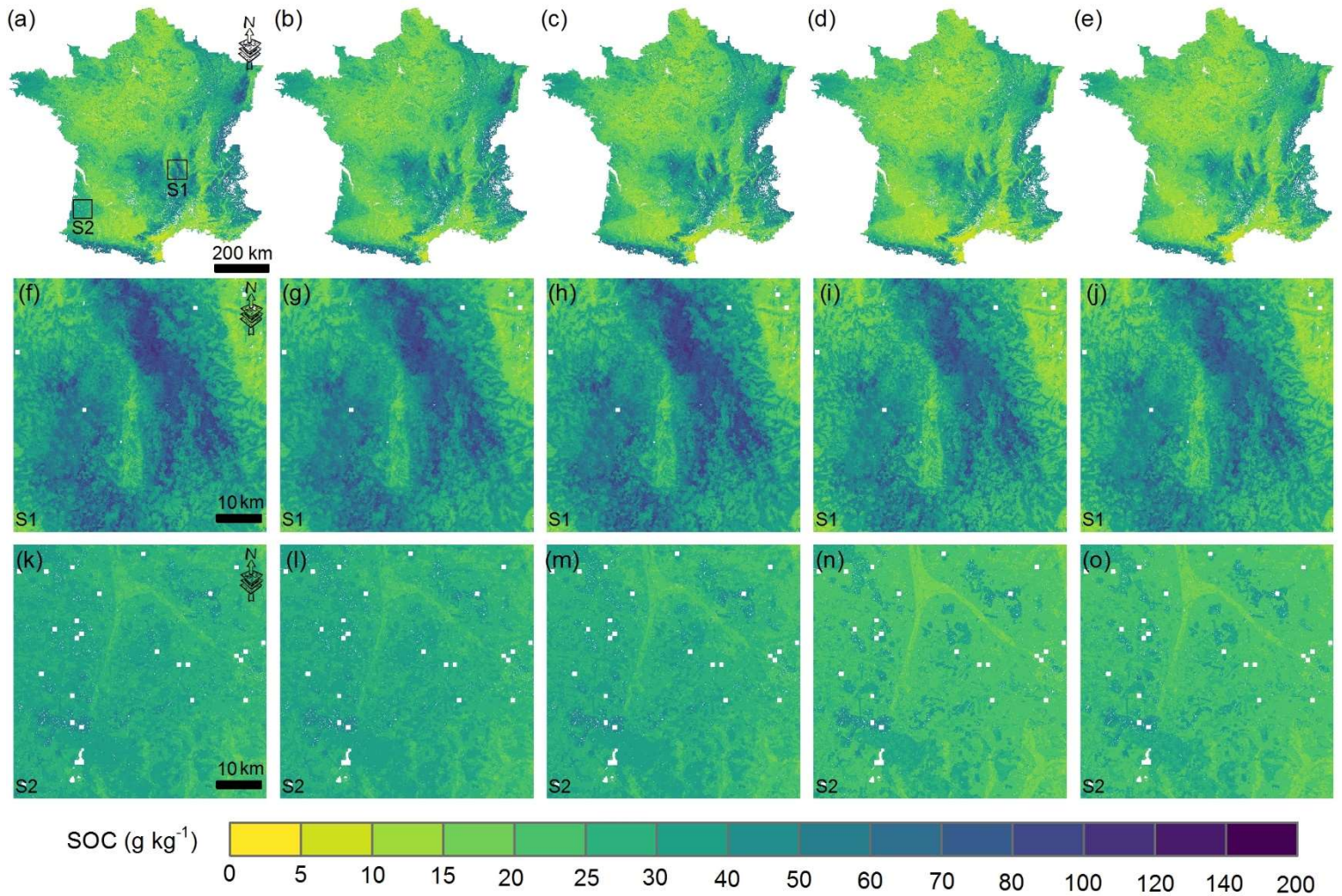
758 Fig. 6 Model performance of the five model averaging approaches using different
 759 calibration sample sizes.



760

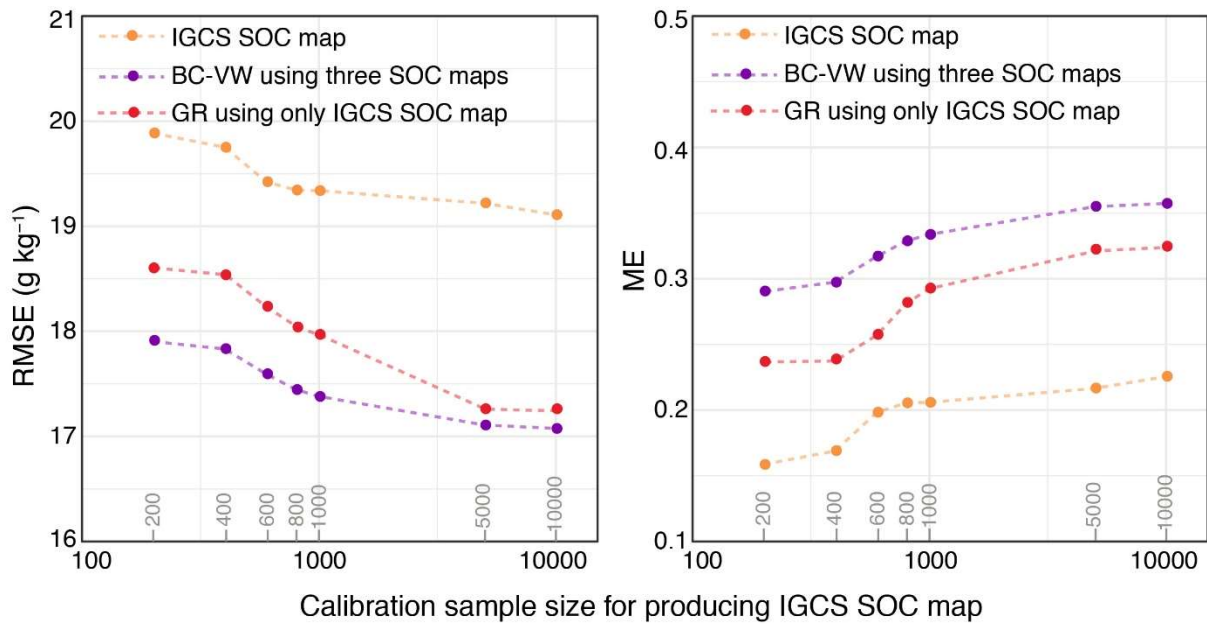
761

762 Fig. 7 SOC maps obtained from the Granger-Ramanathan (a), Bias-corrected
 763 Variance Weighted (b), Bayesian Model averaging (c), Cubist (d) and Residual-
 764 based Cubist (e) model averaging approaches, using all RMQS data for calibration.
 765 Local comparisons in areas S1 (f, g, h, l and j) and S2 (k, l, m , n and o) are also
 766 shown for all five model averaging approaches.



767

768 Fig. 8 Model performance of the Bias-corrected Variance Weighted (BC-VW) model
 769 averaging (using 200 calibration samples for three SOC maps) and Granger-
 770 Ramanathan (GR) model (using the same 200 calibration samples for only calibrating
 771 IGCS SOC map) when using different calibration sample sizes (200 to 10,000) for
 772 generating IGCS SOC map. Using only the LUCAS and SoilGrids SOC maps for BC-
 773 VW model averaging leads to an RMSE of 23.65 g kg⁻¹ and ME of -0.24 (points not
 774 shown). The x-axis is on log10 scale.



775

776

777 **Tables**

778 Table 1 Summary statistics of SOC content (g kg⁻¹) in topsoil (0-20 cm) for IGCS,

779 RMQS and LUCAS datasets.

Dataset	Land use*	N	Min.	Q1	Median	Mean	Q3	Max.	Sk.	SD
IGCS	F & G	5,785	0.39	12.75	19.86	24.88	30.83	373.00	3.42	20.97
	A	24,596	0.09	9.70	13.68	16.66	19.75	354.05	4.92	12.88
RMQS	F & G	985	3.78	18.86	28.37	35.51	44.00	266.60	2.81	26.01
	A	1,011	2.58	11.10	15.40	18.19	22.30	133.00	3.01	11.16
LUCAS	A & G	2,950	1.00	13.20	19.99	26.20	31.30	472.10	6.11	23.93

780 N, dataset size ; Min., minimum; Q1, first quantile; Q3, third quantile; Max., maximum; Sk., skewness;

781 SD, standard deviation. * F, forest; G, permanent grasslands; A, arable.