

**Relationships between particle size distribution and VNIR reflectance spectra are weaker for soils formed from bedrock compared to transported parent materials**

B. G. RAWLINS<sup>a</sup>, S. J. KEMP<sup>a</sup>, A. E. MILODOWSKI<sup>a</sup>

<sup>a</sup>*British Geological Survey, Keyworth, Nottingham NG12 5GG, UK*

Running heading: *Particle size distribution and VNIR reflectance spectra*

Correspondence: B. G. Rawlins. E-mail: [bgr@bgs.ac.uk](mailto:bgr@bgs.ac.uk)

## 1 ABSTRACT

2 The cost of determining particle size distribution (psd) of the soil can be significantly  
3 reduced by using statistical relationships between visible and near infra red diffuse re-  
4 flectance spectra (VNIR-DRS) and the proportions of the three size fractions (sand,  
5 silt and clay). The spectra contain information on the quantities of soil minerals which  
6 occur in these fractions. Statistical models for estimating psd based on a set of soil  
7 samples from common parent materials (PM) – with similar mineralogy – may provide  
8 more accurate predictions than more comprehensive, global models. The aim of this  
9 paper is to compare the performance of statistical models for the prediction of psd from  
10 VNIR-DRS for soils with differing types of parent material; specifically soils derived  
11 directly from bedrock (coal-bearing and mudstone-bearing strata) or from transported  
12 parent materials (glacial till, glacio-lacustrine deposits and alluvium) across eastern  
13 England. We assessed the accuracy of psd predictions using partial least squares re-  
14 gression (PLSR) models between two additive log ratios of the three size fractions  
15 and VNIR-DRS. We also formed a global PLSR model from all five soil groups. We  
16 used mean residual prediction deviation (RPD) from repeated (n=100) cross-validation  
17 to compare the performance of the models because it accounts for the magnitude of  
18 variation in the sample data. The most accurate models for the clay (RPD range  
19 1.82–2.33) and sand fractions (RPD range 1.71–1.94) were for soils developed over the  
20 transported PM; the models for soils developed over bedrock were substantially poorer  
21 (clay RPD range 1.33–1.68; sand RPD range 1.34–1.39). The RPD values for the silt  
22 fraction models were smaller, but the same distinction between transported (better;  
23 RPD range 1.4–1.88) and bedrock derived soils (poorer; RPD range 1.15–1.25) was  
24 observed. The global model had intermediate RPD values for the three size fractions  
25 (clay=1.75, silt=1.76 and sand=1.74). Of the five groups, the soils developed from  
26 glacio-lacustrine deposits had the largest mean sand size fraction (58%), but also the  
27 most accurate models for estimation of clay and sand size fractions. Due to sedimen-  
28 tary transport and deposition, the mineralogy of the soils developed from Quaternary

29 substrates may be more homogeneous than the bedrock-derived soils, which may in  
30 part account for the more accurate models developed for the former. To date we do  
31 not have sufficient evidence to demonstrate this unequivocally.

## 1. Introduction

The ability of soil scientists to map the spatial variation of particle size distribution (psd) accurately at fine scales is important because psd contributes to the soil's hydraulic behaviour and water storage, its handling characteristics under tillage and its susceptibility to erosion. When psd is measured in the laboratory the results are often expressed as the proportions of three grain size-fractions (e.g. clay:  $<2\mu\text{m}$ , silt:  $2\text{--}63\mu\text{m}$  and sand  $63\text{--}2000\mu\text{m}$ ) which sum to 100%. This is an example of compositional data which has constraints for certain statistical analyses (Aitchison, 1986). A large proportion of the spatial variation in psd typically occurs at scales between around 20 and 200 metres (McBratney and Pringle, 1999), so many samples and costly laboratory measurements would be required to map psd accurately using conventional methods. Scientists have shown that remote sensors – ground-based or airborne – can provide effective covariates to aid mapping of soil psd fractions including gamma radiometry (Taylor et al., 2002), geophysical measurements of electrical conductivity (Robinson et al., 2008) and near infra red reflectance spectra (Selige et al., 2006).

The visible and near infra-red (VNIR) diffuse reflectance spectrum (DRS) of a soil sample includes information on the quantities of the mineral phases it contains. It is assumed that each mineral phase or mineral coatings – iron-oxide coatings on silica (Scheidegger et al., 1993) or clay minerals (Carroll, 1958) – occur predominantly in one of the size fractions. The proportion of one size fraction – or all three fractions – in a set of soil samples are estimated from the VNIR spectra using multivariate statistical models by fitting them to laboratory measurements of psd. These multivariate models can then be used to estimate psd for other soils from the local area over which the original samples were collected. Four studies have been published where such models have been successful in predicting the proportions of particles in all three size-fractions (Chang et al., 2001; Shepherd and Walsh, 2002; Cozzolino and Moron, 2003; Sorensen and Dalsgaard, 2005).

By applying laboratory-based VNIR-DRS to a range of samples from the USA,

60 Europe, Africa and Asia, Brown et al. (2006) developed a global model for estimation  
61 of a variety of soil properties including percentage clay fraction using VNIR-DRS. The  
62 authors improved their estimates of percentage clay by including measurements of the  
63 sand size-fraction based on sieving. To date, no published studies have compared the  
64 performance of local and global statistical models for the estimation of all three soil  
65 texture fractions based on VNIR-DRS. Soil scientists need to know whether prediction  
66 accuracies can be substantially improved if statistical models are based on a smaller  
67 subset of soil samples when compared to a regional or global dataset (Sankey et al.,  
68 2008).

69 Remotely-sensed reflectance spectra have also been used to aid mapping of topsoil  
70 texture fractions at fine spatial resolutions (2 to 5 m) at farm scales (Barnes and Baker,  
71 2000) and over small regions using airborne sensors (Selige et al., 2006; Lagacherie et  
72 al., 2008; Gomez et al., 2008). The availability of satellite-based hyperspectral data at  
73 fine spatial (30-m pixel sizes) and spectral (10-nm) resolutions (*e.g.* [www.enmap.org](http://www.enmap.org))  
74 could provide landscape-scale covariates to substantially improve our ability to map  
75 topsoil psd when combined with ground-based measurements, in areas where topsoil  
76 is sufficiently exposed. Soil scientists need to know where in the landscape the rela-  
77 tionships between soil psd and VNIR spectra are likely to be weak or strong – and the  
78 reasons for this – to assess the likely benefits of hyperspectral remote sensing to aid  
79 mapping psd.

80 The prediction of psd fractions from VNIR-DRS is likely to be more accurate if  
81 statistical models are developed and applied to groups of soils with similar mineralogy,  
82 and therefore, VNIR spectra. Of the five soil forming factors, much of the variation  
83 in soil mineralogy – and also VNIR reflectance – is likely to be explained by parent  
84 material (PM) type (Rawlins et al., 2003). This is particularly the case in areas where  
85 Quaternary substrates are the dominant PM type, such as across large parts of north-  
86 ern Europe, where recently formed soils have strong associations with their PM. The  
87 Quaternary parent materials comprise a range of transported materials deposited by

88 glaciers (till), rivers (alluvium) and the wind (aeolian deposits). Where Quaternary  
89 materials are thin or absent, soils develop directly from bedrock. It might be beneficial  
90 to develop statistical models to predict soil texture fractions using VNIR-DRS, based  
91 on their PM type.

92 The aim of this paper is to compare the performance of statistical models for the  
93 prediction of psd from VNIR-DRS for soils with differing types of PM; specifically soils  
94 derived directly from bedrock or from transported PM. We present statistical (partial  
95 least squares regression; PLSR) models used to estimate additive log ratios (Aitchison,  
96 1986) of two texture components (clay:sand and silt:sand) using VNIR-DRS for groups  
97 of soils developed from five PM types in part of agricultural eastern England. We  
98 also establish a single statistical model for all the soils from the five groups. Two  
99 of the groups of PM were sedimentary bedrock; the other three represent a range of  
100 Quaternary (transported) PM types. We compare the VNIR wavelengths which are  
101 significant predictors for the soil texture fractions in the local and global calibration  
102 models and compare their prediction accuracies using independent cross validation  
103 after back transformation to the three size fractions. We seek plausible explanations  
104 to account for the differences in the accuracy of the statistical models in predicting  
105 particle size fraction for soils over the different types of PM.

## 106 **2. Methods**

### 107 *2.1 Study region and soil sampling*

108 The study region is the area of eastern England shown in Figure 1; the spatial dis-  
109 tribution of the soil sampling locations are highlighted. Bedrock in the region ranges  
110 in age from Carboniferous to Cretaceous comprising coals, limestones, sandstone, silt-  
111 stone, mudstone, chalk, marls and ironstones. There are a range of superficial deposits  
112 including glacial till, river and marine alluvium and a large region of lacustrine (lake)  
113 deposits formed by glacial meltwaters which predominantly give rise to Fluvisols. Soils  
114 developed from the two other parent material types are predominantly Cambisols and

115 Gleysols (IUSS Working Group WRB, 2006). Arable agriculture accounts for more  
116 than 90% of land use over the area from which samples were collected.

117 Soil sampling was undertaken across the region at a density of 1 sample per  
118 2 square kilometres in the summers of 1994, 1995 and 1996 as part of a national-  
119 scale geochemical survey (Johnson et al., 2005). Sampling sites were chosen from  
120 alternate kilometre squares of the British National Grid by simple random selection  
121 within each square, subject to the avoidance of roads, tracks, railways, urban land and  
122 other seriously disturbed ground. At each site, surface litter was removed and soil was  
123 sampled from to a depth of 15 cm using five holes at the corners and centre of a square  
124 with a side of length 20 m by a hand auger and combined to form a bulked sample. All  
125 samples of soil were dried and disaggregated. They were sieved to pass 2 mm, coned  
126 and quartered.

127 We selected only those soil sampling sites ( $n=738$ ) over five dominant parent ma-  
128 terial (PM) types (see Figure 1). We did this by assigning to each sampling location a  
129 PM code (based on combinations of solid or superficial geology). These PM codes were  
130 derived from digital versions of the 1:50 000 maps of bedrock geology and superficial  
131 deposits of England, part of DigMap GB of the British Geological Survey (2006). The  
132 number of soil sampling sites in each of the PM groups was as follows. For soils devel-  
133 oped from two different types of bedrock parent material where there was little or no  
134 superficial material above the bedrock (coal-bearing strata; CM  $n=175$  and mudstone-  
135 bearing strata; MDST  $n=47$ ). For soils collected over PM types developed over thick,  
136 superficial deposits: alluvium (both marine and fluvial; ALV  $n=230$ ), glacial till (TILL  
137  $n=186$ ) and lacustrine deposits (LDE  $n=100$ ). The mineralogical composition – based  
138 on X-ray diffraction (XRD) analysis after removal of organic matter – for the different  
139 size fractions of selected soil samples over two of the parent material types (CM and  
140 LDE) are presented in Table 1. It is noteworthy that the soil developed over the la-  
141 custrine deposits has around twice as much kaolinite (33.6%) in the clay size fraction  
142 than the soil over the Coal Measures (16.8%), and that there is chlorite in the clay

143 (13%) and silt (7.1%) fractions of the Coal Measures soil, but this was not detected in  
144 the soil over the lacustrine deposits.

## 145 *2.2 Measurement of diffuse reflectance spectra and redness index*

146 Sub-samples of each soil were scanned in the visible-near infrared region (350–2500 nm)  
147 using an ASD (Analytical Spectral Devices, Boulder, CO) Agri-Spec NIR Spectrometer.  
148 In contrast to the sub-samples which were analysed to determine their psd (see below),  
149 organic matter (OM) was not removed from the sub-samples used for spectral analysis.  
150 The presence of OM – both as particulate carbon and coatings on mineral surfaces –  
151 will influence the VNIR spectra due to the occurrence of organic-related adsorption  
152 features. In some cases, the wavelengths of these adsorption features may coincide  
153 with adsorption features due to minerals in the texture fractions. This would lead to  
154 smaller regression coefficients at these wavelengths in statistical models formed between  
155 the spectra and the texture fractions. In our study, however, by not removing OM  
156 from the soil samples, the main benefit of VNIR-DRS – the rapid and cost-effective  
157 processing of samples – is preserved. In the wider context, remote sensing of soil in  
158 the VNIR region will always include adsorption features of OM in their spectra, so for  
159 its successful application, any interference caused by overlapping adsorption features  
160 must be overcome.

161 A 20-g subsample from each original soil sample was placed in a holder with a  
162 quartz window for scanning. Soils were illuminated and scanned from below using the  
163 spectrometer connected to an ASD muglight with an internal tungsten–quartz–halogen  
164 light source and a 12 mm spot size. Data were collected every 1 nm and every spectrum  
165 was an average of 25 readings. Each sample was scanned twice; the second scan was  
166 made after rotating the sample in its holder through 90° whilst placed on the muglight.  
167 During scanning, a Spectralon 99% reflectance panel was used to optimize and white-  
168 reference the spectrometer after scanning every set of ten samples. We checked that  
169 both sides of the Spectralon panel gave consistent baselines. Before further statistical

170 analysis, we obtained an average of two spectra for each sample.

171 In addition to the spectra, we computed the soil redness index (RI) for each  
172 sample as a potential predictor for particle size fractions. In many soil types, soil  
173 redness is dominated by the occurrence of iron-oxide minerals which form coatings on  
174 clay minerals (Carroll, 1958); so in certain soil types, RI may be strongly correlated  
175 with the proportions of the clay size fraction. The RI was computed as (Mathieu et  
176 al., 1998):

$$RI = \frac{R^2}{(B \times G^3)} \quad (1)$$

177 where R, G, and B represent the reflectance at the wavelength of red, green, and blue  
178 bands (700, 546, and 436 nm, respectively) recorded by the ASD spectrometer.

### 179 *2.3 Particle-size analysis*

180 The protocol for the particle size analysis was recently described in detail by Rawlins et  
181 al. (2009); here we provide a summary of its important features. Organic matter was  
182 removed from all sub-samples prior to psd determination by adding a combination of  
183 hydrogen peroxide and water to each sample and heating the mixture. Calgon solution  
184 was added to the samples to disperse them before analysis by laser granulometry. An  
185 8  $\mu\text{m}$  threshold was used for the upper limit of the clay-sized fraction instead of the  
186 conventional 2  $\mu\text{m}$ ; this corrects for differences in measurements by sedimentation and  
187 laser-based methods for non-spherical particles (Konert and Vandenberghe, 1997). Du-  
188 plicated analyses (n=86) showed that the precision of the method was good; standard  
189 deviations were 2.1% for sand and clay, and 1.2% for silt.

190 The psd for each of the samples in the study (n=738) is shown in Figure 2. There  
191 are substantial differences in the mean sand and clay compositions for each of the five  
192 groups; the mean clay content varies from 23 to 38% and the mean sand content from  
193 21 to 58% (Figure 2 and Table 2). The variation of psd within each of the five groups  
194 (standard deviations shown in Table 2) are quite similar; the TILL and CM group are  
195 somewhat less variable and, as might be expected, the ALV group (alluvial soil parent

196 material) has the most variable psd.

## 197 2.4 Statistical analyses

### 198 2.4.1 Additive log-ratio transformation

199 The compositional constraints on data with distributions that are curtailed at the limits  
200 of 0 and 1 (or 0 and 100%) induces correlations among the components, in this case  
201 the particle size fractions. Linear regression models are not limited in this way, nor are  
202 their predictions constrained to sum to 1.

203 Aitchison (1986) proposed a way out of this difficulty using the additive log ratio  
204 (alr). Suppose we have  $V$  variables, each with values lying between 0 and 1 and  
205 summing to 1, and that we choose  $V - 1$  with values for each unit  $z_1, z_2, \dots, z_{V-1}$ . We  
206 can transform these to

$$q_i = \ln\left(\frac{z_i}{z_V}\right) \quad \text{for all } i = 1, 2, \dots, V - 1, \quad (2)$$

207 where  $z_V$  is the value of the remaining  $V$ th variable. The resulting values over all units  
208 have by definition a logistic normal distribution. This is the additive log ratio (alr)  
209 transform, and it allows us to analyse our compositional data as any other multivariate  
210 normal data.

211 After estimating new values  $\hat{q}_i$ ,  $i = 1, 2, \dots, V - 1$ , we want to return them to  
212 their original scale of composition, and we do so by the inverse transform, the additive  
213 generalized logistic transformation:

$$\begin{aligned} \hat{z}_i &= \frac{\exp(\hat{q}_i)}{1 + \sum_{j=1}^{V-1} \exp(\hat{q}_j)} \quad \text{for all } i = 1, 2, \dots, V - 1 \\ \text{and } \hat{z}_V &= \frac{1}{1 + \sum_{j=1}^{V-1} \exp(\hat{q}_j)}. \end{aligned} \quad (3)$$

214 As Aitchison showed, the results of this back-transformation are the same whichever  
215 variable we select as  $z_V$ .

216 In this study we formed partial least squares regression (PLSR) models for two  
217 alr-transformed variates and back-transformed these values to three size fractions. We  
218 note that the centre of the backtransformed distribution is equivalent to the median

219 on the original distribution (Pawłowsky-Glahn and Olea, 2004) not the mean, and so  
220 the backtransformed values include some bias. We then assessed the accuracy of the  
221 estimates using independent cross-validation.

#### 222 *2.4.2 Partial least squares regression and cross-validation*

223 PLSR is a chemometric technique which is well-suited to multicollinear predictor vari-  
224 ables, such as reflectance measurements in infra red spectroscopy. The predictive re-  
225 gression model can be represented by:

$$Y = b_0 + b_1X_1 + b_kX_k + \epsilon \quad (4)$$

226 where the observed response values ( $Y$ ; in this case the alr ratios) are approximated  
227 by a linear combination of the values of the spectral intensities ( $X$ ), coefficients ( $b$ ) -  
228 referred to as b-coefficients, and an error term ( $\epsilon$ ).

229 To determine the significant wavelengths for prediction of the alr ratios of the  
230 texture fractions, we used both the Variable Importance in the Projection (VIP)  
231 (Chong and Jun, 2005) and the PLS regression coefficients (b-coefficients; Haaland  
232 and Thomas, 1988). For an observed variable  $y$ , the VIP was calculated by:

$$VIP_k(a) = K \sum_a w_{ak}^2 \left( \frac{SSY_a}{SSY_t} \right) \quad (5)$$

233 where  $VIP_k(a)$  gives the importance of the  $k$ th predictor variable based on a model with  
234  $a$  factors,  $w_{ak}$  is the corresponding loading weight of the  $k$ th variable in the  $a$ th PLSR  
235 factor,  $SSY_a$  is the explained sum of squares of  $y$  by a PLSR model with  $a$  factors,  
236  $SSY_t$  is the total sum of squares of  $y$ , and  $K$  is the total number of predictor variables.  
237 The wavelength is considered important if the values of both the b-coefficients and VIP  
238 are sufficiently large. In this study, thresholds for VIP were set to 1 (Chong and Jun,  
239 2005) and the standard deviation of the b-coefficients was applied as their threshold.

240 We used the *pls* package (Mevik and Wehrens, 2007) in the R environment (R  
241 Core Development Team, 2010) to form PLSR models based on the orthogonal scores  
242 algorithm. After taking alr ratios (Equation 2) using the *compositions* package (van  
243 den Boogaart et al., 2008) we fitted models to the two alr ratios for the soils from each  
244 of the five PM groups, and also to all the samples; twelve models in total. We investi-  
245 gated whether spectral pre-processing (first and second derivatives and Savitsky-Golay  
246 smoothing) improved model performance. In each case the untransformed reflectance  
247 data gave the best model performance so we used the original data in fitting all models.  
248 We used a truncated range (450-2450 nm) of wavelengths and the RI as predictors. We  
249 used cross validation to select the optimum number of components from which to form  
250 the models and also calculated the coefficient of determination ( $R^2$ ) to assess model  
251 performance. Prior to forming each model, 10% of the samples were selected randomly  
252 and were not used in model fitting. These samples were then used to assess the model  
253 performance by forming predictions, backtransforming the alr components to propor-  
254 tions of the compositions (Equation 3) and calculating the root-mean-squared-error of  
255 cross validation (RMSE-CV) for between 1 and 12 model components. The RMSE-CV  
256 is calculated as:

$$\text{RMSE-CV} = \sqrt{\frac{1}{n_V} \sum_{i=1}^{n_V} (\hat{z}_i - z_i)^2}, \quad (6)$$

257 where  $z_i$  is the measured proportion of a particle fraction and  $\hat{z}_i$  is its predicted value.  
258 We selected the optimum number of components for each PLSR model based on min-  
259 imisation of the RMSE-CV.

260 To assess the accuracy of the selected models more thoroughly, we undertook  
261 repeated ( $n=100$ ) cross-validation by randomly selecting 10% of the samples from  
262 each group and calculating the mean of the RMSE-CV and the mean of the residual  
263 prediction deviation (RPD); the ratio of the standard deviation of the validation sample  
264 set and the standard error of prediction (Equation 6). This statistic provides a useful

265 indication of the quality of the model because it accounts for the variation in the size  
266 fractions in the validation dataset.

### 267 **3. Results and their interpretation**

#### 268 *3.1 Regression models and psd estimation accuracy*

269 Summary data for the PLSR models fitted to the two alr size-fraction ratios for each  
270 of the soils grouped by PM and all soils grouped together are presented in Table 3.  
271 The position of those wavelengths in the PLSR models identified as having significant  
272 predictive power – based on large b-coefficients and VIP scores – are presented in Fig-  
273 ure 3. For both size-fraction ratios, the global group and alluvial sediment group have  
274 the largest number ( $n=10$  or  $11$ ) of orthogonal model components probably because  
275 they represent a greater diversity of soil types than the other individual soil groups.  
276 It is notable that RI was only a significant predictor for the CM soil group – this  
277 may partly be related to the colour associated with the range and age of iron-bearing  
278 mineral phases in the coal-bearing strata. The mineral pyrite is abundant (Spears et  
279 al., 1999) in the bedrock from which these soils formed and over time this weathers to  
280 form a range of iron-bearing minerals (hematite and goethite; also present in the host  
281 rock) of varying age. Soil samples containing differing proportions of these minerals  
282 will have quite different redness features which may account for its significance as a  
283 predictor in this soil group. This may in part be due to different ageing of iron oxyhy-  
284 droxides giving rise to differences in colour wavelengths (yellow, orange and red). For  
285 example, goethite reddens as it ages to hematite whilst lepidocrocite is dark brown or  
286 black and amorphous iron-oxide is between yellow and orange in colour. We cannot  
287 explain the significant wavelengths between 950 and 1050nm for the TILL and ALV  
288 models as these do not appear to relate to known absorption features in the near infra  
289 red spectrum.

290 The significant predictive wavelengths are consistent with colour in the visible  
291 light range (450-700nm). The H<sub>2</sub>O adsorption bands at 750, 975, 1900–1950 and

292 2200nm are present as significant predictors in many of the PLSR models presumably  
293 associated with certain water-absorbing clay mineral phases. The H<sub>2</sub>O adsorption band  
294 at 1400 nm is absent from all but the global clay:silt ratio model; this demonstrates  
295 the potential problem in forming global calibration models which may be based on  
296 predictive wavelengths that would not be justified based on a groups of local models.  
297 Significant wavelengths for other adsorption bands associated with certain clay mineral  
298 phases which are common in British soils include 2204–2211 nm (illite, kaolinite),  
299 2340 nm (illite) and 2207 nm (smectite). Adsorption bands commonly associated with  
300 hematite (920 nm) and goethite (880 nm) are absent from all the models, although  
301 their effects through soil colour may be more significant in the visible wavelength  
302 range. There is substantially greater ( $n=58$  wavelengths) overlap in the significant  
303 predictive wavelengths for the clay:silt size fraction ratio models (Figure 3a) for the  
304 five individual soil groups compared with the silt:sand ratio size fraction groups ( $n=7$   
305 wavelengths; Figure 3b). The overlapping wavelengths in the clay:silt size fraction  
306 models are dominated by wavelengths centred around the water absorption feature  
307 at 1900-1950nm (illite) and smectite (2004-2211nm); the absorption feature which is  
308 common to all models between 2407 and 2420nm may be related to adsorption features  
309 associated with carbonate minerals; chalk bedrock underlies the TILL soils in the north-  
310 east of the study region (Figure 1).

311 The results of repeated ( $n=100$ ) 10% cross-validation for each of the models  
312 applied to the soil groups are presented in Table 2; the RMSE-CV and RPDs were  
313 calculated after back-transformation to the original three size-fractions. Overall model  
314 performance – based on RPD – is poorest for the silt size-fraction. This may be because  
315 this fraction shares a boundary with the two other size-fractions, whilst they each share  
316 only one. The intermediate, silt size-fraction is likely to comprise a larger proportion of  
317 uncommon minerals than the two other size-fractions. In terms of overall performance,  
318 the RPD values are larger (more accurate estimates) for the models relating to soils  
319 developed over the transported PM compared to those formed from bedrock. For the

320 clay size-fraction, the RPDs decline in the order: LDE > ALV > TILL > MDST >  
321 CM. With the exception of the last two groups which swap places in the order stated  
322 above, the same pattern also applies to the sand size-fraction. It is noteworthy that  
323 even though the LDE group has the largest mean sand-size content (58%; Table 2)  
324 it has the best overall model performance for the sand size-fraction. If, as is often  
325 the case, the sand fraction is dominated by quartz which has no absorption features  
326 in the NIR range (350-2500 nm), we might have expected model performance to be  
327 poor relative to the other groups. In the case of the silt size-fraction, soils over the  
328 transported PM again have the largest RPD values compared to those derived from  
329 bedrock: ALV > TILL > LDE > CM > MDST.

330 In each size-fraction, the RPD for the PLSR models developed for all soils (global)  
331 generally has an intermediate value; greater than the soils over bedrock but less than  
332 those over transported PM types. In the case of estimating clay and sand size-fractions  
333 for soil over the transported PM types, if we rely on a global PLSR model the average  
334 error of our predictions would be substantially larger than if we had developed models  
335 for each PM group (Table 2). In the case of the silt size-fraction, only the overall  
336 model performance (RPD= 1.88) for the alluvial soils is greater than that of the global  
337 model (RPD= 1.76), with particularly poor overall performance for the MDST and  
338 CM models.

#### 339 **4. Discussion**

340 Previous research has highlighted the importance of PM when estimating cation ex-  
341 change capacity of soil using VNIR-DRS across another part of eastern England (Sav-  
342 vides et al., 2010). Our analysis has shown that there are substantial differences in  
343 the performance of statistical models for prediction of particle size fractions based on  
344 VNIR-DRS for soils developed over different PM types across a large area of Eastern  
345 England. This highlights the importance of existing maps of soil PM or soil type to  
346 enhance the application of sensor-based covariates for producing digital soil maps.

347 We currently have no direct evidence to account for the observed differences in

348 VNIR-DRS model performance for the different soil groups. Here we give further  
349 consideration to three possible reasons. The first is that soils derived from transported  
350 (allochthonous) PM types are likely to be more mineralogically homogeneous than  
351 those derived from *in-situ* weathering of bedrock because the former have undergone  
352 sorting processes associated with transport and redeposition, whilst the latter have only  
353 been subject to weathering *in situ*. For example, we know from direct observation and  
354 mineralogical analyses that there is considerable variation in the lithologies (mudstones,  
355 siltstone, sandstone, coal-bearing strata) which comprise the coal-bearing strata (CM)  
356 in our study area.

357 A second explanation to account for the observed differences in model perfor-  
358 mance for the different soil groups concerns the relative abundance of iron-oxide coat-  
359 ings of soil clays which might lead to bias in the reflectance spectra; more or different  
360 types of iron-oxide coatings could diminish the clay mineral signatures and hinder  
361 VNIR-DRS model performance. We assume that total soil iron content – for which  
362 we have measurements for each soil sample from XRF analysis (see Rawlins et al.,  
363 2009) is strongly correlated with iron-oxide concentrations. Soils developed over the  
364 coal-bearing strata contained the greatest median concentration of total iron (4.71%)  
365 and had the poorest overall model performance for size fraction prediction. However,  
366 the alluvial soils group was the next most enriched in total iron (median=3.91%) and  
367 VNIR-DRS model performance was reasonable, which confounds this theory if we as-  
368 sume that iron-oxides have similar associations with soil minerals in each of the soil  
369 groups (i.e. the proportion of iron-oxide coated clay mineral particles has a linear  
370 relationship with total iron content). The relationship between iron oxide coatings of  
371 minerals and VNIR-DRS warrants further investigation.

372 Thirdly, the presence of widely differing proportions of quartz – which has no  
373 spectral adsorption features in the VNIR range (350-2500 nm; Ferraro, 1982) – in the  
374 different size fractions would weaken the statistical relationships with VNIR-DR spec-  
375 tra. Although this limitation of size fraction estimation from VNIR-DRS is recognised,

376 it has not been widely referred to in the soil science literature. Our analyses show  
377 that soils with greater, average sand-size (quartz-dominated) fractions do not in gen-  
378 eral have weaker correlations between VNIR-DR spectra and each texture fraction; the  
379 group with the largest mean sand size fraction had amongst the best (largest) RPD  
380 values. Further research is required to provide an unequivocal, evidence-based expla-  
381 nation for the observed differences in the strength of relationships between psd and  
382 VNIR-DRS for soils developed over transported and bedrock-derived PM types.

383 Our results suggest that accuracy in mapping soil psd based on VNIR-DRS is  
384 likely to be substantially improved if local statistical models are developed compared  
385 to regional or global approaches. This is also likely to apply to the use of hyperspectral,  
386 remotely sensed data to map size fractions (Lagacherie et al., 2008) because the strength  
387 of the statistical relationships between the size fractions and spectral signatures will  
388 be similar.

389 Our previous research showed that the total concentration of five elements in the  
390 soil (Al, Fe, Ni, Ti and Zr) could be used to accurately estimate soil psd across the same  
391 study region with substantially smaller RMSE-CV compared to the local VNIR-DRS  
392 models (clay 4.9% versus a range from 6.1 to 8.8 %, respectively; sand 8.8% versus  
393 a range from 9.7 to 13.1% respectively). Although these estimates are more accurate  
394 the cost associated with acquisition of geochemical data makes it far more costly when  
395 compared to the spectral approach. The latter also has the advantage for the potential  
396 application of exhaustive, remotely sensed data to improve estimates at fine scales  
397 (< 10 m).

## 398 **5. Conclusions**

399 Our analyses have shown that there are substantially stronger relationships between  
400 psd and VNIR-DR spectra for topsoils developed from transported parent materials  
401 (three groups) compared to those developed directly from bedrock (two groups) at the  
402 regional scale. Based on RPD values from repeated cross-validation, the statistical  
403 models developed between the additive log ratios of the psd fractions and VNIR-DR

404 spectra for topsoils over the transported parent materials generally perform better  
405 than a global model developed for all five soil groups soils across the region. This  
406 has important implications for optimal strategies for mapping psd using field-based  
407 VNIR-DRS and the use of remotely sensed, hyperspectral data.

## 408 **Acknowledgements**

409 This paper is published with the permission of the Executive Director of the British  
410 Geological Survey (Natural Environment Research Council). We would like to thank:  
411 (i) staff from the British Geological Survey involved in collecting and preparing the  
412 soil samples, (ii) Andy Tye, Oliver Gould and Yinka Oshokoya who organised and  
413 undertook the spectral and psd analyses, and (iii) Doris Wagner and Ian Mounteney  
414 for the XRD analyses. We also thank an anonymous reviewer for helpful comments on  
415 the original manuscript.

## 416 **References**

- 417 Aitchison, J. 1986 *The Statistical Analysis of Compositional Data*. Chapman and  
418 Hall, London.
- 419 Barnes, E. M., Baker, M. G., 2000. Multispectral data for mapping soil texture:  
420 possibilities and limitations. *Applied Engineering in Agriculture* 16, 731–741.
- 421 British Geological Survey 2006. *Digital Geological Map of Great Britain 1:50 000*  
422 *scale (DiGMapGB-50) data [CD-ROM] Version 3.14*. British Geological Survey,  
423 Keyworth, Nottingham.
- 424 Brown, D.J., Shepherd, K.D., Walsh, M.G., Dewayne Mays, M., Reinsch, T.G., 2006.  
425 Global soil characterization with VNIR diffuse reflectance spectroscopy. *Geo-*  
426 *derma* 132, 273–290.
- 427 Carroll, D., 1958. Role of clay minerals in the transportation of iron. *Geochimica Et*  
428 *Cosmochimica Acta*, 14, 1–16.

- 429 Chang, C.-W., Laird, D.A., Mausbach, M.J. and Hurburgh, C.R., Jr., 2001. Near-  
430 Infrared Reflectance Spectroscopy-Principal Components Regression Analyses of  
431 Soil Properties. *Soil Science Society of America Journal* 65, 480–490.
- 432 Chong, I.-G., Jun, C.-H., 2005. Performance of some variable selection methods when  
433 multicollinearity is present. *Chemometrics and Intelligent Laboratory Systems*  
434 78, 103–112.
- 435 Cozzolino, D., Moron, A., 2003. The potential of near-infrared reflectance spec-  
436 troscopy to analyse soil chemical and physical characteristics. *Journal of Agri-  
437 cultural Science* 140, 65–71.
- 438 Ferraro, J.R., 1982. *The Sadtler infrared spectra handbook of minerals and clays.*  
439 Heyden & Son Ltd. London. 440 pp.
- 440 Gomez, C., Lagacherie, P., Coulouma, G., 2008. Continuum removal versus PLSR  
441 method for clay and calcium carbonate content estimation from laboratory and  
442 airborne hyperspectral measurements. *Geoderma* 148, 141–148.
- 443 Haaland, D.M., Thomas, E.V., 1988. Partial least-squares methods for spectral anal-  
444 yses.1. Relation to other quantitative calibration methods and the extraction of  
445 qualitative information. *Analytical Chemistry* 60, 1193–1202.
- 446 Hodgson, J.M., 1974. *Soil Survey Field Handbook Describing and Sampling Soil*  
447 *Profiles, Soil Survey of England Wales, Harpenden, Hertfordshire (1974).*
- 448 IUSS Working Group WRB, 2006. *World reference base for soil resources 2006. 2nd*  
449 *edition.* World Soil Resources Reports No. 103. FAO, Rome.
- 450 Johnson, C.C., Breward, N., Ander, E.L., Ault, L., 2005. G-BASE: Baseline geochem-  
451 ical mapping of Great Britain and Northern Ireland. *Geochemistry: Exploration-  
452 Environment-Analysis*, 5, 1–13.

- 453 Konert, M., Vandenberghe, J., 1997. Comparison of laser grain size analysis with  
454 pipette and sieve analysis: a solution for the underestimation of the clay fraction.  
455 *Sedimentology* 44, 523–535.
- 456 Lagacherie, P., Baret, F., Feret, J.B., Netto, J.M. and Robbez-Masson, J.M., 2008.  
457 Estimation of soil clay and calcium carbonate using laboratory, field and airborne  
458 hyperspectral measurements. *Remote Sensing Of Environment* 112, 825–835.
- 459 Mathieu, R., Pouget, M., Cervelle, B., Escadafal, R., 1998. Relationships between  
460 satellite-based radiometric indices simulated using laboratory reflectance data  
461 and typic soil color of an arid environment. *Remote Sensing of Environment* 66,  
462 17–28.
- 463 McBratney, A.B. and Pringle, M.J., 1999. Estimating average and proportional vari-  
464 ograms of soil properties and their potential use in precision agriculture. *Precision*  
465 *Agriculture* 1, 219–236.
- 466 Mevik, B.H. and Wehrens, R., 2007. The pls Package: Principal Component and  
467 Partial Least Squares Regression in R. *Journal of Statistical Software* 18, 1–24.
- 468 Pawlowsky-Glahn, V., Olea, R. A., 2004. *Geostatistical Analysis of Compositional*  
469 *Data*. Oxford University Press, Oxford.
- 470 R Development Core Team, 2010. *R: A Language and Environment for Statis-*  
471 *tical Computing.*, R Foundation for Statistical Computing, Vienna, Austria,  
472 <http://www.R-project.org>.
- 473 Rawlins, B.G., Webster, R., Lawley, R., Tye, A.M., O’Hara, S.O., 2009 Estimating  
474 particle-size fractions of soil dominated by silicate minerals from geochemistry.  
475 *European Journal of Soil Science* 60, 116–126.
- 476 Rawlins, B.G., Webster, R., Lister, T.R., 2003. The influence of parent material on  
477 top soil geochemistry in eastern England. *Earth Surface Processes and Landforms*  
478 28, 1389–1409.

- 479 Robinson, D.A., Abdu, H., Jones, S.B., Seyfried, M., Lebron, I., Knight, R., 2008.  
480 Eco-Geophysical Imaging of Watershed-Scale Soil Patterns Links with Plant  
481 Community Spatial Patterns. *Vadose Zone Journal* 7, 1132–1138.
- 482 Sankey, J. B., Brown, D. J., Bernard, M. L., Lawrence, R. L., 2008. Comparing local  
483 vs. global visible and near-infrared (VisNIR) diffuse reflectance spectroscopy  
484 (DRS) calibrations for the prediction of soil clay, organic C and inorganic C.  
485 *Geoderma*, 148, 149–158.
- 486 Savvides, A. , Corstanje, R. , Baxter, S.J. , Rawlins, B.G. and Lark, R.M., 2010.  
487 The Relationship between Diffuse Spectral Reflectance of the Soil and Its Cation  
488 Exchange Capacity Is Scale-Dependent, *Geoderma* 154, 353–358.
- 489 Scheidegger, A., Borkovec, M., Sticher, H., 1993. Coating of silica sand with goethite:  
490 preparation and analytical identification. *Geoderma*, 58, 43–65.
- 491 Selige, T., Bohner, J., Schmidhalter, U., 2006. High resolution topsoil mapping using  
492 hyperspectral image and field data in multivariate regression modeling proce-  
493 dures. *Geoderma* 136, 235–244.
- 494 Shepherd, K.D., Walsh, M.G., 2002. Development of reflectance spectral libraries for  
495 characterization of soil properties. *Soil Science Society of America Journal* 66,  
496 988–998.
- 497 Sorensen, L. K., Dalsgaard, S., 2005. Determination of clay and other soil properties  
498 by near infrared spectroscopy. *Soil Science Society of America Journal* 69, 159–  
499 167.
- 500 Spears, D.A., Manzanares-Papayanopoulos, L. I., Booth, C.A., 1999. The distribution  
501 and origin of trace elements in a UK coal; the importance of pyrite. *Fuel* 78,  
502 1671–1677.
- 503 Taylor, M.J., Smettem, K., Pracilio, G., Verboom, W., 2002. Relationships between

504 soil properties and high-resolution radiometrics, central eastern Wheatbelt, West-  
505 ern Australia. *Exploration Geophysics* 33, 95–102.

506 van den Boogaart, K. G., Tolosana, R., Bren, M., 2008. *compositions: Compositional*  
507 *Data Analysis*. R package version 1.01-1. <http://www.stat.boogaart.de/compositions>.

508 **Figure captions**

509 **Figure 1** Parent material and soil sampling locations across the study region.

510 **Figure 2** Particle-size distribution for groups of soils developed from five different  
511 parent material types across the study region. Figure 1 shows their spatial dis-  
512 tribution. The partitions of the triangle and class names are those in the Field  
513 Handbook of the Soil Survey of England and Wales compiled by Hodgson (1974).

514 **Figure 3** Wavelengths at which both variable importance in the projection (VIP)  
515 scores and regression (beta) coefficients are significant in partial least squares  
516 models of reflectance spectra for prediction of additive log ratios of texture frac-  
517 tions for soils grouped by parent material and for all soils: a) clay:silt, b) silt:sand.  
518 Bands which are significant in all five soil model groups are shown in grey.

519 Table 1 X-ray diffraction analysis of estimated mineralogic composition for three size  
 520 fraction separates for single selected topsoil samples developed over coal-bearing strata  
 521 (CM) bedrock and lacustrine deposits (LDE).

Size fraction Parent material	Sand		Silt		Clay	
	CM	LDE	CM	LDE	CM	LDE
proportion of total mass (%)	8	58	23	26	69	15
albite	na	3.2	3.9	5.4	nd	<0.5
522 anatase	na	nd	0.5	<0.5	0.7	<0.5
*kaolin	na	nd	7.5	9.1	16.8	33.6
K-feldspar	na	7.2	3.4	5.6	nd	<0.5
†mica	na	nd	27.5	13.7	58.2	57.7
quartz	na	89.6	49.9	66.1	11.3	8.2
chlorite	na	nd	7.1	nd	13	nd

523 \* kaolin: undifferentiated kaolin group minerals possibly including kaolinite, halloysite

524 †mica: undifferentiated mica species, possibly including muscovite, biotite, illite and

525 illite/smectite

526 nd = not detected

527 na = not analysed

528

529 Table 2 Mean (%) and standard deviation (St.Dev) of three size fractions for groups of five soil samples and all samples (Global).  
530 Root mean square error of cross-validation (RMSE-CV: %) based on  $n=100$  repeated random selection of 10% of samples based on  
531 PLSR models fitted to each group and mean residual prediction deviation (RPD).

	Clay			Silt			Sand					
	Mean	St. Dev.	RMSE-CV	RPD	Mean	St. Dev.	RMSE-CV	RPD	Mean	St. Dev.	RMSE-CV	RPD
Global	33	14	7.9	1.75	31	13	7.6	1.76	36	23	13.1	1.74
ALV	38	15	7.1	2.01	41	15	8.0	1.88	21	22	11.9	1.87
LDE	23	14	6.1	2.33	18	11	7.0	1.4	58	23	11.7	1.94
TILL	26	11	5.7	1.82	22	8	5.3	1.62	52	17	9.7	1.71
CM	38	11	8.8	1.33	32	7	5.8	1.25	29	16	11.7	1.39
MDST	33	13	7.0	1.68	34	10	7.8	1.15	32	17	11.5	1.34

533 Table 3 Features of partial least squares regression (PLSR) models formed between two  
 534 alr ratio size fractions and VNIR-DRS for five sets of soil grouped by parent material  
 535 type: a) number ( $n$ ) of orthogonal PLSR components, b) coefficient of determination  
 536 (adjusted  $R^2$ ), c) number ( $n$ ) of wavelengths ( $\lambda$ ) in PLSR model which have VIP scores  
 537 and beta coefficients greater than significance thresholds (see text).

538

Group ( $n$ )	alr clay:silt			alr silt:sand		
	$n$ components	$R^2$	$n(\lambda)$	$n$ components	$R^2$	$n(\lambda)$
Global (738)	11	0.64	523	11	0.60	217
ALV (230)	10	0.80	413	10	0.75	405
LDE (100)	8	0.86	411	8	0.77	225
TILL (186)	8	0.69	432	8	0.67	534
*CM (175)	7	0.50	304	7	0.51	245
MDST (47)	5	0.59	379	5	0.60	441

540 \* includes redness index (RI) as a significant predictor – without the RI the maximum  
 541  $R^2$  values of PLSR models between VNIR-DRS and the clay:silt and silt:sand fractions  
 542 of the CM soils were 0.33 and 0.38, respectively.

Figure 1:

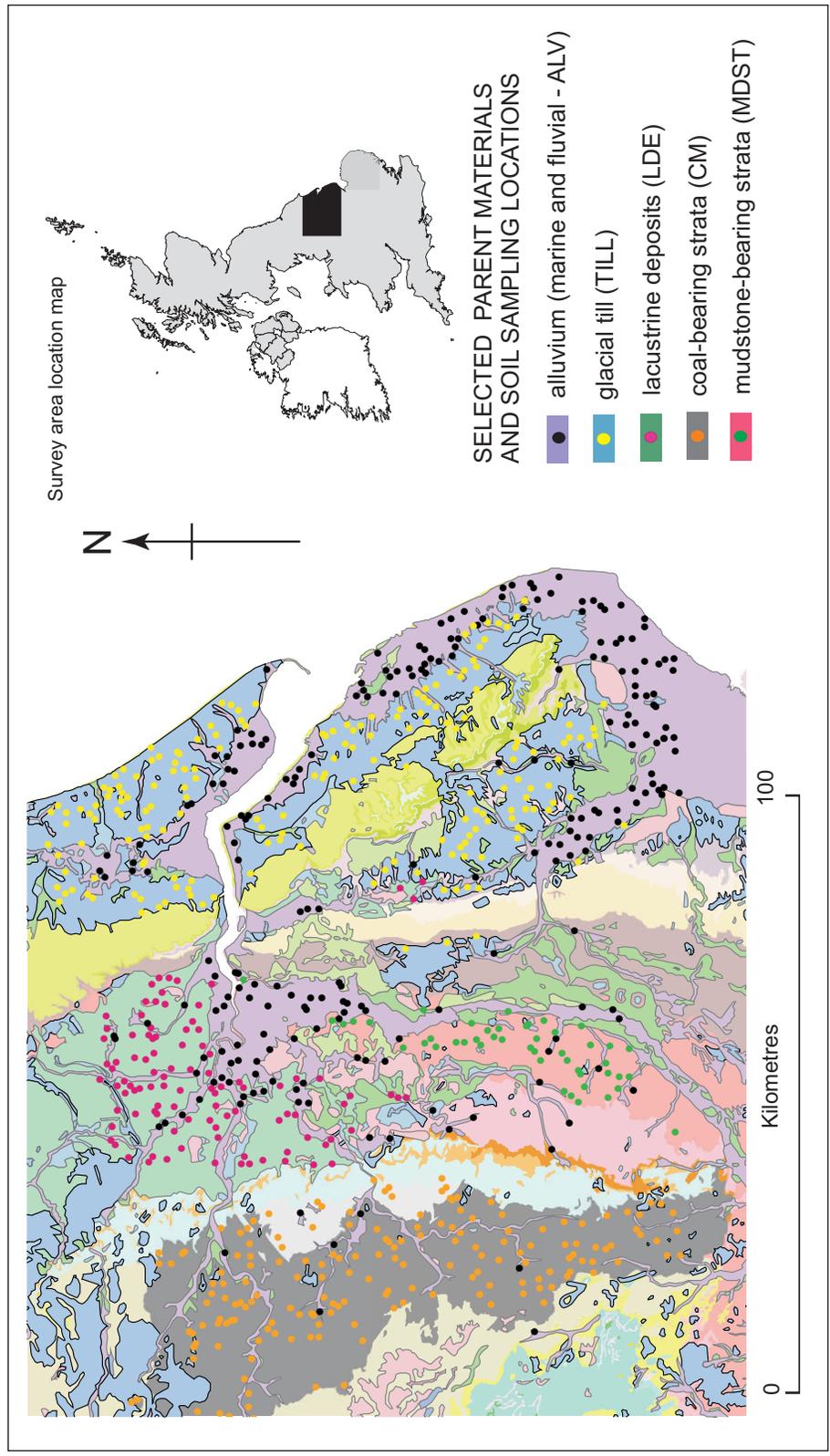


Figure 2:

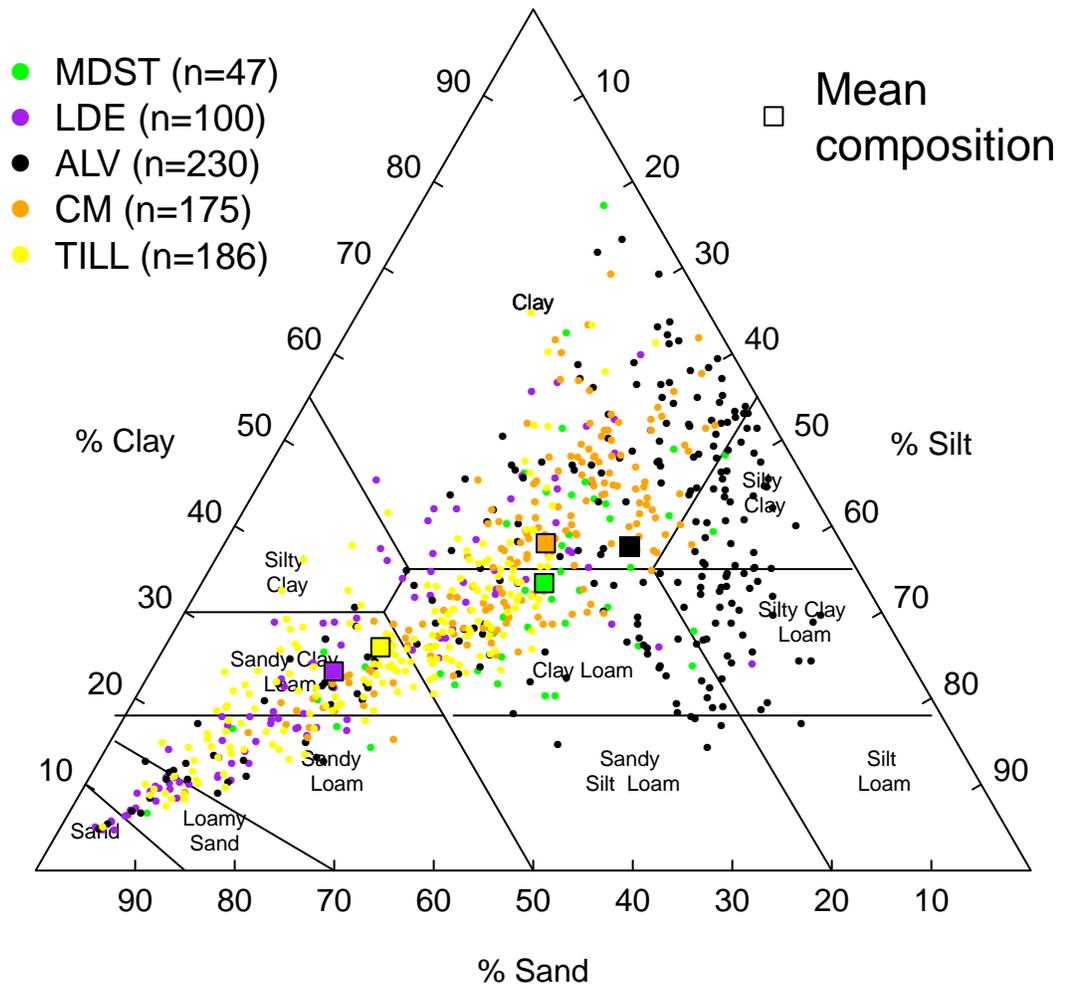


Figure 3:

