

Effect of soil horizon stratigraphy on the microbial ecology of alpine paleosols

Young, J. M., Skvortsov, T., Kelleher, B. P., Mahaney, W. C., Somelar, P., & Allen, C. C. R. (2019). Effect of soil horizon stratigraphy on the microbial ecology of alpine paleosols. *Science of the Total Environment*, *657*, 1183-1193. https://doi.org/10.1016/j.scitotenv.2018.11.442

Published in:

Science of the Total Environment

Document Version: Peer reviewed version

Queen's University Belfast - Research Portal: Link to publication record in Queen's University Belfast Research Portal

Publisher rights

© 2018 Elsevier B.V. This manuscript version is made available under the CC-BY-NC-ND 4.0 license http://creativecommons.org/licenses/by-nc-nd/4.0/,which permits distribution and reproduction for noncommercial purposes, provided the author and source are cited

General rights

Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Open Access

This research has been made openly available by Queen's academics and its Open Research team. We would love to hear how access to this research benefits you. – Share your feedback with us: http://go.qub.ac.uk/oa-feedback

Effect of soil horizon stratigraphy on the microbial ecology of alpine paleosols

4

Jonathan M. Young^{1*}, Timofey Skvortsov², Brian P. Kelleher³, William C. Mahaney⁴, Peeter Somelar⁵,
 Christopher C.R. Allen¹

7 ¹School of Biological Sciences, Queen's University Belfast, UK.

8 ²School of Pharmacy, Queen's University Belfast, UK.

9 ³School of Chemical Sciences, Dublin City University, Ireland.

10 ⁴ Quaternary Surveys, 26 Thornhill Ave, Thornhill, Ontario, Canada

- 11 ⁵ Department of Geology, Tartu University, Tartu, Estonia
- 12

13 *Author for correspondence.

14 **1 ABSTRACT**

15 There is remarkable potential for research at the interface between the earth sciences and 16 environmental microbiology that may lead to advances in our understanding of the role of bacterial 17 communities in the surface or subsurface environment of our planet. One mainstay of sedimentary 18 classification is the concept of differential soil and/or paleosol horizons being the result of primarily physical and chemical weathering, with relatively little understanding of how microbial communities 19 20 between these stratified horizons differ, if at all. In this study we evaluate the differences in 21 microbial community taxonomy and biogeochemical functional potential between stratified soil 22 horizons in an alpine paleosol environment using next-generation sequencing (NGS) shotgun 23 sequencing. Paleosols represent a unique environment to study the effect of differences soil horizon 24 environments on the microbial community due to their relative isolation, and the fact that three 25 distinct stratified soil horizons can be identified within the top 30cm of the soil. This enables us to 26 assess variation in microbial community composition that will be relatively distinct from variation

27 due to distance alone. We test the hypothesis that variation in soil community composition is linked 28 to variation in the physical and chemical parameters that define stratigraphy. Multivariate statistical 29 analysis of sequencing reads from soil horizons across five sampling sites revealed that 1223 30 microbial genera vary significantly and consistently in abundance across stratified soil horizons at 31 class level. Specifically Ktedonobacter, Bacilli and Betaproteobacteria responded most strongly to 32 soil depth. Alpha diversity showed a positive correlation with soil depth. Beta diversity, however, did not differ significantly between horizons. Genes involved in carbohydrate and nitrogen metabolism 33 34 were found to be more abundant in Ah horizon samples. Closer inspection of carbohydrate 35 metabolism genes revealed that genes involved in CO₂ fixation, fermentation and saccharide 36 metabolism decreased in abundance with depth while one-carbon metabolism increased down 37 profile.

38 **2** INTRODUCTION

39

40 Soil microbiology is a rapidly growing discipline, driven by the advances in high throughput 41 sequencing and bioinformatic techniques. Microbiological studies on alpine soils to date 42 have focused mainly on alpine forest, meadow and grassland soils (Ding et al., 2015; Yashiro 43 et al., 2016; Zhang et al., 2013), comparatively little work has been carried out on cold, high 44 altitude alpine soils and paleosols above the timberline (termed the alpine zone; Mahaney et al., 2016). Studies which have focused on soils/paleosols within the alpine zone, such as 45 glacier forefields, have utilised techniques such as amplicon sequencing of the 16S rRNA 46 gene, restriction fragment length polymorphism (RFLP) and/or GeoChip microarrays (He et 47 48 al., 2007) to assess microbial populations (Lazzaro et al., 2015; Zhang et al., 2013). These

49 techniques are subject to significant limitations and potential biases, that can now be
50 addressed using shotgun sequencing approaches coupled with bioinformatics.

51 Several studies have shown that the number of high alpine ecosystems may be expanding 52 worldwide, due to an increase in glacier and ice cap thawing rates (King et al., 2011; Byers, 53 2008; Zemp et al., 2006)), although there is a chance that increased temperature may 54 reverse this trend by expanding forest to higher elevations (Körner and Paulsen 2004). Regardless, concerns are rising regarding the potential implications this may have for global 55 56 atmospheric carbon levels as soils are a major carbon sink: soil organic matter contains 57 approximately 3 times the size of the atmospheric pool of carbon and 4.5 times that of the biotic pool (Lal, 2004). Much more organic carbon exists in soils than in vegetation and the 58 59 atmosphere combined, and the global carbon budget can be strongly influenced by changes 60 in soil carbon content (Fierer et al, 2009; Serna-Chavez et al, 2013). Thawing permafrost in soils newly exposed to increased temperature may increase the rate at which soil microbes 61 62 degrade soil organic carbon thus leading to a net increase in CO₂ emissions (Davidson and Janssens, 2006). Additionally, the rate of methanogenesis could potentially increase, due to 63 an increase in available metabolites for methanogens (such as acetate), generated as by-64 products from organisms carrying out soil organic carbon degradation, which could lead to a 65 66 net increase in methane emissions. Though these sediments are not completely anaerobic, soils in general become more anoxic with increasing depth (Yu et al., 2006). Additionally 67 68 even in well oxygenated soils, methanogenesis can still occur due to the presence of anoxic 69 microenvironments within soil aggregates (Fierer 2017). Without a more thorough understanding of the functional potential of the microbial communities present in these 70 soils we cannot make predictions as to the possible effects of an increase in exposed alpine 71 72 soils.

The alpine paleosols studied in this report have developed from the weathering of parent 73 74 material deposited during two periods following the last glaciation: The Bølling – Allerød warming of ~15-12.8 ka and the Younger Dryas (YD) from 12.8-~11.5 ka (Mahaney et al., 75 2017). Moraines (deposits of glacial debris described here) from both pre-YD and YD time 76 77 are evident at elevations greater than 2400m above sea level in this region of France/Italy. The paleosols in these deposits are classified as Cryochrepts (i.e. cold-climate pedons) with 78 Ah/Bw/Cox profiles over fresh undifferentiated, unweathered substrates. In geological 79 80 terms the soils are identified as paleosols, and they possess distinct stratigraphy with 81 'classically' identified Ah/Bw/Cox soil horizons (Mahaney et al., 2013). The soils in this study 82 are generally acidic with pH ranging from 4.0-6.0 (supplementary data file 1). 83 How microbial communities vary with respect to these soil horizons still remains an 84 underexplored avenue of research in soil-paleosol microbial ecology. Of very few relevant published studies, those which compare soil-paleosol horizons either employed 16S rRNA 85 86 gene amplicon sequencing (Baldrian et al., 2012; Mahaney et al. 2016) or focused 87 specifically on certain functional enzymes, such as dehalogenation enzymes (Weigold et al., 2016). No studies have focused on changes in overall microbial community structure in 88 relation to paleosol stratigraphy. 89 90 The alpine paleosols studied here present a unique opportunity to address this apparent 91 gap in knowledge due to: i. The fact that they possess well stratified horizons, being distinguishable horizons in 92 93 physical terms. ii. 94 The wealth of geochemical data already available for the region (e.g. Mahaney et al.,

95 2016; Mahaney and Keiser, 2013).

96 iii. Their isolation from the Anthropocene: as paleosols these soils are relatively
97 undisturbed by human activity until now, allowing us to assess the effect of >10,000 years of
98 isolated soil development. Their evolutionary history initiated in the Late Pleistocene
99 coming up to their present state in the Late Holocene entitles them to the 'paleo'
100 designation.

There are a number of soil classification systems in current use (Eswaran, 2003), these varying by country and utility (e.g. assessing the suitability of the soil for agricultural use, or inferring its geological origin). In this study soils have been classified according to the United States Department of Agriculture (USDA) soil taxonomy system (USDA, 1999). In this system, a distinct soil horizon is classified as a layer of minerals and organic material, which differs from the parent material in mineral, physical, chemical, morphological, and biological features.

In this study, using high throughput shotgun sequencing, we first aimed to test the
hypothesis that microbial communities vary significantly and consistently between stratified
soil horizons in an alpine paleosol environment, both in terms of their taxonomic profile and
functional potential. We then test the hypothesis that variations in the soil microbiome can
be linked to variations in the physical and chemical properties of the soil horizons. Finally,
we investigated which taxa and functional gene categories showed the strongest changes in
abundance *between* soil horizons.

115 **<u>3 Methods</u>**

Sites were selected on the basis of air photo interpretation (1:20,000 scale) of deposits.
Clasts embedded in major landform surfaces and profiles (Fig. 1) were sampled following

excavation of sections to depths of ~0.8 m. Profiles extend to depths of ~40-50 cm ± 5 cm. 118 Paleosol descriptions follow standard nomenclature (NSSC, 1995) and Birkeland (1999). The 119 'Cox' horizon designation, originally defined by Birkeland (1999), is applied to strata with 120 detectable levels of secondary Fe hydroxides and oxides, whereas 'Cu' refers to 121 unweathered parent material (Hodgson, 1976). The 'Ah' horizon designation is applied 122 when surface color is stronger than 10YR 3/1, an indication of appreciable organic carbon 123 124 accumulation (Canada Soil Survey Comm., 1998). Soil colors were assigned using Oyama and 125 Takehara's (1970) soil colour chips. Bulk samples (250-300 g) were collected from paleosol 126 horizons for particle size, clay mineral, geochemical and microbiological analyses. DNA was extracted from alpine paleosol samples using a PowerSoil DNA extraction kit (Mo 127 Bio). For each DNA extraction 0.25g of sample was used per extraction, extractions were 128 129 performed in triplicate and pooled. Extractions were performed according to the manufacturer's protocol with the following modifications: samples were homogenised using 130 131 a FastPrep 120 cell disrupter system (Thermo-Fisher) at 5.5 m.s⁻¹ for 2 minutes, rather than a standard bench top vortex. The eluted and pooled DNA was further purified via two 132 rounds of ethanol precipitation; DNA solution was suspended in 3 volumes of ice cold 100% 133 ethanol, 0.1 volumes 5 M sodium acetate solution (pH 5.2) and 2 µL of linearized 134 135 polyacrylamide (LPA), the solution was then incubated overnight at -20°C and centrifuged at 18000 x g at 4°C for 30 minutes. The supernatant was discarded and the pellet washed in 136 137 70% ice cold ethanol and again centrifuged at 18000 x g at 4°C for 5 minutes. Finally, the supernatant was discarded and the pellet allowed to air dry for 15 minutes before 138 resuspension in 50 µL molecular grade H₂0. Final DNA concentrations were measured using 139 a Quantus Fluorometer (Promega) in conjunction with the Quantiflour DsDNA dye system 140

(Promega). Soil physical and chemical variables were measured as previously described 141 (Mahaney et al., 2016) 142

143 Sequencing read quality, length and adapter contamination were initially assessed using 144 Fastqc (Andrews, 2010). Adapter trimming was performed using bbduk from the bbtools 145 package (Bushnell, 2015) using the provided library of Illumina adapters. Quality trimming 146 was also performed using bbduk form the bbtools package, reads were trimmed to a minimum quality score of 20 over a sliding window of 10 bases, additionally the flag 147 148 'modulo = 5' was used to remove trailing single odd numbered bases (i.e. the 301st base in a 300bp library, or the 101st base in 100bp library), which are common error bases in Illumina 149 datasets. Read merging was performed using bbmerge from the bbtools package. Human 150 DNA contaminants were removed by using removehuman.sh from the bbtools package. 151 152 Briefly, raw reads are mapped onto a prebuilt index of the human genome which had been masked to hide 1) any low complexity repeat regions and 2) any regions which showed > 153 154 85% identity to any sequence in the Silva rRNA database (Quast et al., 2013) over a 70 basepair window. This method removes potential human DNA contamination while minimising 155 false positive hits to low complexity regions and ribosomal RNA sequences 156

157

164

Assembly benchmarking was performed using three *de-novo* genome assemblers optimized 158 for metagenomic data, Megahit (Li et al., 2016), SPAdes (Bankevich et al., 2012) (with the 159 160 flag: - meta), and IDBA-UD (Peng et al., 2012), using kmers ranging from 27-127 in steps of 161 10. Benchmarking was performed against un-normalized raw reads and reads normalized to minimum kmer depth of 3 and maximum kmer depth of 100 for kmers of size 32. 162 Normalization was performed using bbnorm form the bbtools package, with the flags 163 min=3, max=100, k = 32. Read recruitment for each assembly was estimated using bbmap

165 with the flags *k=13 vslow=t*. Annotation of quality trimmed reads was achieved using Kaiju (Menzel et al., 2016). A database of all proteins from all bacteria, archaea, single celled 166 167 eukaryotes and viruses in the NCBI non-redundant protein database 168 (NCBI Resource Coordinators, 2017) was constructed using the makeDB.sh script supplied with Kaiju. Nucleotide reads were translated into amino acid sequence in all 6 reading 169 frames and taxonomically annotated by alignment against this database using Kaiju, with a 170 171 maximum of 5 mismatches allowed and a minimum bit score of 60. Reads were projected 172 onto all taxonomic ranks from phylum to species and per-sample abundances were 173 compiled into single data tables for downstream statistical analysis using a custom bash 174 script. Functional annotations were assigned by mapping NCBI accession IDs from the Kaiju analysis onto functional classifications from the SEED subsystems protein hierarchy 175 (Overbeek et al., 2005) at levels 1, 2 and 3 using MEGAN 6 (Huson et al., 2016). 176 177

178 All taxonomic and functional annotations and read counts were concatenated and downstream analysis was performed using R 3.4.1 (R Core Team, 2011). Taxonomic and 179 functional abundances were summed at each taxonomic and functional rank and 180 normalized by dividing counts for each sample by the total number of reads that were 181 182 annotated for that sample using the *aggregate* function. Prior to statistical analysis, taxonomic and functional abundance values were Hellinger transformed as described in 183 184 Legendre and Gallagher, (2001) using the function *decostand* from the package vegan 185 (Oksanen et al., 2016). Soil abiotic variables were log transformed and standardized such that each variable had a mean of zero and standard deviation of 1 across all samples. 186 187

188 Euclidean distance dissimilarity matrices were produced for normalized abundance and abiotic variable tables using the *veqdist* function from vegan. Principal component analysis 189 was performed using the *capscale* function in vegan, samples were plotted in two 190 191 dimensional space using their first two principal components, plots were produced for 192 taxonomic abundances at the class and genus levels, for functional abundances, and for soil 193 abioitic variables. For taxonomic annotation tables, genus-level richness and Shannon-194 Weaver diversity indices were computed using the diversity function from the package 195 vegan. Beta diversity for each soil horizon was inferred for all taxonomic and functional ranks as the distance to the group centroid in Euclidean space using the betadisper function 196 from the package vegan. Analysis of variance (ANOVA) tests were performed to test for 197 significant differences between group beta dispersions using the function anova. Analysis of 198 similarities (ANOSIM) (Clarke, 1993) was performed on all samples grouped by biome type, 199 200 sample site (for Alpine soils only) and soil horizon using the function anosim from the 201 package vegan. Analysis of variance using distance matrices was performed on all samples 202 grouped by, sample site and soil horizon using the *adonis* function from the package vegan 203 with 9999 permutations. Taxa and functions whose abundances differed significantly 204 between horizons depths were identified by applying the Kruskal-Wallis H test (Kruskal and 205 Wallis, 1952) using the function kruskal.test, p-values were adjusted for multiple testing using the Benjamin-Hochberg false discovery rate method (Benjamini and Hochberg, 1995) 206 with the function *p.adjust*, only taxa with an adjusted p-value < 0.05 were considered 207 208 significant.

209

210 Data availability

Raw Sequence data files are available in the NCBI sequence read archive under Bioprojectnumber PRJNA39461.

213 **4 Results**

214 **DNA EXTRACTION, SEQUENCING DATA QC AND PROCESSING**

215 Vertical sections of weathered sediment (paleosols) in a series of five moraines in the Guil 216 river valley of the French Alps, (G1, G2, G3a, G9, and G11) were sampled at three discrete stratified soil horizons (Ah, Bw and Cox) and were classified as previously described 217 218 (Mahaney et al., 2016; 2017) (Figure 1). The geology of these paleosols is an active area of glacial, cosmic and microbiological research as described in several related geological papers 219 to date (Mahaney et al., 2013, 2016, 2017; Mahaney and Keiser, 2013). 220 221 Total DNA per gram of soil generally decreased with increasing depth, indicating a lower 222 total biomass at lower depths in the soil profile (supplementary table S3). Average sequence 223 read lengths were as expected (300bp), average sequence quality was found to be > 20 for the full length of each read for all samples. Adapter contamination was also found to be low, 224 on average < 1% per library. Although merging of the paired end reads was attempted, on 225 average < 30% of reads could be merged, indicating that the average insert size of the 226

227 sequencing libraries exceeded 600bp, therefore downstream processing was performed on

228 interleaved PE reads.

Attempts at assembly of the metagenomic libraries revealed generally poor assemblies (supplementary table S4) and a significant loss of information (i.e. the percentage of raw reads which could be mapped onto the assembly). Since the goal here was to characterize these soil samples in a general sense it was decided that the loss of information during assembly and

overall poor assembly statistics were deemed unacceptable. Therefore annotation was
performed using only the raw PE read files, since the reads were on average 300bp in length,
a reasonable level of taxonomic resolution can still be obtained, and in most cases
assignments could be made down to the genus level without much further loss of information
(supplementary figure S1).

On average ~ 80% of all reads could be assigned to at least the Phylum level, the functional gene assignment rates were significantly lower at ranging from 17-26% depending on the SEED subsystem level used. As taxonomic resolution increases so the number of reads which can be confidently assigned to a taxon decreases, with a significant drop-off in the number of reads which can be assigned to the species level, therefor for downstream analysis, the genus level was the minimum taxonomic rank which was analysed..

244 **TAXONOMIC / FUNCTIONAL ANNOTATION**

At the class level it is apparent that the community is dominated by four major taxa (Figure 245 2), the Actinobacteria, Alphaproteobacteria, Betaproteobacteria and Gammaproteobacteria 246 all of which are common genera of soil microbes which appear to be ubiquitous in soils 247 (Barberán et al., 2014). The Actinobacteria account for roughly 30% of the total annotated 248 249 reads in all samples. The Bacilli and Acidobacteriia also appear to be reasonably abundant, 250 as might be expected for paleosols with low pH 4.0-6.0. Certain horizon dependent patterns 251 in the data also emerge even in these boxplots, for example, the Betaproteobacteria and 252 Gammaproteobacteria appear to be generally more abundant in the Bw and Cox horizon samples while the *Ktedonobacteria* appear to be generally more abundant in the Ah horizon 253 254 samples.

255 ALPHA / BETA DIVERSITY ANALYSIS

256 Alpha diversity was estimated for each sample at each taxonomic level using the Shannon-Weaver diversity index (Shannon and Weaver, 1964), Inverse Simpson diversity index 257 258 (Simpson, 1949) and genus-level richness (i.e. number of unique genera per sample). Linear regression of alpha diversity measures against sample depth reveals a fairly strong, 259 statistically significant linear relationship between soil depth and alpha diversity (Figure 3). 260 261 Interestingly while richness appears to decrease with depth, the diversity indices appear to increase. Both inverse Simpson and Shannon diversity indices account for taxon abundances, 262 where a lower value indicates more uneven taxon abundances, suggesting that while the 263 264 number of unique genera decrease while moving down a soil profile, the relative abundances of these genera become more even. 265

In order to test the statistical significance of differences in community structure between 266 horizons, ANOSIM (Figure 3) and Adonis tests (supplementary table S2) were applied. It was 267 necessary to test for homogeneity of group dispersions, i.e. whether or not the multivariate 268 269 spread of the samples from their group centroids are significantly different (Figure 3 panel 270 B). Homogeneity of group dispersion tests were performed for samples grouped by both soil horizon and sample site, results were not statistically significant in any case, meaning that 271 downstream ANOSIM tests may be interpreted confidently without any caveats. The results 272 of both ANOSIM and Adonis tests indicate that soil horizon has a much stronger effect on 273 274 microbial community than sampling site: the ADONIS R statistic is reasonably high for all ranks 275 when grouped by soil horizon and much lower when grouped by sampling site, indicating that the effect is genuine. Statistical significance is strong with all ranks showing three star 276 277 significance (P < 0.001) when grouped by soil horizon, apart from the phylum level which

shows two star significance (P < 0.01) (supplementary table S3). When grouped by sampling
site, there is a marginally significant effect at the phylum level although this is not seen at
lower levels of taxonomic resolution, or at the functional gene level.

281

CORRELATIONS WITH SOIL ABIOTIC VARIABLES

The effect of soil horizon classification on soil abiotic variables was also assessed (Figure 4). 282 Out of the 27 variables which were assessed here [supplementary data file 1], 12 were 283 found to correlate significantly with soil depth and between stratified soil horizons. These 284 12 include 5 of the 6 essential elements for life (C, H, N, P, S), it can also be presumed that 285 286 available oxygen also decreases with soil depth. Where C,N and P were compered total rather than soluble levels were used, because we wanted to consider insoluble sources of 287 these elements (e.g. lignin). Given this information, it is unsurprising that the total DNA per 288 289 g of soil decreases with depth and that the microbial communities between soil horizons are 290 significantly different from one another (ANOSIM significance = 3.9e⁻⁴) (Figure 3). 291 Conversely, the concentration of numerous elements important for life (Ca, Cu, Na, K and Mn) appears to increase with soil depth – perhaps due to leaching. This is typical during 292 weathering of parent material, when the elements are realised they are leached down 293 profile (Retallack, 2001). The exception is K that follows decreasing trend toward the parent 294 295 rock. The K is most likely incorporated into secondary clay minerals like Illite, illite-smectite 296 and/or illite-vermiculite (Mahaney et al., 2016).

297 PRINCIPAL COMPONENT ANALYSIS

Principle component analysis was performed on taxon abundances (at the class and genus
levels), functional abundances (at SEED subsystems level 1) and abiotic variables (Figure 5),

300 revealing strikingly similar patterns between samples in multivariate space. There is a clear and consistent separation of soil horizons along the first principal component axis in all cases. 301 There is also a slight overlap between Bw and Cox horizons in the case of the Genus level 302 303 abundance and chemical variables, indicating that the Bw and Cox horizons are less strongly 304 differentiated from each other. This corresponds with the sedimentology - which indicates that the Bw horizon is only weakly differentiated and still developing (hence the w 305 classification), and therefore it is not surprising that these data suggest a closer degree of 306 307 similarity between population in the Bw and Cox horizons (Mahaney et al., 2013). In all cases, 308 the first two principle component axes cumulatively explain > 70% of the total variation (> 309 80% in the case of the taxonomic variation).

In order to identify the taxa and functions which contribute most to the differences between soil horizons the Kruskal-Wallace H test (Kruskal and Wallis, 1952) for differences between groups was applied to taxonomic abundances at the class level, and functional abundances at level 1 Of the SEED subsystems hierarchy. For clarity, plots were only produced for taxa and functions which occurred at an abundance greater than 0.1% in at least one sample.

315 It is clear that the abundances of many classes of bacteria vary along the depth profile (Figure 316 6). Nine of the top ten most significantly variable taxa appear to increase with soil depth, while the only one which decreases are the Ktedonobacteria. Likewise, classes which are 317 318 known to comprise many obligate and facultative anaerobes (the Bacilli, Clostridia and 319 Negativicutes) increase significantly with depth. The Betaproteobacteria are one of the most 320 abundant classes present in the dataset and show a near 3% increase in median abundance between the Ah and Cox horizons. A total of 57 Classes were shown to vary significantly 321 322 between soil horizons, applying this analysis to genus level abundances revealed that 1223

genera vary significantly between soil horizons. Overall, what has become clear from these
data is that a large number taxa vary significantly and consistently between soil horizons over
all six sample sites, even when measured at higher taxonomic ranks.

Similarly, many functional gene categories vary in abundance with respect to soil horizon, 326 although the effect is much less pronounced than for taxon abundances, with only 6 major 327 328 categories showing significant variation (Figure 7). Nitrogen metabolism and carbohydrate metabolism are the two which can be most obviously linked to ecosystem level processes. 329 330 Both decrease with increasing depth, corresponding to the decrease in soil nitrogen and organic carbon with depth recorded by Mahaney, et al. (2016). The low measured N 331 (supplementary data file 1) concentrations in these soils would suggest that the environment 332 is extremely nitrogen limited. Therefore we would expect that Ammonia oxidation may be an 333 334 extremely important component of the nitrogen cycle in these soils. Genes relating to central metabolism (i.e. metabolic pathways essential for organism survival) and secondary 335 336 metabolism (i.e. pathways which are non-essential) show a near inverse relationship, with secondary metabolism decreasing with depth while central metabolism genes increase with 337 depth. This may be explained by the fact that in lower soil horizons, where nutrients are less 338 abundant (Fierer, 2017), functions which are non-essential for survival are less vital and thus 339 are selected against over time - core metabolism genes become more important. 340

341

Additionally, in the Ah horizon where C N and P are more abundant we might expect antimicrobial secondary metabolites biosynthesis genes to be more prevalent due to the higher abundance of microbes competing for these nutrients. In lower soil horizons with less microbial activity, competition for nutrients is likely to be lower leading to a drop in the

346 abundance of anti-microbial biosynthesis genes, although this cannot be confirmed from the results of this study. This does however suggest the possibility that nutrient-rich surface soil 347 348 horizons may be ideal locations to screen for the presence of novel antimicrobial compounds 349 and future work should endeavour to test this hypothesis. Fatty acid and isoprenoid synthesis 350 genes are more abundant in the Ah horizon, this may be due to the fact these genes are 351 heavily involved in cell membrane maintenance and manufacture (Kaneda, 1991), and may indicate that the Ah horizon hosts the most actively growing microbial cells. Finally, 352 353 membrane transport related genes appear to be most abundant in the Cox horizons. The 354 reason for this trend is not immediately obvious, though one could speculate that it may be 355 connected to an increase in chemo-lithotrophic lifestyles in the lower mineral soil horizons where alternative electron donors (such as sulfide, ferric Iron or ammonia) must be 356 transported across the cell membrane before they can be used by the cell for ATP synthesis 357 358 (Peck 2003).

359 Carbohydrate metabolism is one of the largest SEED subsystem categories and contains many 360 subcategories of genes which are linked to the carbon cycle, therefore carbohydrate metabolism genes were analyzed in more detail (Figure 8). After multiple test corrections 361 were applied 5 categories within carbohydrate metabolism show significant variation. Poly, 362 363 mono, di and oligo saccharide metabolism genes all decrease moving down the soil profile, as do carbon dioxide fixation genes. This is hardly surprising due to the fact that organic carbon, 364 365 and carbon dioxide are generally present at higher concentrations in upper soil horizons. 366 However, one-carbon metabolism genes appear to increase with soil depth, as might be expected given the importance of the methane cycle in low organic carbon environments 367 (Serrano-Silva et al., 2014). A picture of the carbon cycle in these soils then begins to emerge 368 369 whereby organic carbon is primary energy and provides carbon sources for many microbes in

upper soil horizons as might be expected, while organisms in lower soil horizons gain energy
and carbon by utilizing reduced one-carbon compounds (e.g. methane) as energy sources
which are liberated from the anaerobic turnover of waste products from the upper horizons.

373 **5 Discussion**

It is clear that the abundances of many bacterial vary consistently between stratified soil 374 375 horizons (Figure 6). We found that the *Ktedonobacter* showed the largest differential 376 abundance between stratified soil horizon and was generally most abundant in the Ah horizons, members of this class are thought to be aerobic filamentous, spore-forming, gram-377 positive, heterotrophic bacteria (Yabe et al. 2017). Conversely, many classes of obligate and 378 facultative anaerobes decreased in abundance with soil depth (eg. the *Bacilli, Clastridiales*) 379 (Figure 3), suggesting a strong influence of oxygen availability on the community composition. 380 381 Interestingly, the most abundant classes of microbes across all samples (eg: Acinobacteria, 382 Proteobacteria) (Figure 2) did not show a statistically significant variation in abundance between soil horizons (Figure 6), perhaps suggesting that the more successful microbes in 383 384 these environments are ones with mixotrophic lifestyles who can adapt to variable conditions between soil horizons. 385

386

When comparing functional profiles between soil horizons, differences in certain
biogeochemical cycling genes become clear. It is evident that many gene abundances
correlate either positively or negatively with depth, which is understandable given that most
key nutrients will vary down profile. Nitrogen metabolism appears to decrease with
increasing depth, and likewise, carbohydrate metabolism genes appear to be significantly

392 more abundant in the Ah horizon, corresponding to the zone of highest organic carbon turnover (evident from CNP measurements, Mahaney et al., 2016). Taking a more detailed 393 394 look at carbohydrate metabolism (Figure 8) revealed that there are significant variations 395 even within this category. Carbon dioxide (autotrophy) and central carbohydrate 396 metabolism genes appear to be relatively abundant in the Ah horizon, while genes for other 397 one-carbon metabolism becomes more abundant in the Cox horizon. This is consistent with increasing abundance of methane (associated with both methylotrophy and 398 399 methanogenesis) turnover processes in deeper soils (Dunfield, 2007). This is likely due to the fact that, as complex plant polymers are broken down by microbes in the upper soil 400 401 horizons, reduced carbon compounds are produced as waste products and filter down to the lower horizons where they are then utilized by methylotrophic microbes. 402 403

404 It should be noted that there is significant co-linearity between the many abiotic and biotic 405 variables measured and compared here. As an example, looking specifically at almost all other 406 abiotic and biotic parameters correlate with soil depth (in cm) to some degree (Figure 5). 407 Therefore it was not possible to investigate and compare correlations between individual taxa/functions and abiotic variables as disentangling this co-linearity to reveal the true source 408 of variation is not possible without carefully designed laboratory experiments to control for 409 410 the effect of specific co-linear variables identified here. Indeed, any attempt to do so may be statistically dubious and may therefore be considered as an example of "p-hacking" (Head et 411 al., 2015). However, this analysis makes it abundantly clear that soil horizon matters: in terms 412 of chemistry, taxonomy and functionality of the community. 413

414

415 There are several important caveats to the results presented here that must also be considered. Firstly, metagenomes are not expected to be static over time, so caution must 416 417 be taken when making broad statements about a microbial community. In extreme cases, 418 we might expect drastic changes in microbial community structure to (Mahaney et al., 419 2017). Otherwise, previous studies have established that temporal variation within soil 420 microbiomes is generally much lower than spatial variation (Lauber et al., 2013), so the 421 patterns observed here are likely to persist across seasons, at least in the short term. 422 Secondly, for any enzyme-catalysed processes to occur, transcription and translation of protein coding genes are essential requirements. Therefore, a positive relationship between 423 the abundance of gene or transcripts and corresponding process rate is not always 424 necessarily true, though it is often presumed. Rocca et al., (2015) indicated that functional 425 gene abundances are only weakly correlated to process rates, but are consistently 426 427 correlated across multiple environments. Finally, the relatively low depth of sequencing per 428 sample means that complete genomes cannot be resolved, (a common limitation for diverse 429 soil metagenomes (Nesme et al. 2016) and annotations are based on ~300bp fragments. 430 Much more robust results could be obtained with more complete genomes and functional genes, allowing investigation of functional genes in their genomic context, and analysis of 431 complete metabolic pathways within certain species. The relatively short read length used 432 here may be useful for taxonomic classification but full length genes are preferable for 433 functional annotation. Despite this limitation, there is a growing body of evidence to suggest 434 that shallower sequencing across many samples with suitable replication levels is sufficient 435 to answer key questions about microbial ecology (Knight *et al.*, 2012). 436

437

Fundamentally we know that soil stratigraphy is a function of soil chemistry (USDA 1999).
Our full metagenome data, for the first time, show that soil stratigraphy is strongly
correlated with variation in microbial community composition – as shown in figures 3 and 5.
Furthermore we find that the variation seen between stratified soil horizons is greater than
the variation seen between soil sampling sites when taken as a whole (Figure 3).

443

Ultimately, the physical and chemical properties of a soil are the true drivers of microbial 444 445 diversity and soil horizons are a useful method for categorizing soils with similar physicochemical characteristics. While the gene analysis presented here has provided an 446 447 overview of microbial community structure and functionality in alpine soils, it has also highlighted a need for more specific methods in order to make definitive statements about 448 biogeochemical processes occurring between soil horizons. The use of more specific and 449 450 informative functional gene ontologies or the use of meta-transcriptomic information will 451 certainly add a further level of understanding to this field. One correlation that we did not 452 investigate in this study, is the relationship between oxygen concentration and/or redox 453 potential with metagenome composition down profile. In general we would expect these to decrease with depth. Clearly there is abundant microbial activity in the sediments, and 454 oxygen – as the ultimate acceptor – would certainly become limiting to microbial 455 metabolism within a cm or two of the Ah horizon (Birkeland, 1999). 456

457

458 CONCLUSION

In conclusion, our study begins to demonstrate a clear relationship between the physicochemical parameters that are typically used to define soil pedons, and parallel the structure
of complex microbial communities. The ability to do so has only become available to us in

462 recent years due to the advent of shotgun sequencing and metagenomic analysis. The limitations of culture-dependant microbial analysis, have previously not allowed any such 463 464 connection to be established. However, in this study some dependency is evident: the reduction between organic carbon availability and the related genetic profile of functional 465 genes involved in C1 turnover is especially clear and is quite rational. The results here may 466 lead us to consider that the formation of paleosol horizons is inextricably linked to the 467 biological, not just physical and chemical transformations that occur. Critically, we 468 469 emphasise here a key finding in this study – that variation in microbial gene populations between palesol samples sites, sometimes hundreds of metres apart, is significantly less 470 471 than that seen between soil pedons, just centimetres apart at specific sample sites. This remarkable metagenome diversity may be exploited for gene mining applications. 472

473 Acknowledgements

474 We acknowledge the REMEDIATE KTN (CCRA, BK), Invest Northern Ireland (TS), and

475 Quaternary Surveys (WM) for financial support.

476

477 **6 REFERENCES**

479 http://www.bioinformatics.babraham.ac.uk/projects/fastqc.¶

480

483 ISME J 6: 248–258.¶

⁴⁷⁸ Andrews S. (2010). FastQC: a quality control tool for high throughput sequence data.

Baldrian P, Kolařík M, Štursová M, Kopecký J, Valášková V, Větrovský T, et al. (2012). Active and total
 microbial communities in forest soil are largely different and highly stratified during decomposition.

485 486 487	Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. (2012). SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. J Comput Biol 19: 455–477.¶
488	
489 490	Barberán A, Ramirez KS, Leff JW, Bradford MA, Wall DH, Fierer N. (2014). Why are some microbes more ubiquitous than others? Predicting the habitat breadth of soil bacteria. Ecol Lett 17: 794–802.¶
491	
492	Birkeland, P.W., 1999. Soils and Geomorphology. Oxford University Press, Oxford, UK, p. 430.¶
493	
494 495	Bushnell B. (2015). BBMap short-read aligner, and other bioinformatics tools. Available from sourceforge net/projects/bbmap.¶
496	
497 498	Byers AC. (2008). An assessment of contemporary glacier fluctuations in Nepal's Khumbu Himal using repeat photography. Himal J Sci 4: 21-26.¶
499	
500 501	Canada Soil Survey Committee (CSSC), (1998). The Canadian System of Soil Classification, 637. NRC Research Press, Ottawa, Canada, p. 187.¶
502	
503 504	Clarke KR. (1993). Non-parametric multivariate analyses of changes in community structure. Austral Ecol 18: 117–143. ¶
505	
506 507	Davidson EA, Janssens IA. (2006). Temperature sensitivity of soil carbon decomposition and feedbacks to climate change. Nature 440: 165–173.¶
508	
509 510	Ding J, Zhang Y, Deng Y, Cong J, Lu H, Sun X, et al. (2015). Integrated metagenomics and network analysis of soil microbial community of the forest timberline. Sci Rep 5: 7994.¶
511	
512	Dunfield, PF. (2007). 10 The Soil Methane Sink. Greenhouse gas sinks, p.152.¶
513	
514 515	Fierer N, Strickland MS, Liptzin D, Bradford MA, & Cleveland CC. 2009. Global Patterns in Belowground Communities. Ecology Letters 12: 1238–49.¶
516	
517	Fierer N. (2017). Embracing the unknown: disentangling the complexities of the soil microbiome. Nat

518 Rev Microbiol 15: 579–590.¶

519 520	Head ML, Holman L, Lanfear R, Kahn AT, Jennions MD. (2015). The extent and consequences of p- hacking in science. PLoS Biol 13: e1002106.¶
521	
522 523	He S, Wurtzel O, Singh K, Froula JL, Yilmaz S, Tringe SG, et al. (2010). Validation of two ribosomal RNA removal methods for microbial metatranscriptomics. Nat Methods 7: 807–812.¶
524	
525 526	Hodgson, JM. (1976). Soil Survey Field Handbook d Soil Survey Tech, Monograph, No. 5. 773. Rothamsted Experimental Stn, Harpenden, Herts, p. 99.¶
527	
528 529 530	Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, Mitra S, et al. (2016). MEGAN Community Edition - Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing Data. PLOS Comput Biol 12: e1004957.¶
531	
532 533	Kaneda T. (1991). Iso- and anteiso-fatty acids in bacteria: biosynthesis, function, and taxonomic significance. Microbiol Rev 55: 288–302.¶
534	
535 536 537	King AJ, Karki D, Nagy L, Racoviteanu A, Schmidt SK. (2011). Microbial biomass and activity in high elevation (> 5100 meters) soils from the Annapurna and Sagarmatha regions of the Nepalese Himalayas. Himal J Sci 6: 11-18¶
538	
539 540	Knight R, Jansson J, Field D, Fierer N, Desai N, Fuhrman JA, et al. (2012). Unlocking the potential of metagenomics through replicated experimental design. Nat Biotechnol 30: 513–20.¶
541	
542 543	Körner, Christian, and Jens Paulsen. (2004). A World-Wide Study of High Altitude Treeline Temperatures. Journal of Biogeography 31: 713–32.¶
544	
545 546	Kruskal WH, Wallis WA. (1952). Use of ranks in one-criterion variance analysis. J Am Stat Assoc 47: 583–621.¶
547	
548 549	Lal, R. (2004). Soil carbon sequestration impacts on global climate change and food security. Science 304: 1623-1627.
550	
551 552	Lauber CL, Ramirez KS, Aanderud Z, Lennon J, Fierer N. (2013). Temporal variability in soil microbial communities across land-use types. ISME J 7: 1641–1650.¶
553	

554 555	Lazzaro A, Hilfiker D, Zeyer J. (2015). Structures of Microbial Communities in Alpine Soils: Seasonal and Elevational Effects. Front Microbiol 6: 1330.¶
556	
557 558	Legendre P, Gallagher ED. (2001). Ecologically meaningful transformations for ordination of species data. Oecologia 129: 271–280.¶
559	
560 561 562	Li D, Luo R, Liu CM, Leung CM, Ting HF, Sadakane K, Yamashita H, Lam TW. (2016). MEGAHIT v1. 0: A fast and scalable metagenome assembler driven by advanced methodologies and community practices. Methods102: 3-11.
563	
564 565	Mahaney WC, Keiser L. (2013). Weathering rinds-unlikely host clasts for an impact- induced event. Geomorphology 184: 74-83¶
566	
567 568 569 570	Mahaney, WC, Somelar, P, West A, Krinsley D, Allen CCR, Pentlavalli P, Young JM, Dohm JM, LeCompte M, Kelleher B, Sean J, Pulleyblank C, Dirszowsky R, Costa P. (2017) Evidence for cosmic airburst/impact in the Western Alps archived in Late Glacial Paleosols. Quaternary International, 438: 68-80.¶
571	
572 573 574	Mahaney WC, Keiser L, Krinsley DH, Pentlavalli P, Allen CCR, Somelar P, et al. (2013). Weathering rinds as mirror images of paleosols: examples from the Western Alps with correlation to Antarctica and Mars. J Geol Soc London 2012: 150.
575	
576	
577 578 579 580	Mahaney WC, Somelar P, Dirszowsky RW, Kelleher B, Pentlavalli P, McLaughlin S, et al. (2016). A Microbial Link to Weathering of Postglacial Rocks and Sediments, Mount Viso Area, Western Alps, Demonstrated through Analysis of a Soil/Paleosol Bio/Chronosequence. Journal of Geology. 124: 149–169.¶
581	
582 583	Menzel P, Ng KL, Krogh A, Marth G, Lipman D. (2016). Fast and sensitive taxonomic classification for metagenomics with Kaiju. Nat Commun 7: 11257.¶
584	
585 586 587	Mitra S, Rupek P, Richter DC, Urich T, Gilbert JA, Meyer F, Wilke A, Huson DH. (2011) Functional analysis of metagenomes and metatranscriptomes using SEED and KEGG. BMC bioinformatics. 12: S21.
588	
589 590	National Soil Survey Center (NSSC), 1995. Soil Survey Laboratory Information Manual. Soil Survey Investigations Report No. 45. Version 1.00. USDA, Washington, DC., p. 305¶

591	
592 593	NCBI Resource Coordinators. (2017). Database Resources of the National Center for Biotechnology Information. Nucleic Acids Res 45: D12–D17.¶
594	
595 596	Neira, J., Ortiz, M., Morales, L., & Acevedo, E. (2015). Oxygen diffusion in soils: Understanding the factors and processes needed for modeling. Chilean journal of agricultural research 75: 35-44.¶
597	
598 599	Nesme J, Achouak W, Agathos SN, Bailey M, Baldrian P, Brunel D, et al. (2016). Back to the Future of Soil Metagenomics. Front Microbiol 7: 73.¶
600	
601 602	Oksanen J, Blanchet FG, Kindt R, Legendre P, Minchin PR, O'hara RB, Simpson GL, Solymos P, Stevens MH, Wagner H (2011). vegan: Community ecology package. R package version. 2011:117-8.
603	
604 605 606	Overbeek R, Begley T, Butler RM, Choudhuri J V, Chuang H-Y, Cohoon M, et al. (2005). The Subsystems Approach to Genome Annotation and its Use in the Project to Annotate 1000 Genomes. Nucleic Acids Res 33: 5691–5702.¶
607	
608 609	Oyama, M., Takehara, H., (1970). Standard Soil Color Charts. Japan Research Council for Agriculture, Forestry and Fisheries, Tokyo, Japan.
610	
611 612	Peck Jr, H. D. (1968). Energy-coupling mechanisms in chemolithotrophic bacteria. Annual Reviews in Microbiology 22: 489-518.¶
613	
614 615	Peng Y, Leung HCM, Yiu SM, Chin FYL. (2012). IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. Bioinformatics 28: 1420–1428.¶
616	
617 618 619	Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. (2013). The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. Nucleic Acids Res 41: D590–D596.¶
620	
621 622	R Core Team. (2011). R: A Language and Environment for Statistical Computing. R Found Stat Comput 1: 409.¶
623	
624 625	Retallack GJ, (2001). Soils of The Past: An Introduction to Paleopedology, 2nd ed., Blackwell, Oxford, UK. 600 p.
626	

627 628 629	Rocca JD, Hall EK, Lennon JT, Evans SE, Waldrop MP, Cotner JB, et al. (2015). Relationships between protein-encoding gene abundance and corresponding process are commonly assumed yet rarely observed. ISME J 9: 1693–9.¶
630	
631 632	Serna-Chavez HM, Fierer N, van Bodegom PM. (2013). Global drivers and patterns of microbial abundance in soil. Glob Ecol Biogeogr 22: 1162–1172.¶
633	
634 635	Serrano-Silva N, Sarria-Guzan Y, Dendooven L, Luna-Guido M. (2014). Methanogenesis and Methanotrophy in Soil: A Review. Pedosphere 24: 291–307.¶
636	
637 638	Shannon CE, Weaver W. (1964). The mathematical theory of communication. University of Illinois Press. \P
639	
640	Simpson, E.H. (1949). Measurement of diversity. Nature.¶
641	
642	Tukey JW (John W. (1977). Exploratory data analysis. Addison-Wesley Pub. Co.¶
643	
644 645	USDA. (1999). Soil taxonomy: A basic system of soil classification for making and interpreting soil surveys. 2nd edition.¶
646 647 648	Weigold P, El-Hadidi M, Ruecker A, Huson DH, Scholten T, Jochmann M, et al. (2016). A metagenomic-based survey of microbial (de)halogenation potential in a German forest soil. Sci Rep 6: 28958.¶
649	
650 651	Yabe S, Sakai Y, Abe K, Yokota A. (2017). Diversity of Ktedonobacteria with Actinomycetes- Like Morphology in Terrestrial Environments. Microbes Env 32: 61–70.¶
652	
653 654 655	Yashiro E, Pinto-Figueroa E, Buri A, Spangenberg JE, Adatte T, Niculita-Hirzel H, et al. (2016). Local Environmental Factors Drive Divergent Grassland Soil Bacterial Communities in the Western Swiss Alps. Appl Environ Microbiol 82: 6303–6316.¶
656	
657 658 659	Yu, K., Faulkner, S. P., & Patrick Jr, W. H. (2006). Redox potential characterization and soil greenhouse gas concentration across a hydrological gradient in a Gulf coast forest. Chemosphere 62: 905-914.¶
660	
661 662	Zemp M, Haeberli W, Hoelzle M, Paul F. (2006). Alpine glaciers to disappear within decades? Geophys Res Lett 33: L13504¶

Zhang Y, Lu Z, Liu S, Yang Y, He Z, Ren Z, et al. (2013). Geochip-based analysis of microbial
communities in alpine meadow soils in the Qinghai-Tibetan plateau. BMC Microbiol 13: 72.

667 FIGURE LEGENDS

- Figure 1: Overview of Sampling sites and soil stratigraphy. Panel A: Satellite map of sampling
 area. Panel B: Stratigraphy of soil profiles for each sample (from Mahaney *et al.*, 2016).
- 670 Sampling sites are shown at the top of each profiles, depth in centimeters is displayed on
- the left and soil horizon classification is displayed on the right of each profile.
- 672 Figure 2: Boxplot showing the relative abundances of the 20 most abundant Taxa across all

673 samples at the Class level. Jittered points overlaid on boxplots represent individual sample

and are coloured by sample soil horizon classification. The bottom and top of the boxes

represent the first and third quartiles, with the central band representing the median,

676 whiskers represent 1.5 times the interquartile range according to Tukey's schematic boxplot

677 method (Tukey, 1977)

Figure 3: Diversity analysis of alpine soils. Panel A: Alpha diversity plots for taxonomic 678 679 abundance tables at each taxonomic level. Y axis represents inverse Simpson diversity index, 680 Shannon-Weaver diversity index, and Genus level richness (i.e. the number of unique genera 681 detected per sample) respectively. X axis represents depth in cm of soil sample. Blue lines represent linear regressions, shaded area represents the 95% confidence interval for the 682 regression analysis. R² values and P-values for linear regressions along with significance 683 684 codes are displayed for each plot. Panel B: homogeneity of group dispersions for samples 685 grouped by soil horizon, Y axis represents distance to the group centroid in 2d Euclidean space, standard error of the mean across all samples, analysis of variance (ANOVA) test P 686

value for differences between soil horizons value is displayed on the plot. Panel C: Boxplot
of analysis of similarities (ANOSIM) results, Y axis represents ranked order of dissimilarities
ANOSIM R statistic and significance level are shown on the plot, significance of the R statistic
was assessed by permutation for 9999 replicates. The bottom and top of the boxes
represent the first and third quartiles, with the central band representing the median,
whiskers represent 1.5 time the interquartile range according to Tukeys schematic boxplot
method (Tukey, 1977).

694 Figure 4: Analysis of effect of paleosol depth and paleosol horizon on paleosol abiotic 695 variables. Panel A: Spearmans rank correlation matrix for pairwise correlations between soil geochemical variables. Points are coloured and scaled according to the value of the 696 correlation coefficient, only correlations with a P value < 0.05 are displayed. Variables which 697 698 correlate significantly with soil depth are highlighted in red. PanelB: Boxplots of geochemical variables which vary significantly between stratified soil horizons. The bottom 699 700 and top of the boxes represent the first and third quartiles, with the central band 701 representing the median, whiskers represent 1.5 time the interquartile range according to Tukeys schematic boxplot method (Tukey, 1977). 702

Figure 5: PCoA plots of abundance tables at each taxonomic and functional rank. Constrained analysis of principal components was performed on dissimilarity matrices using the function capscale from the package Vegan (Oksanen et al., 2016) in R version 3.4.1 (R Core Team, 2011). In each case the x and y axis represent the first and second principal component axis respectively. Points are coloured by soil horizon and the convex hulls are drawn and highlighted for each horizon grouping.

Figure 6: Boxplots showing relative abundance of microbial classes which differed significantly
between soil horizons. Relative abundance is expressed as a percentage of the total annotated reads.
The bottom and top of the boxes represent the first and third quartiles, with the central band

representing the median, whiskers represent 1.5 time the interquartile range according to Tukey's
boxplot method (Tukey, 1977). All data points are also plotted as points.

Figure 7: Boxplots of relative abundance of SEED level 1 functional categories. Relative abundance is expressed as a percentage of the total annotated reads. The bottom and top of the boxes represent the first and third quartiles, with the central band representing the median, whiskers represent 1.5 times the interquartile range according to Tukey's boxplot method (Tukey, 1977). All data points are also plotted as jittered points.

Figure 8: Abundance of SEED level 2 categories related to carbon metabolism. Relative abundance is expressed as a percentage of the total annotated reads. The bottom and top of the boxes represent the first and third quartiles, with the central band representing the median, whiskers represent 1.5 time the interquartile range according to Tukey's boxplot method (Tukey, 1977). All data points are also plotted as points.





725 Figure 1: Overview of Sampling sites and soil stratigraphy. Panel A: Satellite map of sampling

area. Panel B: Stratigraphy of soil profiles for each sample (from Mahaney *et al.*, 2016).

727 Sampling sites are shown at the top of each profiles, depth in centimetres is displayed on

the left and soil horizon classification is displayed on the right of each profile.



Taxon (Class Level)

Figure 2: Boxplot showing the relative abundances of the 20 most abundant Taxa across all
samples at the Class level. Jittered points overlaid on boxplots represent individual sample
and are coloured by sample soil horizon classification. The bottom and top of the boxes
represent the first and third quartiles, with the central band representing the median,

- whiskers represent 1.5 times the interquartile range according to Tukey's schematic boxplot
- 735 method (Tukey, 1977)







753 **Figure 4:** Analysis of effect of paleosol depth and paleosol horizon on paleosol abiotic

variables. Panel A: Spearmans rank correlation matrix for pairwise correlations between soil

755 geochemical variables. Points are coloured and scaled according to the value of the

correlation coefficient, only correlations with a P value < 0.05 are displayed. Variables which

correlate significantly with soil depth are highlighted in red. **PanelB:** Boxplots of

758 geochemical variables which vary significantly between stratified soil horizons. The bottom

and top of the boxes represent the first and third quartiles, with the central band

representing the median, whiskers represent 1.5 time the interquartile range according to

761 Tukeys schematic boxplot method (Tukey, 1977).



762

Figure 5: PCoA plots of abundance tables at each taxonomic and functional rank. Constrained analysis
 of principal components was performed on dissimilarity matrices using the function capscale from the
 package Vegan (Oksanen et al., 2016) in R version 3.4.1 (R Core Team, 2011). In each case the x and y
 axis represent the first and second principal component axis respectively. Points are coloured by soil

horizon and the convex hulls are drawn and highlighted for each horizon grouping.



768

Figure 6: Boxplots showing relative abundance of microbial classes which differed significantly 769 770 between soil horizons. Relative abundance is expressed as a percentage of the total annotated reads. 771 The bottom and top of the boxes represent the first and third quartiles, with the central band 772 representing the median, whiskers represent 1.5 time the interquartile range according to Tukey's 773 boxplot method (Tukey, 1977). All data points are also plotted as points.









Figure 8: Abundance of SEED level 2 categories related to carbon metabolism. Relative abundance is
expressed as a percentage of the total annotated reads. The bottom and top of the boxes represent
the first and third quartiles, with the central band representing the median, whiskers represent 1.5
time the interquartile range according to Tukey's boxplot method (Tukey, 1977). All data points are
also plotted as points.