



**HAL**  
open science

## Do climate and land use affect the pool of total silicon concentration? A digital soil mapping approach of French topsoils

Amélia Landré, Sophie S. Cornu, J.-D Meunier, Annie Guérin, Dominique Arrouays, Manon Caubet, Céline Ratié, Nicolas P A Saby

### ► To cite this version:

Amélia Landré, Sophie S. Cornu, J.-D Meunier, Annie Guérin, Dominique Arrouays, et al.. Do climate and land use affect the pool of total silicon concentration? A digital soil mapping approach of French topsoils. *Geoderma*, 2020, 364 (Avril), pp.114175. 10.1016/j.geoderma.2020.114175 . hal-02503457

**HAL Id: hal-02503457**

**<https://hal.science/hal-02503457>**

Submitted on 10 Mar 2020

**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.

1 DO CLIMATE AND LAND  
2 USE AFFECT THE POOL OF  
3 TOTAL SILICON  
4 CONCENTRATION? A  
5 DIGITAL SOIL MAPPING  
6 APPROACH OF FRENCH  
7 TOPSOILS.

8 Landré, A.<sup>a,b</sup>, Cornu, S.<sup>c</sup>, Meunier, J.-D.<sup>c</sup>, Guerin A.<sup>d</sup>, Arrouays D.<sup>a</sup>, Caubet M.<sup>a</sup>, Ratié C.<sup>a</sup>,  
9 Saby, N.P.A.<sup>a</sup>

10 a) INRA, Infosol, US 1106, Orléans, France

11 b) INRA - INPT-ENSAT - INPT-EI-Purpan, UMR 1248 AGIR AGroécologie,  
12 Innovations, teRritoires. Centre de recherche Occitanie-Toulouse, Auzeville, France.

13 c) Aix-Marseille Univ, CNRS, IRD, Coll de France, INRA, CEREGE, Aix-en-Provence,  
14 France

15 d) INRA, Laboratoire d'Analyses des Sols US, Arras, France

16

## 17 Abstract

18 Silicon (Si) is the second most abundant element in the Earth's crust after O. Its  
19 concentration in soils is highly variable from less than 1 % to greater than 45 %. Parent  
20 material is well known to be a major parameter for explaining this variability. In this  
21 study, we proposed to analyze the impact of climate and land use on the total Si  
22 concentration in soils and to explore the link between total Si and plant available Si  
23 (PAS). To do so, we based our analysis on the French soil monitoring network  
24 considering the upper soil horizon that was thought to be the most impacted by both  
25 the effect of land use and climate and was also the most important horizon in terms of  
26 plant availability. In order to extract the impact of climate and land use and for digital  
27 mapping purposes, we stratified the database by parent material and soil-types. This  
28 stratification was based on the classification used in the 1:100 000 French soil map  
29 and 1:100 000 French soil parent material map. For non carbonated soils, we showed  
30 that Si concentrations was decreasing with annual rainfall, evidencing a climatic effect  
31 on the total Si concentration of French topsoils. No significant effect of the land used  
32 could be identified. At last, we showed that PAS (by the  $\text{CaCl}_2$  method) is negatively  
33 weakly correlated to total Si concentration. This relationship is however variable among  
34 soil classes.

35

## 36 Keywords

37 Silicon, silica, RMQS, Digital Soil Mapping, France.

## 38 1 Introduction

39 Silicon is the second most abundant element in the Earth's crust after O, and it has an  
40 average mass concentration of 28 % (Wedepohl, 1995). The SiO<sub>4</sub> tetrahedron is the  
41 elemental brick that constitutes the basic structure of Si in nature from solid (silicates)  
42 to soluble (silicic acid) states. In soils, the variability of Si is large, from less than 1 wt%  
43 in Histosols to greater than 45 wt% in some Podzols (Sommer et al., 2006). Si is  
44 included in a large number of minerals issued either from the parent material as for  
45 instance quartz, feldspar, and phyllosilicates (muscovite, biotite for instance), or from  
46 the transformation of the parent material minerals, along weathering and pedogenesis,  
47 in clay minerals and amorphous silica particles, notably phytoliths for the part of Si  
48 recycled by vegetation (Sommer et al., 2006).

49 The soil concentration of total Si is intimately link to the concentration of Si from parent  
50 rocks (Berner and Berner, 1996; [Gray et al., 2016](#)). [Gray et al. \(2016\)](#) showed that sand  
51 dunes are shown to be the richest parent material (close to 100 %) while limestones  
52 are the poorest (a few percent). Among the parent material rich in Si (> 60 %), quartz  
53 and arkose sandstones, granite, rhyotite, greywacke, granolite, dacite, shale are  
54 encountered, while diorite, andesite, basalt and peridotite have Si content ranging from  
55 40 to 60 % and laterite and bauxite from 10 to 15 % (Gray et al., 2016 and references  
56 included). In addition, weathering and pedogenesis modify the initial Si pool both in  
57 concentration and status (Lucas et al., 1993; White et al., 2012). The upper soil horizon  
58 is the most impacted by these processes being both depleted by weathering and clay  
59 translocation, that may have contradictory effect on the total Si concentration of the  
60 horizon. Weathering and pedogenesis are influenced by both climate and vegetation  
61 as demonstrated by river dissolved Si (DSi) (Bluth and Kump, 1994; Dürr et al., 2009)  
62 at the territory scale. Using DSi in the streams as an indicator of total chemical

63 weathering is however not straightforward because of the fraction of Si that is retained  
64 in the soil for forming clay minerals or the one that is retained in biogenic silica in soils  
65 and streams (Frings et al., 2015). Analysis of climosequences have proven to be  
66 helpful for documenting some modification of the soil Si pools. On a volcanic  
67 climosequence, secondary Si phases (produced by weathering) increase under a more  
68 humid climate (Taboada et al., 2019). On intrusive rock climosequence, smectite  
69 distribution is correlated to Si loss when rainfall (or precipitation) increases (Egli et al.,  
70 2003), but the analysis of the variation of soil surface Si is not well addressed. As  
71 shown by Dere et al. (2016), the presence of poorly reactive quartz in the parent rock  
72 may be greater at the surface resulting from a relative enrichment during weathering.  
73 Therefore, the impact of rainfall on the soil Si deserves to be more fully documented.

74 While vegetation on the weathering scale is strongly correlated to climate, human  
75 activity has strongly modified the land use and thus the biological cycle of most of the  
76 elements. Some of the main cultivated crop are Si accumulators and long-term  
77 exportation may act on Si budget. It seems well documented that conversion of forest  
78 to cropland lead to a decrease of the phytolith pools (Struyf et al., 2010; Guntzer et al.,  
79 2012; Vandevenne et al. (2015). Struyf et al. (2010) suggest that the transformation of  
80 the european temperate forests into cultivated land has also lead to a decrease of  
81 weathering of silicate minerals. However, Yang and Zhang (2018) came to the opposite  
82 conclusion (agricultural activities increase silicate weathering) based on a geochemical  
83 analysis of rain and streams in subtropical China. In the subtropical southern Brazil,  
84 Ameijeiras-Marino et al. (2018) share Yang and Zhang (2018)'s conclusion using the  
85 Ge/Si tracer. Accordingly, the question of the intensity of the impact of land use on soil  
86 Si loss remains open.

87 The study of the status of Si in soils and its biogeochemical cycle is a subject of  
88 increasing interest because of the growing body of evidence showing that Si can  
89 improve crop development (Coskun et al., 2019; Guntzer et al., 2012; Liang et al.,  
90 2015; Rodrigues and Datnoff, 2015). Plant available Si (PAS) depends on the reserve  
91 of weatherable Si-bearing minerals (Cornelis and Delvaux, 2016) and can be estimated  
92 using various protocols (Narayanaswamy and Prakash, 2010; Meunier et al.; 2017). The  
93 relationship between soil total Si and PAS is complex depending on soil type, parent  
94 material and degree of weathering. On various agricultural soils of Asia, a negative  
95 significant correlation between PAS and soil Si content was found (Yanai et al., 2016;  
96 Meunier et al. (2017). This correlation was due to the presence of low solubility Si  
97 minerals such as quartz for explaining low PAS values. On volcanic soils containing  
98 no quartz, Henriot et al. (2008) founded a highly significant positive correlation (n= 6)  
99 between soil Si content and Si extracted by the 0.01M CaCl<sub>2</sub> method due to a decrease  
100 of the reserve of weathered Si minerals along weathering.

101 The objective of this study was to (i) determine the impact of climate and land-use on  
102 the total topsoil Si concentrations (Si<sub>tot</sub>) in French soils and (ii) explore the relationship  
103 between Si<sub>tot</sub> and PAS. Topsoil horizon was chosen as it was considered as the most  
104 weathered and pedogenized horizon where the roots are mainly concentrated. The  
105 French territory was chosen as (i) France is one of the countries with the largest soil  
106 diversity in the world (Minasny et al., 2010), that offers notably a good  
107 representativeness of European soils notably; (ii) France processes a soil monitoring  
108 network (RMQS) for which Si<sub>tot</sub> is available along with all the typical soil characteristics  
109 (Landré et al., 2018). Indeed, soil Si is rarely analyzed at the territory scale, with the  
110 notable exception of the soil Si map at the European scale with only one site every  
111 2500 km<sup>2</sup> (De Vos et al., 2006; Reimann et al., 2014) and the French soil monitoring

112 network (RMQS) with one point at each node of 16 to 16 km grid (Jolivet et al., 2006;  
113 Landré et al., 2018).

114 Since geology and soil type (pedogenesis) are the main Si drivers, we first stratified  
115 the database in geological-pedological classes. Then we considered the impact of  
116 climate and land-use in both the obtained classes and at the national scale with a  
117 modern digital soil mapping (DSM) approach using the geological-pedological classes  
118 defined as a covariate. The strong input of this type of covariate in the DSM model was  
119 clearly demonstrated by Gray et al., (2016).

## 120 2 Materials and Methods

### 121 2.1 The data

122 Most of the data used in this study were obtained from 2111 sites from the first  
123 campaign of the RMQS sampled between 2001 and 2009 (Jolivet et al., 2006), which  
124 covered all the mainland of France (around 550 000 km<sup>2</sup>) based on a 16 km × 16 km  
125 grid. Among these sites, two had no soil type description, 98 had no parent material  
126 description, six were anthroposols and three histosols. All these sites were not further  
127 considered for our analysis, ending with database of 2004 sites.

128 Composite samples of the top horizon (0-30 cm) were sampled, air dried, and sieved  
129 to 2 mm before laboratory analysis at the Soil Analysis Laboratory of INRA in Arras,  
130 France. The following parameters were measured: (i) the total organic C (OC) content  
131 measured by dry combustion (NF ISO 10694), (ii) particle size distribution by wet  
132 sieving and the pipette method (NF X 31-107), (iii) cation exchange capacity and  
133 exchangeable cations (hexamminecobalt method NF ISO 23470), (iv) pH in water (1  
134 to 5 soil to water ratio; NF ISO 10390), (v) calcium carbonate using the volumetric

135 method (NF ISO 10693) ( $\text{CaCO}_3$ ), (vi) and total P, K, Ca, Mg, Fe, and Al determined  
136 by ICP-AES after dissolution with hydrofluoric and perchloric acids (NF X 31-147).

137 Landré et al. (2018) analyzed a subset of 673 samples for  $\text{Si}_{\text{tot}}$ . For the sites not  
138 analyzed, we estimated  $\text{Si}_{\text{tot}}$  according to the following conceptual equation:

$$139 \quad \text{Si} = f(\text{Al}, \text{Fe}, \text{K}, \text{Na}, \text{Ca}_{nc}, \text{Mg}_{nc}, \text{P}, \text{Mn}, \text{OC}, \text{CaCO}_3, \text{residual water}) \quad (1)$$

140 where  $\text{Ca}_{nc}$  and  $\text{Mg}_{nc}$  are the fractions of Ca and Mg that are not included in carbonate  
141 minerals nor adsorbed on the exchangeable surfaces, and OC is the organic matter  
142 percentage. The total Si concentration in French topsoil was modelled using a Cubist  
143 modelling algorithm (see Supplementary File 1 for more details).

144 The database was completed by PAS measurement ( $\text{Si}_{\text{CaCl}_2}$ ) using the 0.01 M  $\text{CaCl}_2$   
145 method (Haysom and Chapman 1975) on 1986 sites. This widely used method  
146 (Meunier et al., 2017) allows estimating the pool of Si that is readily soluble. Briefly, 2g  
147 of dry soil was mixed with 20 ml of the solution and shaken during 16h in polyethylene  
148 tubes. The solutions were then filtered at 0.45  $\mu\text{m}$ , and Si concentration was measured  
149 using Inductively Coupled Plasma Atomic Emission Spectroscopy (axial ICP-AES; 720  
150 ES, Varian).

## 151 2.2 Stratification of the database in homogeneous geo- 152 pedological classes

153 Since  $\text{Si}_{\text{tot}}$  in soil is known to be linked to parent material and being anti-correlated to  
154 carbonate concentration in soils (Landré et al., 2018), we first stratified the database  
155 by the type of parent material and by the content in carbonates. In French soils,  
156 carbonates were encountered in soils developed on sedimentary rocks with a limit for  
157 carbonated soils classically considered at a  $\text{CaCO}_3$  concentration of 5 % (Baize, 2000).



158 Since a large part of the sedimentary rocks were defined on the basis of  
159 geomorphological processes (e.g. alluvions, colluvions, terraces...), we were not able  
160 for DSM purposes to relate directly parent material to carbonate content. In addition,  
161 since we were working on the topsoil horizon, depending on the stage of pedogenesis,  
162 soils developed on carbonated rocks may have been completely decarbonized (e.g.  
163 Calcisols notably). Other geological domains were sorted in igneous extrusive or  
164 intrusive and metamorphic. The classification used did not allow further classifying the  
165 parent materials in more geochemical meaning classes as made by Gray et al. (2016),  
166 since (i) a large part of the sites had no more information than igneous extrusive or  
167 intrusive or metamorphic rocks; (ii) acidic extrusive igneous rocks or basic intrusive  
168 igneous rocks were too poorly represented in the database (no more than two to three  
169 individuals). Podzols were also separated, these soils being mainly composed of  
170 quartz were supposed to have high  $Si_{tot}$ . We then sorted the soils in the different  
171 geological groups a step further with the exception of the soils on igneous intrusive  
172 rocks that were poorly represented. Soils on metamorphic parent materials and on  
173 igneous extrusive rocks were sorted by soil types, ranging from poorly differentiated  
174 soils (regosols, lithosols and rankers) to strongly differentiated soils (planosols, luvisols  
175 and podzoluvisols) for soils on metamorphic rocks. Podzols were sorted between those  
176 developed on sedimentary rock and those developed on other parent materials. For  
177 non-carbonated soils on sediment, sorting was performed both by soil types and a  
178 more detailed description of the considered sediment.

## 179 2.3 Data treatment

### 180 2.2.1. Mapping total topsoil Si concentrations in French soils

181 The DSM approach was based on the work of McBratney et al. (2003), who proposed  
182 a quantitative relationship between soil properties and the soil forming factors plus a  
183 spatially correlated residual element ( $e$ ), as follows:

$$184 \quad \text{Soil} = f(s, c, o, r, p, a, n) + e \quad (2)$$

185 where *Soil* is a soil property. The *s* refers to other or previously measured properties  
186 of the soil at a point either from prior soil maps or from remote or proximal sensing  
187 data; *c* refers to the climatic properties of the environment at a point; *o* refers to  
188 organisms, including vegetation or fauna, or human activity; *r* refers to relief; *p* refers  
189 to the parent material or lithology; *a* refers to the soil age; and *n* refers to the space or  
190 spatial position. Finally, *e* is the locally varying, spatially dependent residuals from *f*.

191 We selected a set of available environmental covariates describing the *scorpan* factors  
192 for the whole French territory. We harmonized them at 90 m resolution (Table 1). The  
193 parental material and soil covariates were generated by the results of the PCA step.  
194 We also used a very popular spectral index, namely the normalized difference  
195 vegetation index (NDVI) (Huete et al., 2002). We focused on the yearly changes in the  
196 NDVI computed from a time series of remote sensing data to describe the  
197 photosynthetic capacity of the vegetation cover. This variable was interpreted as a  
198 vegetation growth dynamics proxy. The underlying assumption was that this  
199 information on changes in vegetation may reflect various behaviors linked to climate,  
200 land management or soil properties (see Loiseau et al., 2019).

201 The spatial predictive model between  $Si_{\text{tot}}$  and *scorpan* covariates was first constructed  
202 using the ensemble learning method Random Forests (RF) (Breiman, 2001). We used  
203 the RF implementation provided by the package `randomForest` in R (Liaw and Wiener,

204 2002; R. Core Team, 2016). Three parameters should be defined in the RF model,  
205 namely the number of trees to grow ( $n_{tree}$ ), the number of variables randomly sampled  
206 as candidates at each split ( $m_{try}$ ), and the minimum size of the terminal nodes  
207 (nodesize; Liaw et al., 2002). The default values were used for  $n_{tree}$  and nodesize,  
208 which were 500 and 5, respectively. The optimal value of  $m_{try}$  was set at 2 by the lowest  
209 out-of-bag error estimate.

210 To extract useful information from this large set (22) of potentially correlated  
211 covariates, we ran a preliminary step of variable selection using the Boruta algorithm  
212 (Kursa and Rudnicki, 2010). This algorithm was a wrapper built around the RF  
213 classification algorithm implemented in the R package randomForest, and it used a Z  
214 score computed by dividing the average loss by its standard deviation. The algorithm  
215 was used in feature selection (Kursa and Rudnicki, 2010), but was also applied to  
216 support the model establishment of the RF regression (Kursa, 2014).

217 The residuals of the model in equation 2) were computed as the difference between  
218 the RF predictions and the measured values at the observed location and then  
219 interpolated by ordinary kriging (Matheron, 1971). The R package gstat (Pebesma and  
220 Graeler, 2019) was used to select variogram models and perform the kriging  
221 procedure. The final predictions summed the RF predictions and the Kriging outputs  
222 (Keskin and Grunwald, 2018). The model performance was evaluated by 30-folds  
223 cross-validation.

224 To rank the influences of the final list of scorpan factors on  $Si_{tot}$ , we calculated covariate  
225 importance from the RF algorithm as the mean increase in accuracy (%IncMSE). This  
226 indicator was constructed by permuting the values of each variable of the validation  
227 set, recording the prediction error, and comparing the set with the un-permuted

228 validation set prediction of the variable (normalized by the standard error). It calculated  
229 the average increase in the squared residuals of the validation set when the variable  
230 was permuted. A higher %IncMSE value represented higher variable importance.

### 231 2.2.2. Other statistical analysis

232 Since Gray et al. (2016) showed that pedological parameters and  $Si_{tot}$  were highly  
233 dependent from parent material composition, the database was stratified in  
234 homogeneous subgroup defined on a combination of geological and pedological  
235 classes, in order to use the obtained classes for DSM approach. The parent material  
236 and soil type classifications used for the stratification were those used at the French  
237 territory scale that is the EUSIS classification (King et al, 1994) for the parent material  
238 and the FAO (1985) classification for the soil types.

239 We then performed multiple comparison tests: multiple Mann-Whitney tests, Kruskal-  
240 Wallis tests, and a post-hoc test from Siegel and Castellan (1988). These tests were  
241 performed using the medians of  $Si_{tot}$  within the geo-pedological groups identified in  
242 order to discuss the discriminating power of the obtained classification in terms of total  
243 Si topsoil concentrations.

244 In order to understand the meaning of the relations between total  $Si_{tot}$  and  $Si_{CaCl_2}$   
245 obtained for the different parent material domains, we ran a principal component  
246 analysis (PCA) on all the soil characteristics available in the RMQS database, with total  
247 Si was considered a passive variable.

## 248 3 Results and Discussion

### 249 3.1 Stratification of the database in homogeneous geo- 250 pedological classes

251 For soils developed on sedimentary rocks, we observed as expected a strong negative  
252 correlation between  $Si_{tot}$  and carbonate contents for soils with carbonates  
253 concentration higher than 5 % (Figure 1a). This negative correlation was interpreted  
254 as a diluting effect of silicate minerals by carbonates. We then sorted geo-pedological  
255 classes for soils developed on sediment in carbonated soils developed on sediments  
256 and non-carbonated soils developed on sediment (Figure 1b) in order to use this  
257 classification in the DSM approach. Based on this separation and on other parent  
258 material classes (igneous intrusive and extrusive rocks, metamorphic rock), a boxplot  
259 analysis confirmed that this first separation was meaningful (Figure 2) with the highest  
260 concentrations in  $Si_{tot}$  for Podzol topsoils and the lowest for topsoils on igneous  
261 extrusive rocks. This low concentration obtained for igneous extrusive rocks may seem  
262 surprising as Gray et al. (2016) showed that both Si rich (rhyotite) and Si poor (basalt,  
263 andesite) rocks could be found in that group. However, in the case of France, igneous  
264 extrusive rocks that are poorly abundant (only 28 sites over 2000) consist mainly in  
265 basalt, slag, tuff and other basic volcanic rocks while rhyolite and other acid volcanic  
266 rocks are rare (only two sites over the 28 classified as igneous extrusive volcanic  
267 rocks).

268 At last, some parent material classes while significantly different from others, still  
269 exhibit a very large Si concentration variability, justifying the further classification  
270 performed on parent material and soil type criteria (Fig. 3). In this figure, we observe  
271 that Podzol topsoils on sedimentary parent materials exhibited significantly higher  $Si_{tot}$   
272 ( $p$ -value =  $8.573e-06$ ; Fig. 3). These Podzols are located in two main areas in France,

273 namely the Landes of Gascony (in southwest France) and an area in central France  
274 north of the Massif Central mountains (Figure 4). The Landes of Gascony was entirely  
275 covered by aeolian sands (their texture is usually more than 95 % sand; Augusto et  
276 al., 2010) that were nearly pure quartz grains (300-550  $\mu\text{m}$  in size) and therefore could  
277 be considered as a reference for nearly pure Si topsoils. The other area in central  
278 France corresponded to an ancient delta of the Loire River, with very sandy material  
279 also quartz rich originating from the erosion of the ancient Massif Central and its  
280 northern foothills. Our results are in a good agreement with those of Gray et al. (2016)  
281 showing that dune sands were among the richest parent material in Si.

282 On igneous extrusive rock, andosol topsoils exhibited significantly lower  $\text{Si}_{\text{tot}}$  than  
283 those of other soil types ( $p$ -value < 0.001; Figure 3). For topsoils on metamorphic rock,  
284 poorly differentiated soils (Lithosols, Regosols, and Rankers) differed from well  
285 differentiated soils (Planosols, Luvisols, and Podzoluvisols), with Cambisols and  
286 Fluvisols being intermediate (Figure 3). These three groups of topsoils significantly  
287 differed in their  $\text{Si}_{\text{tot}}$  ( $p$ -value < 0.001 and post-hoc test with a  $p$ -value < 0.05). Soils on  
288 sedimentary rocks were the more abundant, and their situation was more complex.  
289 Nineteen and seventeen groups were identified either on the basis of both the parent  
290 material and the soil type for the carbonated and non-carbonated soils, respectively.  
291 These groups differed in their topsoil  $\text{Si}_{\text{tot}}$  ( $p$ -value < 0.001 for the carbonated soils and  
292  $p$ -value < 0.001 for the non-carbonated soils; Figure 3). As an example, very low  
293 concentrations of topsoil  $\text{Si}_{\text{tot}}$  were found for poorly differentiated soils on chalk (Fig 3)  
294 i.e. nearly pure  $\text{CaCO}_3$  with  $\text{Si}_{\text{tot}}$  close to zero. These soils are encountered mainly in  
295 two areas in France: in the Champagne area and in the Charentes (in western France)  
296 where soils are locally called Champagne's soils (Figure 4).

297 Interestingly, this analysis showed that some carbonated soils could exhibit very high  
298  $Si_{tot}$  in their topsoils, despite the relative absence of Si in carbonates (e.g. calcic  
299 cambisols on other sediments and solonchack; Figure 3). This classification was used  
300 in the DSM approach as a covariate along with climate, vegetation using land use, and  
301 relief variables.

### 302 3.2 Impact of environmental factors (climate and land 303 use) on the total Si concentrations of the French top 304 soils

305 The covariate importance from the RF algorithm are presented in Figure 5. Two groups  
306 of covariates stood out with a threshold of around 20 %. The first group contained six  
307 environmental covariables. The first covariate, called PG, corresponded as expected  
308 to the pedogeological classes previously defined. This covariable made a significant  
309 contribution in the model construction, as it exhibited a mean increase in accuracy of  
310 55 %. Next was the elevation (“srtm”) with an importance of 33 %, followed by the NDVI  
311 covariate (“PC1\_NDVI”) and the type of climate (“typo”) with an importance of 27 %  
312 and 24 %, respectively. Finally, the annual mean evapotranspiration (“etp\_mean”) and  
313 the net primary production (“NPP\_max”) presented an importance of 22 % and 21 %,  
314 respectively. NDVI covariate is not an independent variable. Indeed, as demonstrated  
315 by Loiseau et al. (2019), beside land use, NDVI is partially linked to pedology and  
316 climate. Elevation is also not an independent variable in the case of France. For  
317 instance, in the case of the Massif Central, igneous extrusive rocks only located at the  
318 highest altitude Elevation can be a proxy for both climate and geology in France. Thus,  
319 after parent material and soil type, climate seemed to be an important factor in  
320 determining the total topsoil Si concentration.

321 Figure 6 presents  $Si_{tot}$  as a function of the annual rainfall for the different parent  
322 material domains. For carbonated soils on sediments,  $Si_{tot}$  was independent from  
323 annual rainfall as well as non-carbonated soils developed on sedimentary parent  
324 materials. A significant negative correlation was obtained for soils on other parent  
325 materials, with higher  $Si_{tot}$  for lower annual rainfall. This correlation was highly  
326 significant with the exception of the soils developed on igneous extrusive rocks. The  
327 latter being represented by fewer individuals ( $n=28$ ) than the other groups, a  
328 correlation significant only at a 5 % confidence level was acceptable. Such a negative  
329 correlation between the concentration in  $Si_{tot}$  and the average annual rainfall was  
330 interpreted as resulting from the increasing losses of Si with an increase in weathering  
331 intensity owing to the larger water flow through the soil (higher rainfall). Therefore, the  
332 analysis of surface soils may be used as an alternative to the river chemistry (Bluth  
333 and Kump (1994) to document the effect of climate on Si weathering.

334 The impact of land use on  $Si_{tot}$  in French topsoils was tested on the geo-pedological  
335 classes defined on Figure 3. To do so, for each class sites were sorted in two land use  
336 groups, namely permanent vegetation (forests and pasture) and arable land. We then  
337 performed a two-way ANOVA with an interaction term between land use and geology  
338 followed by a tukey HSD test. The results showed that the interaction effect was slightly  
339 significant ( $p = 0.012$ ). We investigated also the model coefficients and we found that  
340 for most of them, the 95 % confidence intervals overlapped except for a few cases.  
341 (see  $p$ -values reported in supplementary material).

342 Lastly, neither the relief nor the landuse could be clearly identify as a factor impacting  
343 total Si concentration in soils at least at this scale. Concerning relief, all the proxies of  
344 the landforms had a very low weight in this analysis (Figure 4d). For landuse, our data



345 could not document the positive or negative effect of agriculture, at this scale, on the  
346 intensity of silicate weathering that has been suggested in the literature (Struyf et al.,  
347 2010; Yang and Zhang, 2018).

### 348 3.3 Relationship between $Si_{tot}$ and $Si_{CaCl_2}$ in French 349 topsoils.

350 When considering the total dataset (Table 2), a significantly different from 0, negative  
351 correlation between  $Si_{CaCl_2}$  and  $Si_{tot}$  was given but the coefficient is weak ( $r = -0.32$ ).  
352 While considering the different classes,  $Si_{CaCl_2}$  and  $Si_{tot}$  were also weakly negatively  
353 correlated and significantly different from 0 for soils on metamorphic rocks. The  
354 negative correlation was moderate ( $r = -0.41$ ) for non-carbonated soils on sediments,  
355 very weak ( $r = -0.16$ ) and poorly significant for soils on igneous intrusive rocks and non  
356 significant for soils on igneous extrusive rocks. Such a negative trend was already  
357 observed in other continents by Yanai et al. (2016) and Meunier et al. (2017).  
358 Nevertheless this negative correlation between  $Si_{CaCl_2}$  by  $Si_{tot}$  was not robust enough  
359 to be used as a proxy for PAS. However, the correlation for podzol was stronger than  
360 the previous ones ( $r = -0.75$ ) and for carbonated soils on sediments, a positive non-  
361 significant correlation was obtained. PCA (Figure 7) allowed disentangling the  
362 observed correlation between  $Si_{tot}$  and various soil characteristics for the different  
363 soil/parent material classes.  $Si_{tot}$  was correlated positively with sand and negatively  
364 with clay for Podzols showing that the content of sand was mainly composed of quartz,  
365 a poorly soluble mineral. Therefore, the best negative correlation given for podzols  
366 between  $Si_{CaCl_2}$  and  $Si_{tot}$  was interpreted as the presence of quartz as the dominant  
367 mineral, which control the poorly soluble pool of  $Si_{tot}$ . For soils on igneous and  
368 metamorphic rocks and non-carbonated soils on sediments, sands were associated  
369 with elements such as K or Na, which may be interpreted as the presence of feldspars,

370 more soluble than quartz and explain the weak correlation between  $Si_{CaCl_2}$  and  $Si_{tot}$ .  
371 Besides, the elements opposed to  $Si_{tot}$  in the PCA may indicate pools of minerals,  
372 which are more favorable to PAS, such as the finer fractions and Al and Fe oxides  
373 (Yanai et al., 2016). For metamorphic rocks for instance,  $Si_{tot}$  was negatively correlated  
374 to Fe and Al concentration as seen by the second axis of the PCA (21 % of the  
375 variance), but independent from sand content. This axis expressed a quartz content  
376 opposed to a secondary mineral content (represented by Al and Fe) and was  
377 independent from the sand content since loess deposits are frequent on metamorphic  
378 rocks in the Brittany region, one of the largest area of metamorphic rocks in France  
379 (Lemercier et al., 2011). For carbonated soils on sediments, the PCA showed that  $Si_{tot}$   
380 were negatively correlated to Ca (carbonates) and positively correlated to the sand  
381 content along the second axis that represented 20 % of the variability. This axis  
382 showed that Si is probably mainly contained in feldspar rather than quartz for  
383 explaining the positive trend between topsoil  $Si_{CaCl_2}$  and  $Si_{tot}$ . Thus,  $Si_{tot}$  may be a good  
384 proxy for  $Si_{CaCl_2}$  in the cases of podzols, but for the other soil classes a more detailed  
385 analysis of the parameters that control  $Si_{CaCl_2}$  should be done.

## 386 4 Conclusions

387 In a territory as diverse as France, we showed that the total topsoil Si concentrations  
388 were highly variable with values ranging from: 22.8 to 456 g.kg<sup>-1</sup> and thus covering  
389 almost the entire range of soil Si concentration recorded in the literature so far.

390 The spatial variability in total topsoil Si concentrations was as expected due to the  
391 diversity of parent material and soil types. However, climate, notably through the  
392 impact of rainfall on weathering intensity, was identified as a driving factor for non-  
393 carbonated soils. No impact of land use and relief could be identified, as well as no

394 impact of the relief at least at this scale. Further work based on a paired site approach  
395 in different pedo-geological contexts is needed to better conclude on the impact of  
396 these factors on the total topsoil Si concentrations.

397 Lastly, topsoil  $Si_{CaCl_2}$  tended to decrease when topsoil  $Si_{tot}$  increase. The relationship  
398 was the strongest in Podzols where Si was mainly contained in quartz. For the other  
399 soil classes, we suggest that  $Si_{tot}$  may only be considered as a proxy for bioavailable  
400 Si if the mineralogical composition is well constrained.

## 401 Acknowledgements

402 This work was performed in the frame of the French ANR BioSiSol project (ANR-14-  
403 CE01-0002). RMQS soil sampling and physico-chemical analyses were supported by  
404 the GIS Sol, which is a scientific group of interest on soils involving the French Ministry  
405 for ecology and sustainable development and Ministry of agriculture, the French  
406 National forest inventory (IFN), ADEME (Agence de l'environnement et de la maîtrise  
407 de l'énergie, which is a French government agency concerned with environmental  
408 protection and energy management), IRD (Institut de recherche pour le  
409 développement, which is a French public research organization dedicated to southern  
410 countries) and INRA (Institut national de la recherche agronomique, which is a French  
411 public research organization dedicated to agriculture s.l.).

412

## 413 1-References

414 Achache, J., Debeglia, N., Grandjean, G., Guillen, A., Le Bel, L., Ledru, P., Renaud,  
415 X., Autran, A., Bonijoly, D., Calcagno, P., Pluchery, E., Guennoc, P., Truffert, C., Rossi,

416 P., Vairon, J., Avouac, J.P., Poli, E., Senechal, G., Brun, J.P., Galdeano, A., Diament,  
417 M., Tarits, P., Mervier, J., Paul, A., Poupinet, G., Marquis, G., Bayer, R., Chautra, J.M.,  
418 1997. GEOFRANCE 3D: l'imagerie géologique et géophysique 3D du sous-sol de la  
419 France. Mém. Société Géologique 53–71.

420 Ameijeiras-Marino, Y., Opfergelt, S., Derry, L.A., Robinet, J., Govers, G., Minella,  
421 J.P.G., Delmelle, P., 2018. Ge/Si ratios point to increased contribution from deeper  
422 mineral weathering to streams after forest conversion to cropland. Applied Geochem.  
423 96, 24-34.

424 Baize, D., 2000. Guide des analyses en pédologie, Quae. ed. INRA Editions, Paris  
425 (France).

426 Barão, L., Clymans, W., Vandevenne, F., Meire, P., Conley, D.J., Struyf, E., 2014.  
427 Pedogenic and biogenic alkaline-extracted silicon distributions along a temperate land-  
428 use gradient. Eur. J. Soil Sci. 65, 693–705.

429 Berner, E.K., Berner, R.A., 1996. Global environment: water, air, and geochemical  
430 cycles. Printice Hall, Upper Saddle River, New Jersey.

431 Bluth, G., J., S., Kumpr, L.R., 1994. Lithologic and climatologic controls of river  
432 chemistry. Geochim. Cosmochim. Acta 58, 2341-2359.

433 Breiman, L., 2001. Random Forests. Mach. Learn. 45, 5–32.  
434 <https://doi.org/10.1023/A:1010933404324>

435 Cerdan, O., Govers, G., Le Bissonnais, Y., Van Oost, K., Poesen, J., Saby, N., Gobin,  
436 A., Vacca, A., Quinton, J., Auerswald, K., Klik, A., Kwaad, F.J.P.M., Raclot, D., Ionita,  
437 I., Rejman, J., Rousseva, S., Muxart, T., Roxo, M.J., Dostal, T., 2010. Rates and spatial

438 variations of soil erosion in Europe: A study based on erosion plot data.  
439 *Geomorphology* 122, 167–177. <https://doi.org/10.1016/j.geomorph.2010.06.011>

440 Cornelis, J.T., Delvaux, B., 2016. The functional role of silicon in plant biology. Soil  
441 processes drive the biological silicon feedback loop. *Functional Ecology* 30, 1298-  
442 1310.

443 Coskun, D., Deshmukh, R., Sonah, H., Menzies, J.G., Reynolds, O., Ma, J.F.,  
444 Kronzucker, H.J., Bélanger, R.R., 2019. The controversies of silicon's role in plant  
445 biology. *New Phytologist* 221, 67-85.

446 De Vos, W., Tarvainen, T., Salminen, R., Reeder, S., De Vivo, B., Demetriades, A.,  
447 Pirc, S., Batista, M.J., Marsina, K., Ottesen, R.T., 2006. *Geochemical atlas of Europe.*  
448 *Part 2. Interpret. Geochem. Maps Addit. Tables Fig. Maps Relat. Publ. Geol. Surv. Finl.*  
449 *Espoo.*

450 Dere, A.L., White, T.S., April, R.H., Brantley, S.L., 2016. Mineralogical transformations  
451 and soil development in shale across a latitudinal climosequence. *Soil Sci. Soc. Am.*  
452 *J.* 80, 623-636.

453 Dürr, H.H., Meybeck, M., Hartmann, J., Laruelle, G.G., Roubéix, V., 2011. Global  
454 spatial distribution of natural riverine silica inputs to the coastal zone. *Biogeosciences*  
455 8, 5978-620

456 Egli, M., Mirabella, A., Sartori, G., Fitze, P., 2003. Weathering rates as a function of  
457 climate: results from a climosequence of the Val Genova (Trentino, Italian Alps),  
458 *Geoderma* 111, 99-121.

459 European Environment Agency, 2007. CLC2006 technical guidelines. Publications  
460 Office, Luxembourg.

461 Faroux, S., Kaptué Tchuenté, A.T., Roujean, J.-L., Masson, V., Martin, E., Moigne,  
462 P.L., 2013. ECOCLIMAP-II/Europe: A twofold database of ecosystems and surface  
463 parameters at 1 km resolution based on satellite information for use in land surface,  
464 meteorological and climate models. *Geosci. Model Dev.* 6, 563–582.

465 Frings, P.J., Clymans, W., Fontorbe, G., Gray, W., Chakrapani, G.J., Conley, D.J., De  
466 La Rocha, C., 2015. Silicate weathering in the Ganges alluvial plain. *Earth Planet. Sci.*  
467 *Let.* 427, 136-148.

468 Gray, J.M., Bishop, T.F.A., Wilford, J.R., 2016. Lithology and soil relationship for soil  
469 modelling and mapping. *Catena* 147, 429-440.

470 Lucas, Y., Luizao, F.J., Chauvel, A., Rouiller, J., Nahon, D., 1993. The relation between  
471 biological activity of the rain forest and mineral composition of soils. *Science* 260, 521-  
472 523.

473 Matheron, G., 1971. The Theory of Regionalised Variables and its Applications. *Cah.*  
474 *Cent. Morphol. Math.* 5, 212.

475 Gneiting, T., Balabdaoui, F., Raftery, A.E., 2007. Probabilistic forecasts, calibration  
476 and sharpness. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 69, 243–268.

477 Guntzer, F., Keller, C., Poulton, P.R., McGrath, S.P., Meunier, J.-D., 2012. Long-term  
478 removal of wheat straw decreases soil amorphous silica at Broadbalk, Rothamsted.  
479 *Plant Soil* 352, 173–184.

480 Haysom, M.B.C., Chapman, L.S., 1975. Some aspects of the calcium silicate trials at  
481 Mackay. Proc. Qld. Soc. Sugar Cane Technol. 42, 117-122.

482 Henriot, C., Bodarwé, L., Dorel, M., Draye, X., Delvaux, B., 2008. Leaf silicon content  
483 in banana (*Musa* spp.) reveals the weathering stage of volcanic ash soils in  
484 Guadeloupe. Plant Soil 313, 71.

485 Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., Jarvis, A., 2005. Very high  
486 resolution interpolated climate surfaces for global land areas. Int. J. Climatol. 25, 1965–  
487 1978. <https://doi.org/10.1002/joc.1276>

488 Huete, A., Didan, K., Miura, T., Rodriguez, E.P., Gao, X., Ferreira, L.G., 2002.  
489 Overview of the radiometric and biophysical performance of the MODIS vegetation  
490 indices. Remote Sens. Environ., The Moderate Resolution Imaging Spectroradiometer  
491 (MODIS): a new generation of Land Surface Monitoring 83, 195–213.  
492 [https://doi.org/10.1016/S0034-4257\(02\)00096-2](https://doi.org/10.1016/S0034-4257(02)00096-2)

493 Institut National de l'Information Géographique et Forestière, 2006. Base de Données  
494 Forêt.

495 James, G., Witten, D., Hastie, T., Tibshirani, R., 2013. An Introduction to Statistical  
496 Learning, Springer Texts in Statistics. Springer New York, New York, NY.  
497 <https://doi.org/10.1007/978-1-4614-7138-7>

498 Jolivet, C., Arrouays, D., Boulonne, L., Ratié, C., Saby, N., 2006. Le réseau de  
499 mesures de la qualité des sols de France (RMQS). Etat D'avancement Prem. Résultats  
500 Etude Gest. Sols 13, 149–164.

501 Joly, D., Brossard, T., Cardot, H., Cavailhes, J., Hilal, M., Wavresky, P., 2010. Les  
502 types de climats en France, une construction spatiale. *Cybergeo Eur. J. Geogr.*  
503 <https://doi.org/10.4000/cybergeo.23155>

504 Keskin, H., Grunwald, S., 2018. Regression kriging as a workhorse in the digital soil  
505 mapper's toolbox. *Geoderma* 326, 22–41.  
506 <https://doi.org/10.1016/j.geoderma.2018.04.004>

507 King, D., Jones, R.J.A., Thomasson, A.J., 1995. European land information systems  
508 for agro-environmental monitoring.

509 Kuhn, M., Steve, W., Chris, K., Nathan, C., code), R.Q. (Author of imported C., code),  
510 R.R.P.L. (Copyright holder of imported C., 2017. *Cubist: Rule- And Instance-Based*  
511 *Regression Modeling*.

512 Kursa, M.B., 2014. Robustness of Random Forest-based gene selection methods.  
513 *BMC Bioinformatics* 15, 8. <https://doi.org/10.1186/1471-2105-15-8>

514 Kursa, M.B., Rudnicki, W.R., 2010. Feature selection with the Boruta package. *J Stat*  
515 *Softw* 36, 1–13.

516 Landré, A., Saby, N.P.A., Barthès, B.G., Ratié, C., Guerin, A., Etayo, A., Minasny, B.,  
517 Bardy, M., Meunier, J.-D., Cornu, S., 2018. Prediction of total silicon concentrations in  
518 French soils using pedotransfer functions from mid-infrared spectrum and pedological  
519 attributes. *Geoderma* 331, 70–80.

520 Le Bas, C., Barthès, S., Boutefoy, I., Fort, J., Scheurer, O., Darracq, S., Lacassin, J.,  
521 Sauter, J., 2004. Utilisation des données sols d'I.G.C.S. en France : Un état des lieux.  
522 *Étude Gest. Sols* (2), 8.



523 Liang, Y., Nikolic, M., Bélanger, R., Gong, H., Song, A., 2015. Silicon in agriculture.  
524 Dordr. Springer Doi 10, 978–94.

525 Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest 2, 6.

526 Loiseau, T., Chen, S., Mulder, V.L., Román Dobarco, M., Richer-de-Forges A.C.,  
527 Lehmann, S., Bourenane, H., Saby N.P.A., Martin, M.P., Vaudour, E., Gomez, C.4  
528 Lagacherie, P., Arrouays D. 2019, Satellite data integration for soil clay content  
529 modelling at a national scale, International Journal of Applied Earth Observations and  
530 Geoinformation, 82 , 108205

531 Magidson, J., 1981. Qualitative variance, entropy, and correlation ratios for nominal  
532 dependent variables. Soc. Sci. Res. 10, 177–194. [https://doi.org/10.1016/0049-](https://doi.org/10.1016/0049-089X(81)90003-X)  
533 [089X\(81\)90003-X](https://doi.org/10.1016/0049-089X(81)90003-X)

534 McBratney, A.B., Santos, M.M., Minasny, B., 2003. On digital soil mapping. Geoderma  
535 117, 3–52.

536 Meunier, J.D., Sandhya, K., Prakash, N. B., Borschneck, D., Dussouillez, P., 2018.  
537 pH as a proxy for estimating plant-available Si? A case study in rice fields in Karnataka  
538 (South India). Plant and Soil 432, 143-155.

539 Minasny, B., McBratney, A.B., 2006. A conditioned Latin hypercube method for  
540 sampling in the presence of ancillary information. Comput. Geosci. 32, 1378–1388.

541 Minasny, B., McBratney, A.B., Hartemink, A.E., 2010. Global pedodiversity, taxonomic  
542 distance, and the World Reference Base. Geoderma 155, 132–139.  
543 <https://doi.org/10.1016/j.geoderma.2009.04.024>

544 Minasny, B., Tranter, G., McBratney, A.B., Brough, D.M., Murphy, B.W., 2009.  
545 Regional transferability of mid-infrared diffuse reflectance spectroscopic prediction for  
546 soil chemical properties. *Geoderma* 153, 155–162.  
547 <https://doi.org/10.1016/j.geoderma.2009.07.021>

548 Moriasi, D.N., Arnold, J.G., Liew, M.W.V., Bingner, R.L., Harmel, R.D., Veith, T.L.,  
549 2007. Model Evaluation Guidelines for Systematic Quantification of Accuracy in  
550 Watershed Simulations. *Trans. ASABE* 50, 885–900.  
551 <https://doi.org/10.13031/2013.23153>

552 Mulder, V.L., Lacoste, M., Richer-de-Forges, A.C., Arrouays, D., 2016. GlobalSoilMap  
553 France: High-resolution spatial modelling the soils of France up to two meter depth.  
554 *Sci. Total Environ.* 573, 1352–1369. <https://doi.org/10.1016/j.scitotenv.2016.07.066>

555 Narayanaswamy, C., Prakash, N.B., 2010. Evaluation of selected extractants for plant-  
556 available silicon in rice soils of Southern India. *Comm. Soil Sc. PI Anal.*, 41, 977-989.

557 Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part  
558 I — A discussion of principles. *J. Hydrol.* 10, 282–290. [https://doi.org/10.1016/0022-](https://doi.org/10.1016/0022-1694(70)90255-6)  
559 [1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)

560 NCAR-Research Applications Laboratory, 2015. verification: Weather Forecast  
561 Verification Utilities.

562 Pebesma, E., Graeler, B., 2019. *gstat: Spatial and Spatio-Temporal Geostatistical*  
563 *Modelling, Prediction and Simulation.*

564 Quinlan, J.R., 1992. Learning with continuous classes, in: 5th Australian Joint  
565 Conference on Artificial Intelligence. Singapore, pp. 343–348.

566 Quintana-Seguí, P., Le Moigne, P., Durand, Y., Martin, E., Habets, F., Baillon, M.,  
567 Canellas, C., Franchisteguy, L., Morel, S., 2008. Analysis of Near-Surface  
568 Atmospheric Variables: Validation of the SAFRAN Analysis over France. *J. Appl.*  
569 *Meteorol. Climatol.* 47, 92–107. <https://doi.org/10.1175/2007JAMC1636.1>

570 R. Core Team, 2016. R: A language and environment for statistical computing. Version  
571 3.3. 1. 2016.

572 Reimann, C., Birke, M., Demetriades, A., Filzmoser, P., O'Connor, P., 2014. Chemistry  
573 of Europe's agricultural soils, part A.

574 Rodrigues, F.A., Datnoff, L.E., 2015. Silicon and plant diseases. Springer.

575 Roudier, P., 2017. *clhs: Conditioned Latin Hypercube Sampling*.

576 Siegel, S., Castellan, N.J., 1988. *Nonparametric Statistics for the Behavioral Science*,  
577 2nd edition. ed. McGraw-Hill Education, New York; etc.

578 Sommer, M., Kaczorek, D., Kuzyakov, Y., Breuer, J., 2006. Silicon pools and fluxes in  
579 soils and landscapes—a review. *J. Plant Nutr. Soil Sci.* 169, 310–329.  
580 <https://doi.org/10.1002/jpln.200521981>

581 Struyf, E., Smis, A., Van Damme, S., Garnier, J., Govers, G., Van Wesemael, B.,  
582 Conley, D.J., Batelaan, O., Frot, E., Clymans, W., 2010. Historical land use change  
583 has lowered terrestrial silica mobilization. *Nat. Commun.* 1, 129.

584 Taboada, T., Ferro-Vazquez, C., Stoops, G., Martinez Cortizas, A., Rodriguez Flores,  
585 R., Rodriguez-Laco, L., 2019. Secondary aluminium, iron and silica phases across a

586 volcanic soil climosequence, Galapagos Islands. *Europ. J. Soil Sci.*, 70,  
587 <https://doi.org/10.1111/ejss.12788>

588 USGS, 2006. Shuttle Radar Topography Mission, 1 Arc Second Scene  
589 SRTM\_u03\_n008e004, Unfilled Unfinished 2.0. Shuttle Radar Topography Mission.

590 Vandevenne, F.O., Barao, L., Ronchi, B., Govers, G., Meire, P., Kelly, E.F., Struyf, E.,  
591 2015. Silicon pools in human impacted soils of temperate zones. *Global Geochemical*  
592 *Cycles*, 29, doi:10.1002/2014GB005049.

593 Viscarra Rossel, R.A., Walvoort, D.J.J., McBratney, A.B., Janik, L.J., Skjemstad, J.O.,  
594 2006. Visible, near infrared, mid infrared or combined diffuse reflectance spectroscopy  
595 for simultaneous assessment of various soil properties. *Geoderma* 131, 59–75.  
596 <https://doi.org/10.1016/j.geoderma.2005.03.007>

597 Wedepohl, K.H., 1995. The composition of the continental crust. *Geochim.*  
598 *Cosmochim. Acta* 59, 1217–1232.

599 White, A. F., Vivit, D.V., Shulz, M. S., Bullen, T.D., Evett, R.R., Aagarwal, J., 2012.  
600 Biogenic and pedogenic controls on Si distributions and cycling in grasslands of the  
601 Santa Cruz soil chronosequence, California. *Geochim. Cosmochim. Acta* 94, 72-94.

602 Yanai, J., Taniguchi, H., Nakao, A., 2016. Evaluation of available silicon content and  
603 its determining factors of agricultural soils in Japan. *Soil Sci. Plant Nutr.* 62, 511-518.

604 Yang, J.L., Zhang, G.L., 2018. Silicon cycling by plant and its effects on soil Si  
605 translocation in a typical subtropical area. *Geoderma* 310, 89-98.

606

607

608

## Captions

Figure 1: Variability of  $Si_{tot}$ . (a) Relationship between the concentration in total Si concentration and in carbonate content in French topsoils. The different soils developed on sedimentary rocks are reported in yellow while those developed on igneous and metamorphic rocks are reported in grey; (b) Separation between carbonated soils and non-carbonated soils on sedimentary rocks based on the geo-pedological classification. All the soil groups with mean carbonate content lower than 5 % were considered as non-carbonated soils.

Figure 2: Boxplots of  $Si_{tot}$  for the main geological domains and podzols. The boxes represent the interquartile range, the bold horizontal segment the median  $Si_{tot}$ , the whiskers, the 99 % range and the black dots the outliers. Letters above boxplots represent groups with  $Si_{tot}$  significantly different and numbers, the sample size.

Figure 3: Boxplots of  $Si_{tot}$  for all the considered pedo-geological groups. The boxes represent the interquartile range, the bold horizontal segment the median  $Si_{tot}$ , the whiskers, the 99 % range and the black dots the outliers. Colors represent distinction of main geological domains plus podzols.

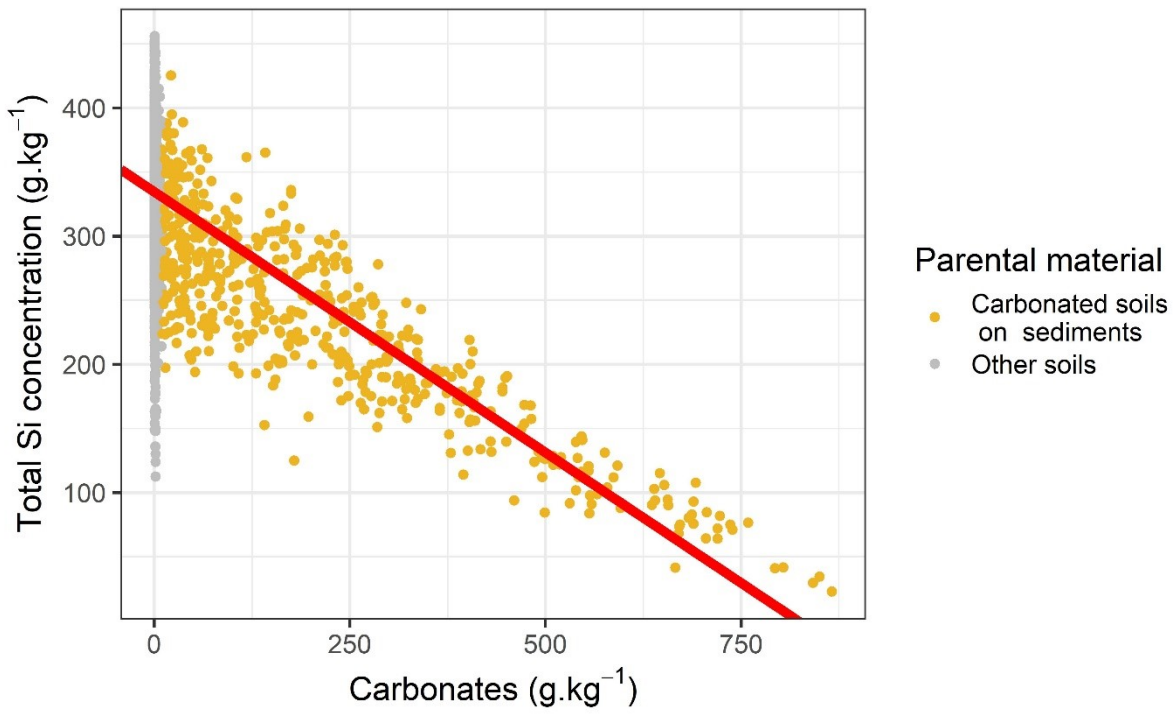
Figure 4:  $Si_{tot}$  spatial distribution from digital soil mapping approach (DSM).

Figure 5: Variable importance in the random forest Model. The X-axis is the Mean Decrease Accuracy computed as the increase percentage in mean squared error of predictions (estimated with out-of-bag-Cross Validation) as a result of variable j being permuted. The definition of the environmental covariates are listed in Table 2.

Figure 6: Plots of the total topsoil Si concentration versus the 30 year average of annual rainfall for the different parent material domains. Soils on sedimentary parent materials were sorted according to the nature of the parent material for the carbonated soils (shalk, marl, limestone and other) and of the texture of the parent material for non-carbonated soils (clay, loam, sand). The black lines represent the fitted linear regression. The Pearson correlation is also provided per parent material domains. Stars represents the level of signification of the correlation: one star being significant at a 5 % confidence level, two stars at 1 % and three stars at 1 %.

Figure 7: Variables correlation circle for the first two components of the Principal Component Analyzes (PCA) performed for each main geological domain and podzols. In blue, Si considered as a passive variable in these PCAs.

a



b

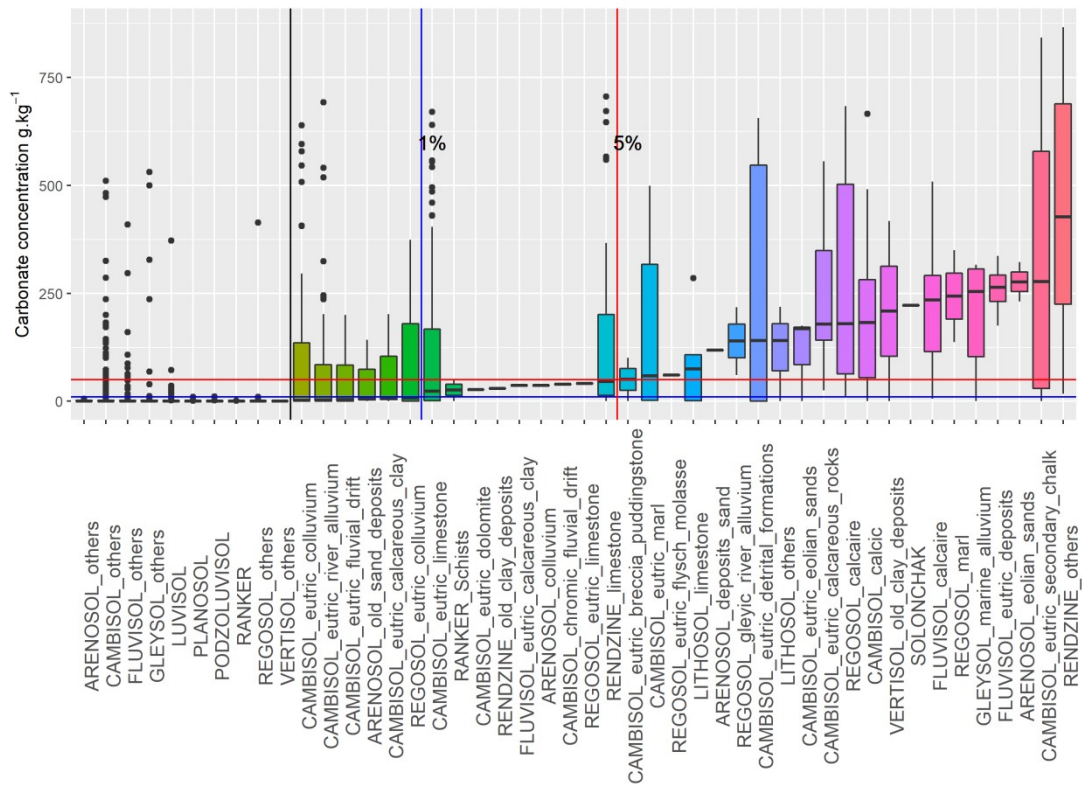


Figure 1

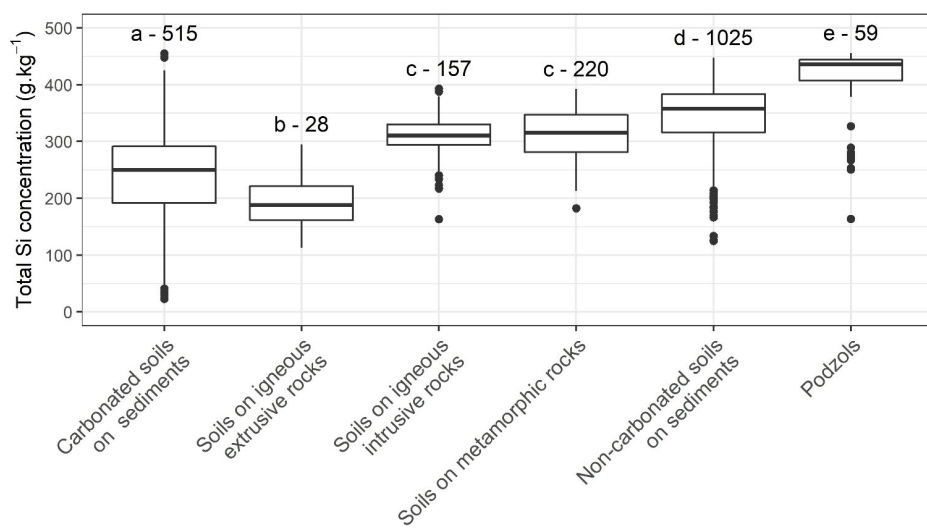


Figure 2



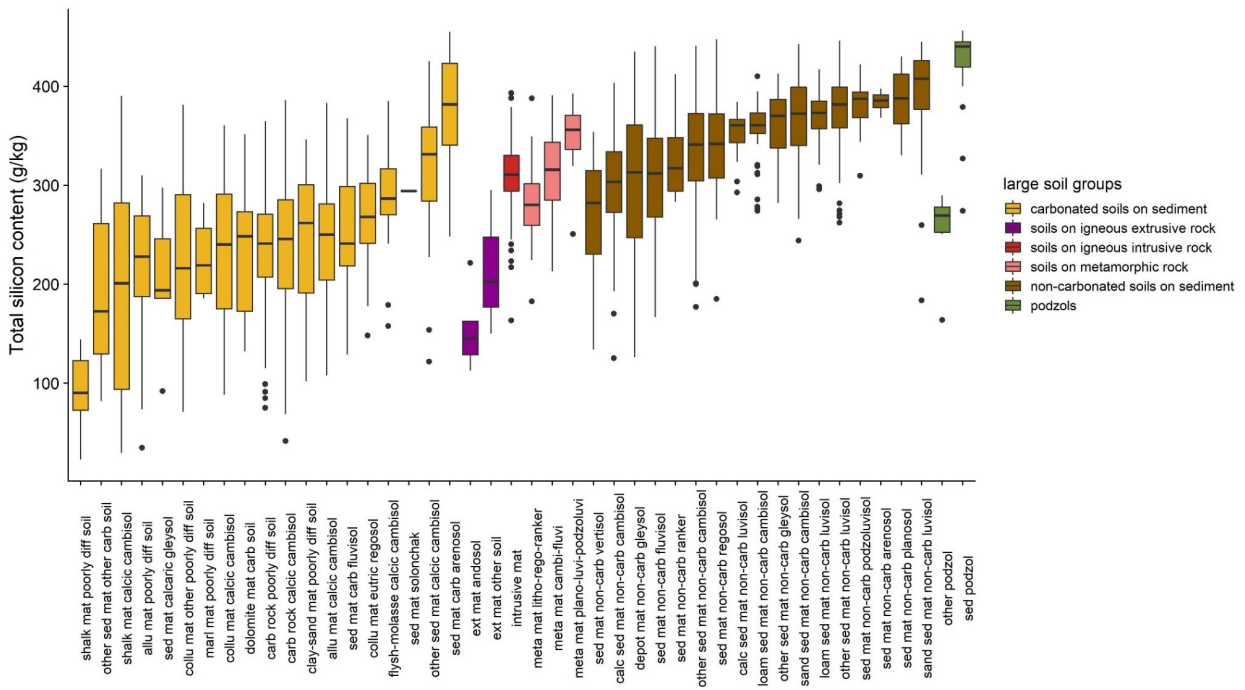


Figure 3

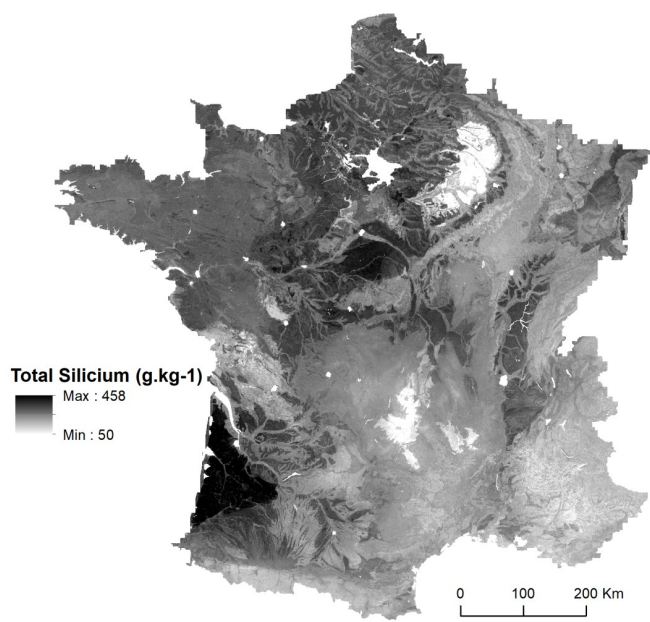


Figure 4

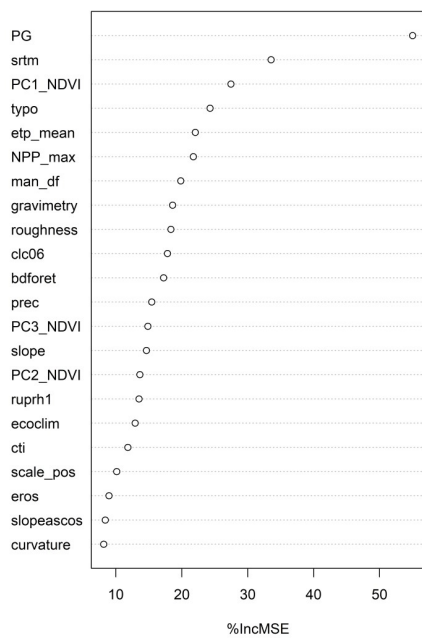


Figure 5

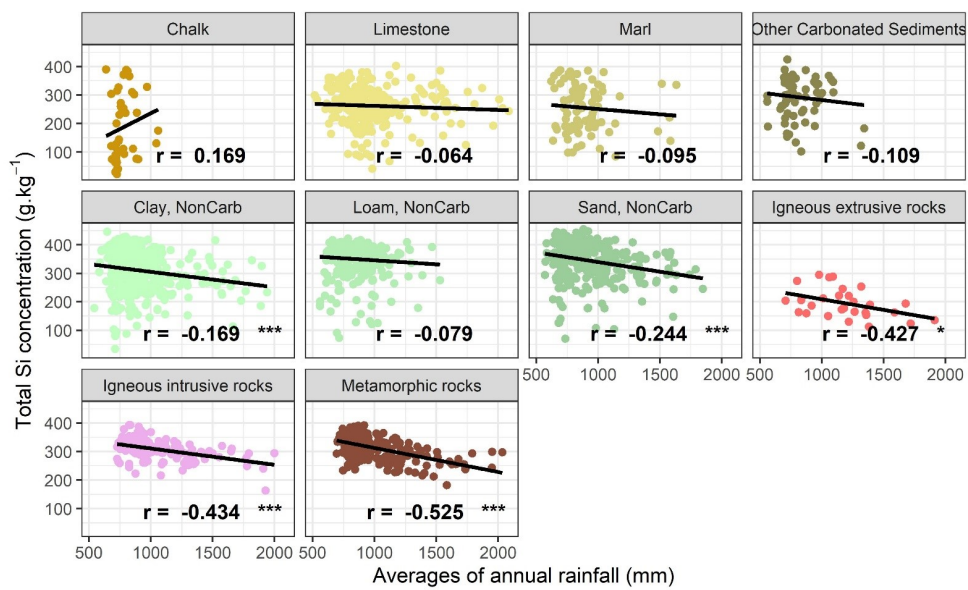


Figure 6

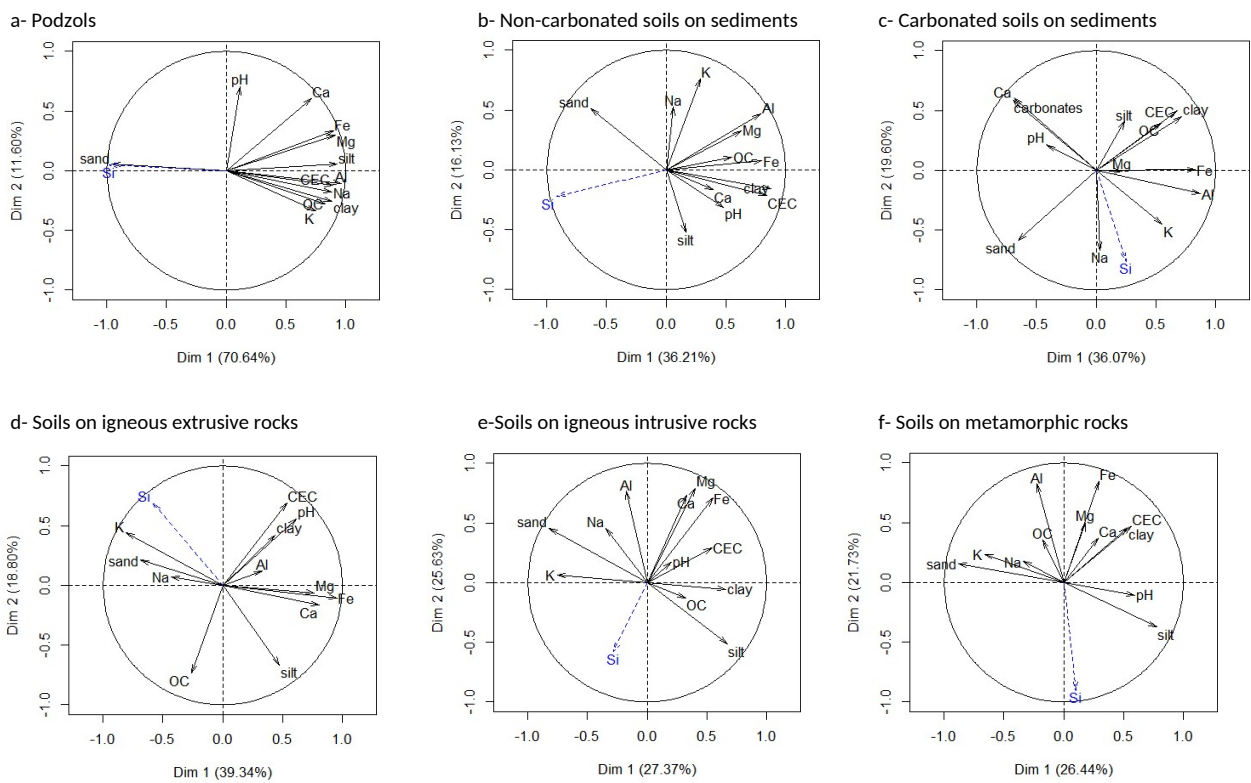


Figure 7

**Table 1: Exhaustive categorical and continuous covariates used for the Digital Soil Mapping approach**

Pedological factor	abbreviates	variables	scale	type	Reference
Vegetation	bdforet	Forest type		category	Institut National de l'Information Géographique et Forestière (2006)
	ecoclim	Ecoclimap land use	1km	category	Faroux et al. (2013)
	NPP_max	Net Primary Production	1km	quantitative	LPDAAC (2001)
	clc06	land use (Corine Land Cover)	250m	category	European Environment Agency (2007)
	PC1_NDVI	1st component of Normalized Difference Vegetation Index	90m	quantitative	Loiseau et al. (2019)
	PC2_NDVI	2nd component of Normalized Difference Vegetation Index	90m	quantitative	Loiseau et al. (2019)
	PC3_NDVI	3rd component of Normalized Difference Vegetation Index	90m	quantitative	Loiseau et al. (2019)
Relief	cti	Compound Topographic Index (SRTM)	90m	quantitative	USGS (2006)
	roughness	roughness (SRTM)	90m	quantitative	USGS (2006)
	curvature	curvature (SRTM)	90m	quantitative	USGS (2006)
	scale_pos	slope position (SRTM)	90m	quantitative	USGS (2006)
	slope	slope (SRTM)	90m	quantitative	USGS (2006)
	slopascos	slope cosinus (SRTM)	90m	quantitative	USGS (2006)
	srtm	elevation DEM (SRTM)	90m	quantitative	USGS (2006)
Climate	eros	Erosion rate	1 :1000000	quantitative	Cerdan et al. (2010)
	prec	mean annual precipitation	~1km	quantitative	Hijmans et al. (2005)
	Typo	climate typology	250 m	category	Joly et al. (2010)
Parent material	etp_mean	mean annual potential evapotranspiration	8km	quantitative	Quintana-Segui et al. (2008)
Soil	gravimetry	free-air Bouguer anomaly	4km	quantitative	Achache et al. (1997)
Soil	ruprh1	available water capacity	1 :1000000	quantitative	Le Bas et al. (2004)
Soil and parent material	PG	Pedological and Geological classification	1 :1000000	category	King et al. (1995)

Table 2: Pearson correlation coefficient between  $Si_{CaCl_2}$  and  $Si_{tot}$  in French topsoils. Significance levels: \*\*\*:  $p$ -value < 0.001; \*\*:  $p$ -value < 0.01; \*:  $p$ -value < 0.05

Classes	n	r
Carbonated soils on sediments	508	0.12***
Soils on igneous extrusive rocks	28	-0.25
Soils on igneous intrusive rocks	155	-0.16*
Soils on metamorphic rocks	217	-0.26***
Non-carbonated soils on sediments	1020	-0.41***
Podzols	58	-0.75***
All data	1986	-0.32***

# DO RAINFALL AND LAND USE AFFECT THE POOL OF TOTAL SILICON CONCENTRATION? A DIGITAL SOIL MAPPING APPROACH OF FRENCH TOPSOILS.

Landré, A.<sup>a,b</sup>, Cornu, S.<sup>c</sup>, Meunier, J.-D.<sup>c</sup>, Guerin A.<sup>d</sup>, Arrouays D.<sup>a</sup>, Caubet M.<sup>a</sup>, Ratié C.<sup>a</sup>, Saby, N.P.A.<sup>a</sup>

a) INRA, Infosol, US 1106, Orléans, France

b) INRA - INPT-ENSAT - INPT-EI-Purpan, UMR 1248 AGIR AGroécologie, Innovations, teRritoires. Centre de recherche Occitanie-Toulouse, Auzeville, France.

c) Aix-Marseille Univ, CNRS, IRD, Coll de France, INRA, CEREGE, Aix-en-Provence, France

d) INRA, Laboratoire d'Analyses des Sols US, Arras, France

**Declarations of interest: none**



# Supplementary file: Determination of the total Si concentration ( $Si_{tot}$ ) based on the concentration in major elements

Landre et al. (2018) analyzed  $Si_{tot}$  for a subset of 673 samples, among the 2007 RMQS samples.  $Si_{tot}$  for the remaining sites was modelled using a Cubist modelling approach, as follow:

$$Si = f(\text{Al, Fe, K, Na, } Ca_{nc}, Mg_{nc}, \text{P, Mn, OC, CaCO}_3, \text{residual water}) \quad (S1)$$

Where  $f$  is the cubist regression trees model,

$Ca_{nc}$  and  $Mg_{nc}$  are the fraction of calcium and magnesium, respectively, that are not included in carbonate minerals nor adsorbed on the exchangeable surfaces. They were estimated as follow

$$Ca_{nc} = Ca_{tot} - \left( \left( \frac{CaCo_3}{10} * 0.401 \right) + Ca_{exch} \right) \quad \text{with} \quad Ca_{nc} \geq 0 \quad (S2)$$

$$Mg_{nc} = Mg_{tot} - (Mg_{ech} + Ca_{excess}) \quad \text{with} \quad Mg_{nc} \geq 0 \quad (S3)$$

Where  $Ca_{excess} = Ca_{nc}$  when  $Ca_{nc} \leq 0$  and  $Ca_{nc} = 0$  otherwise

The Cubist model is a form of regression rules that build unconventional regression trees, with final nodes containing linear models instead of discrete values (Quinlan, 1992). Cubist creates comprehensible rules that describe the relationships between predictive variables (here spectra or soil properties) and the variable of interest.

In order to calibrate and validate the models, we used a repeated cross-validation approach combined to a bootstrap step (James et al., 2013). Repeated cross-validation allows assessing the quality of the prediction. It involves randomly dividing the available set of data into two parts, a calibration set and a validation set. We repeated this operation 10 times using a split of 75%-25% for calibration and validation respectively. This subdivision was performed using the conditioned Latin Hypercube

Sampling (cLHS) method (Minasny and McBratney, 2006). This method is a stratified random procedure that provides an efficient way of sampling variables from their multivariate distributions. The bootstrap step allows assessing the uncertainty of the model predictions. For this, we simulated 100 datasets by random sampling with replacement 95 % of the calibration dataset from the cross-validation step. This procedure provided then 100 Cubist models. The final prediction is obtained by averaging the predictions of 100 bootstrapped models. Our modelling approach involves a large number of calibration and validation operations, and we thus used a parallel processing approach to overcome the computational load. It is implemented in R using the packages `foreach`, `doParallel` or `snow`. We used the Cubist model implemented in the R package `Cubist` (Kuhn et al., 2016), the `cLHS` function implemented in the R package `clhs` (Roudier, 2011) and the `crps` function implemented in the R package `verification` (Laboratory NCAR-Research Applications, 2015).

The model was assessed using three conventional performance indicators: the coefficient of determination ( $R^2$ ), the root mean square error (RMSE) also known as standard error of prediction (SEP) and the bias, which is the mean residual of the model. We also took into account the probabilistic characteristic of model predictions by using the continuous rank probability score average (CRPS, equation 6). The CRPS represents the closeness between the prediction distribution and the corresponding observations (Gneiting et al., 2007).

$$CRPS = \int_{-\infty}^{\infty} BS(y) dy \quad (s4)$$

$$BS(y) = \frac{1}{n} \sum_{i=1}^n \{(F_i(y) - \mathbb{1}(x_i \leq y))\}^2 \quad (s5)$$

Where  $BS(y)$  denotes the Brier (1950) score for probability forecasts of the binary event at the threshold value  $y \in \mathbb{R}$ ,  $x$  the observation and  $y$  the model prediction,  $n$  the number of samples,  $F$  the cumulative distribution function (CDF) of  $X$ , a random variable, such as  $F(y) = P[X \leq y]$  and  $\mathbb{1}$  is the Heaviside step function. This function is a discontinuous function whose value is zero for negative argument and one for positive argument.

The  $Si_{tot}$  prediction function is yielding  $R^2$  greater than 0.98 with a very small variance (0.0012) among repetitions of the cross-validation, RMSE lower than  $11 \text{ g kg}^{-1}$  and CRPS lower than  $7 \text{ g kg}^{-1}$  for measured  $Si_{tot}$  ranging from 22.81 to  $455.8 \text{ g kg}^{-1}$  over the RMQS database with a median equal to  $323.6 \text{ g kg}^{-1}$ . These results are better than the results obtained by Landre et al. (2018) with their MIRS PTF and suggest that our PTF can predict with higher accuracy the Si content.

To better figure out the accuracy of our PTF, we plotted in Figure S1, the predicted *versus* measured  $Si_{tot}$  for only one iteration of the cross-validation step. In this figure, the predicted *versus* measured  $Si_{tot}$  distribution closely follow the one to one line with a smaller prediction uncertainty than the analytical uncertainty. Those results allow us considering that the predictions of this function are as good as Si measurements since the prediction uncertainties are within the range of analytical values. However, as explained by Landre et al. (2018), the analytical uncertainty was not taken into account in the prediction uncertainty calculation, as they were not always available. Despite this, the obtained PTF shows exceptional accuracy that is rarely obtained in environmental fields for PTFs (Minasny et al., 2009; Viscarra Rossel et al., 2006) (see Landre et al. (2018) for more information).

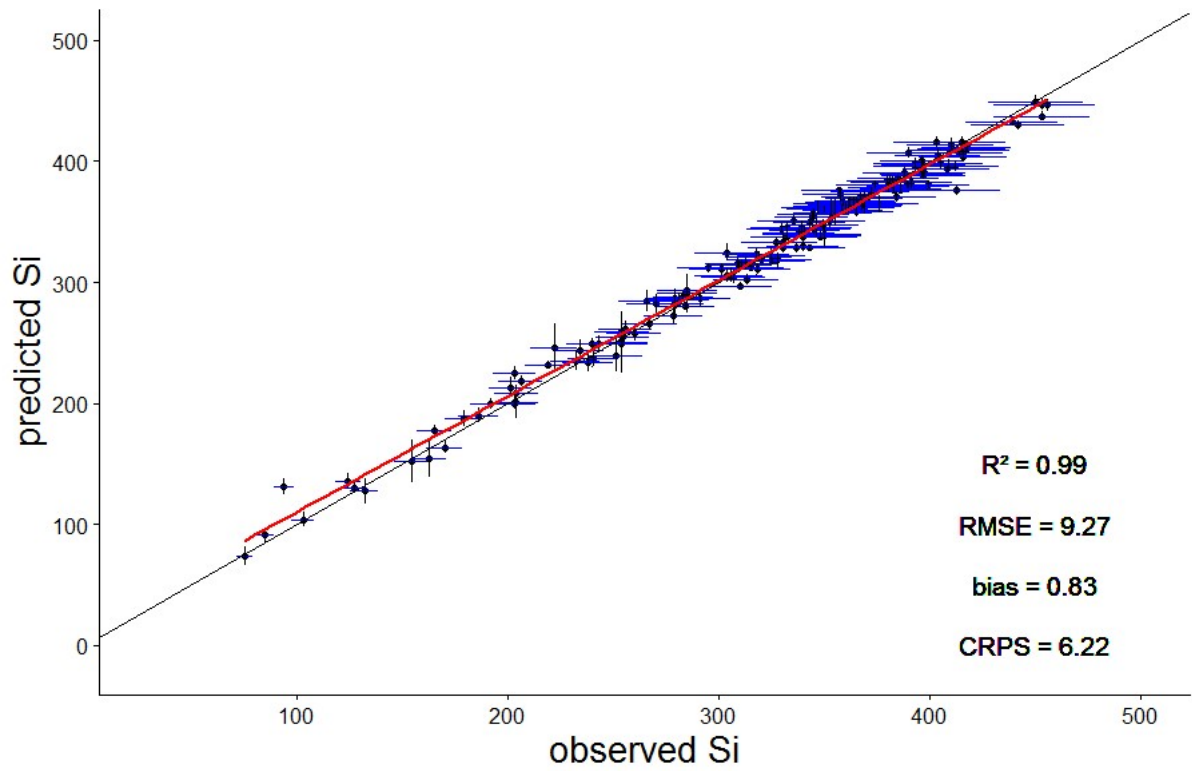


Figure S1: Predicted versus measured Si concentrations (in  $\text{g kg}^{-1}$ ) for the first cross-validation replication of our Si concentration predictive function. In black, the one to one line and in red, the regression line. Black vertical error bars represent the prediction's uncertainty and blue horizontal error bars represent the analytical uncertainty.

# DSM model validation

The model performance was estimated using a cross validation. four conventional performance indicators were computed: the coefficient of determination ( $R^2$ ), the mean square error (MSE), the root mean square error (RMSE, also known as standard error of prediction, SEP) and the bias, which is the mean residual of the model. In addition, a fifth indicator was calculated, the concordance which is a normalized statistic that determines the relative magnitude of the residual variance compared to the measured data variance (Moriasi et al., 2007; Nash and Sutcliffe, 1970). The results of the cross validation of the SCORPAN model are in accordance with previous digital soil mapping exercise of topsoil properties (Mulder et al., 2016). We found an  $R^2$  of 0.45, a MSE of 3647.99, a RMSE of 60.4 g  $\text{kg}^{-1}$ , an important bias of 4.70 and a high concordance of 0.59.