# Article (refereed) - postprint

1  **Detecting macroecological patterns in bacterial communities across independent studies of**

2  **global soils**

3

4  **Authors:** Kelly S Ramirez[*+1], Christopher G. Knight[+2], Mattias Hollander[1], Francis Q. Brearley[3],

5  Bede Constantinides[4], TE Anne Cotton[5], Si Creer[6], Thomas W. Crowther[1], John Davison[7],

6  Manuel Delgado-Baquerizo[8], Ellen Dorrepaal[9], David R. Elliott[3,10], Graeme Fox[3], Rob

7  Griffiths[11], Chris Hale[12], Kyle Hartman[13], Ashley Houlden[14], David L. Jones[6], Eveline J. Krab[9],

8  Fernando T. Maestre[15], Krista L. McGuire[16], Sylvain Monteux[9], Caroline H. Orr[17], Wim H van

9  der Putten[1,18], Ian S. Roberts[14], David A. Robinson[19], Jenny Rocca[20], Jennifer Rowntree[3], Klaus

10  Schlaeppi[13], Matthew Shepherd[21], Brajesh K. Singh[22], Angela Straathof[2], Jennifer M. Talbot[23],

11  Cécile Thion[24], Marcel van der Heijden[13], and Franciska T. de Vries[2]

12

13  * Corresponding author

14  + Joint lead authors

15  [1] Netherlands Institute of Ecology, Droevendaalsesteeg 10 6708 PB Wageningen, The Netherlands

16  [2] School of Earth and Environmental Sciences, d Manchester, M13 9PT, United Kingdom

17  [3] School of Science and the Environment, Manchester Metropolitan University, Chester Street, Manchester, M1 5GD, United Kingdom

18  [4] Evolution and Genomic Sciences, School of Biological Sciences, The University of Manchester, Manchester, M13 9PT, United Kingdom

19  [5] Department of Animal and Plant Sciences, The University of Sheffield, Alfred Denny building, Sheffield, South Yorkshire, S10 2TN, United

20      Kingdom

21  [6] School of the Environment, Natural Resources & Geography, Bangor University, Gwynedd, LL57 2UW, United Kingdom

22  [7] Department of Botany, Institute of Ecology and Earth Sciences, University of Tartu, Lai 40, Tartu 51005, Estonia

23  [8] Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder, CO 80309, USA

24  [9] Climate Impacts Research Centre, Department of Ecology and Environmental Science, Umeå University, Vetenskapensväg 38, 981 07, Abisko,

25      Sweden

26  [10] Environmental Sustainability Research Centre, University of Derby, Kedleston Road, Derby, DE22 1GB, United Kingdom

27  [11] Centre for Ecology and Hydrology, Wallingford, United Kingdom

28  [12] School of Life Sciences, University of Warwick, Coventry, CV4 7AL, United Kingdom

29  [13] Department of Agroecology and Environment, Agroscope, Zurich, Reckenholzstrasse 191, Switzerland

30    [14] Faculty of Biology, Medicine and Health, The University of Manchester, Manchester, M13 9PT, United Kingdom

31    [15] Departamento de Biología y Geología, Física y Química Inorgánica, Escuela Superior de Ciencias Experimentales y Tecnología, Universidad

32         Rey Juan Carlos, Calle Tulipán s/n, 28933 Móstoles, Spain

33    [16] Barnard College, Columbia University Department of Biology; Columbia University, USA

34    [17] School of Science and Engineering, Teesside University, Middlesbrough, TS1 3BX, United Kingdom

35    [18] Laboratory of Nematology, Wageningen University, Droevendaalsesteeg 1, Wageningen 6708 PB, The Netherlands

36    [19] Centre for Ecology and Hydrology, Bangor, LL57 2UW, United Kingdom

37    [20] Department of Biology, Duke University, Durham, NC, 27705, USA

38    [21] Natural England, United Kingdom

39    [22] Hawkesbury Institute for the Environment, Western Sydney University, Richmond 2753 NSW Australia

40    [23] Department of Biology, Boston University, Boston, MA, 02215, USA

41    [24] Institute of Biological and Environmental Sciences, University of Aberdeen, Saint Machar Drive, AB24 3UU, Aberdeen, United Kingdom

42

43 **Keywords:** microbial ecology, soil, diversity, community structure, Illumina sequencing, 16S

44 rRNA gene, biogeography, microbiology, meta-analysis

45

46 **The emergence of high-throughput DNA sequencing methods provides unprecedented**

47 **opportunities to further unravel bacterial biodiversity and its worldwide role from human**

48 **health to ecosystem functioning. However, in spite of the abundance of sequencing studies,**

49 **combining data from multiple individual studies to address macroecological questions of**

50 **bacterial diversity remains methodically challenging and plagued with biases. Here, using a**

51 **machine learning approach that accounts for differences among studies and complex**

52 **interactions among taxa, we merge 30 independent bacterial datasets consisting of 1,998**

53 **soil samples from across 21 countries. While previous meta-analysis efforts have focused on**

54 **bacterial diversity measures or abundances of major taxa, we show that disparate**

55 **amplicon sequence data can be combined at the taxonomy-based level to assess bacterial**

56 **community structure. We find that rarer taxa are more important for structuring soil**

57 **communities than abundant taxa, and that these rarer taxa are better predictors of**

58 **community structure than environmental factors, which are often confounded across**

59 **studies. We conclude that combining data from independent studies can be used to explore**

60 **novel patterns in bacterial communities, identify potential 'indicator' taxa with an**

61 **important role in structuring communities, and propose new hypotheses on the factors that**

62 **shape bacterial biogeography previously overlooked.**

63

64 Soil microbial communities are more diverse and contain more individuals than any species

65 groups on the planet[1,2]. Over the last decade, the use of high-throughput sequencing (HTS)

66 methods has substantially advanced our understanding of the worldwide biogeography and

67 ecology of soil bacterial and fungal communitie[3–6]. Recent work has further demonstrated that

68 inclusion of microbial composition and functional attributes improves earth system models[7,8],

69 which is of paramount importance for predicting effects of global change on ecosystem services

70 such as climate regulation or soil fertility[9,10]. Yet, opposite to the long-standing view that every

71 organism may occur everywhere[11], even at small scales bacterial communities turn out to be

72 more patchy than previously expected[12,13], raising questions regarding dispersal constraints,

73 temporal dynamics, and niche breadth at the global scale[14–17]. Due to these knowledge gaps,

74 combined with practical challenges of exhaustive sample collection and the massive diversity of

75 communities, global assessment of soil microbial diversity remains an ongoing research

76 challenge[18,19].

77

78 For plants and animals, the integration of data from independent studies has been a valuable

79 option for generating an understanding of global biogeography patterns, answering ecological

80    questions (e.g. biodiversity-functioning relationships), and identifying threats to biodiversity

81    from global changes[20–23]. Similarly, our understanding of soil microbial diversity would greatly

82    improve from such worldwide assessments. However, the integration of microbial community

83    HTS data from different studies is not so unlike the merging of museum species records where

84    information and data is constrained by variations in nomenclature over space and time, among

85    many other challenges[24,25]. Like plant and animal records, molecular microbial community

86    records and information can be incomplete, processing and naming varies greatly between

87    studies and over time[26], data storage is inconsistent, and there are few curated databases with

88    high quality data (especially for short read sequences)[27,28]. Further, most microbial community

89    data and metadata are still available only in independently published studies that have been

90    carried out according to their own standards and procedures, and the extent of these confounding

91    factors has never been quantified across studies.

92

93    Regardless of the challenges, as indicated by the many open access data initiatives[29–31], merging

94    microbial sequence data is a potential option to address global scale questions, whether relating

95    to the human microbiome[32], marine systems[33], or predicting the response of soil organisms to

96    global environmental change[34]. For soil systems, the need to merge sequence data is supported

97    by the emerging role of bacterial phyla and classes as indicators of particular soil conditions such

98    as soil pH and nutrient concentrations[35,36]. Until now, attempts to meta-analyze sequence data

99    have been limited to assessing diversity measures or abundances of major taxa, because the

100   merging of community data is constrained by methodological differences between sequencing

101   studies[13,30,37–39]. However, a recent systematic review found that measures of microbial

102   community structure were more often linked to microbial process rates than diversity or

103  presence/absence data[40], and abundance ratios among phyla may be less important than previous

104  believed[41]. Together indicating that information on variation in microbial community structure is

105  potentially more ecologically relevant than measures of diversity and abundances of major taxa.

106

107  Here, we show that, despite the outlined challenges, published microbial community data from

108  independent studies can be analyzed together to address questions about the global structuring of

109  communities. Using a novel machine learning approach, we take methodological and technical

110  biases into account, factor in interactions among taxa, and produce an improved assessment of

111  the abiotic and biotic drivers of soil community structure. The objectives of this study were two-

112  fold: (1) to identify the biases and incompatibilities of microbial community HTS studies (and

113  confounding factors) so as to strengthen our ability to integrate data from disparate studies, and

114  (2) to reveal worldwide soil microbial community patterns by merging independent taxonomy-

115  based datasets.

116

117  **Results and Discussion**

118  **Taxonomy-based merging of disparate amplicon sequence data**

119  We identified 30 individual HTS bacterial studies from 21 countries for our analysis (Figure 1A

120  and Supplementary Table 1). While we aimed to merge HTS data of both soil bacterial and

121  fungal datasets, our approach was only successful for bacterial data (Figure 1B and 1C), and

122  highlights the well-known dilemma of fungal databases, where extremely high diversity

123  combined with high endemism and mismatched taxonomy across continents make merging data

124  by taxonomy difficult and unusable for downstream analyses[4,5,42]. For the bacterial studies, we

125  were able to successfully merge 30 individual OTU tables; using a taxonomy-based approach,

126  datasets were merged using the taxonomic affiliations of individual OTUs. Once filtered, and

127  singletons removed, the final 'taxonomy-based' community contained 1,998 individual soil

128  samples, and 8,287 taxa.  Here 'taxon' is defined as a unique name in the classification; where a

129  name could be a specific phylum, genus, or other taxonomic level. For example, 'Acidovorax'

130  (genus) and Proteobacteria (the phylum containing Acidovorax) were both considered as taxa).

131  To account for variation in sequencing depth between different studies, OTU relative abundances

132  were used per sample, rather than absolute read abundance. To test known biogeographical

133  patterns, metadata (information on geographical location, soil pH and soil core measurements)

134  were compiled for all studies. Technical and methodical information was also collected; all of

135  these 30 studies had conducted amplicon sequencing on hypervariable regions of the 16S rRNA

136  gene in soil samples using either Illumina or (Roche) 454 pyrosequencing (with any primer pair)

137  (Supplementary Table 1). For a validation step we retrieved all usable raw sequence data

138  available, resulting in 417 samples from locations across the globe (approximately 1/5 of all our

139  samples) (Figure 1A). Data not included in this sequence-matched analysis either had an

140  incompatible raw sequence format or simply no longer existed. Available raw sequence data

141  were combined into a single 'sequence-matched' community comprising 44,106 OTUs

142  (Supplementary Figure 1).

143

144  **Machine learning assessment of bacterial community structure**

145  Ordination of the taxonomy-based community reveals large amounts of structure both within and

146  between studies (structure that is removed by permuting taxa among samples (Supplementary

147  Figure 2), without greatly affecting diversity (Supplementary Table 3)), and the observation of

148  the well-established negative relationship between relative abundance of Acidobacteria and soil

149   pH (Figure 1D)[43] confirms our merging method. This visualization also suggests that some of the

150   community variation (e.g. the near absence of Acidobacteria in some studies, even at low pH) is

151   due to technical factors such as the particular primer sets chosen, region sequenced, and

152   sequencing platform (Supplementary Methods and Supplementary Table 2). However, we expect

153   that some taxa are not correlated with technical factors, and are non-randomly distributed with

154   respect to biotic and abiotic factors. Therefore, using a machine learning approach capable of

155   accounting for complex interactions among taxa (Random Forests[TM], see methods), we

156   determined the extent to which individual taxa could influence the community structure of

157   merged independent studies. Here community structure is defined by the presence and relative

158   abundances of individual taxa, along with co-occurrence relationships between those taxa. This

159   was done in two ways: first, we constructed a model that classified the study from which a

160   sample came based on the proportions of the 8,287 taxa it contained (1.5% [± 0.02% CI]

161   classification error, by internal cross-validation). Second, we determined the contribution of each

162   taxon to bacterial community structure by quantifying its importance in a model that separated

163   the observed data from synthetic data randomly drawn from the observed distributions of relative

164   abundances for each taxon [44,45] (*see Methods*).

165

166   Merging of disparate microbial sequence data is known to be plagued with potential biases

167   including: lack of standardization of sample collection, methodological issues regarding DNA

168   extraction and primer choice, incomplete metadata, the technical biases of different sequencing

169   platforms, sequencing depth, PCR Bias, different clustering methods, and the use of different

170   taxonomic classification pipelines[46–52]. We therefore took the novel step to quantify the

171   importance of both technical and environmental factors alongside taxa in the Random Forests

172    models (Figure 2). Of note, 'owner', which encompasses the technical biases and uniqueness of a

173    given dataset, is very effective for differentiating between studies (i.e. the owner is far to the

174    right in Figure 2) yet is entirely uninformative about community structure (i.e. owner is at the far

175    bottom in Figure 2). In fact, *all* technical factors included are better than 98.5% of all taxa to

176    differentiate between studies, indicating that the observed differences among studies in taxon

177    relative abundances are strongly confounded with technical factors. Independent of taxonomy,

178    certain environmental factors, such as country of origin, latitude and longitude, and soil pH, were

179    highly important in differentiating studies but not in determining community structure. By

180    contrast, minimum soil sampling depth was not very important in separating studies, and was

181    more associated with community structure. It is well known that bacterial diversity decreases

182    with soil depth[53] and our results show that in a global assessment, soil depth remains a strong

183    predictor of bacterial community composition. Perhaps most useful for future research, this result

184    highlights that not all environmental factors are equally confounded by technical factors, and

185    shows that by combining data from across many independent studies we may identify previously

186    overlooked taxa and factors relevant for structuring communities.

187

188    **Importance for structuring soil bacterial communities**

189    Although all studies were confounded by technical and environmental covariates, there remained

190    many taxa that were non-randomly distributed and were not confounded with technical

191    differences among studies (upper left in Figure 2). When assessing the role of these different taxa

192    in structuring the community, we found a trade-off between taxon abundance and importance in

193    community structure, such that low abundance taxa are disproportionately important in the non-

194    random structure of communities, where the most important taxa are rarer than expected

8

195     compared to the randomly permuted data (Figure 3). Thus, the importance of taxa for

196     determining community structure is negatively correlated with the average abundance of those

197     taxa, whereas taxon abundance is positively correlated with importance for separating studies ($\rho$

198     $= -0.79$ and $\rho = +0.51$ respectively, rank correlation, cf. null expectations of $\rho = -0.62$ and $-0.12$

199     respectively in permuted data). The taxa most closely associated with differences between

200     studies tend to be those present at or greater than 0.1% relative abundance, but those most

201     important in determining community structure tend to be present at 0.0001% abundance or less

202     (with a null expectation of around 0.01-0.001% in each case, Figure 3). This result is only found

203     by considering the full set of studies and is neither apparent within single studies (Supplementary

204     Fig. 4A-B) nor a subset of studies (whether matched by name or sequence Supplementary Fig.

205     5). It corresponds to the long tail in frequency-abundance distributions of soil microbial

206     communities[54], where many taxa in the soil are known to occur at low abundance. Thus, if rarer

207     taxa tend to be more important for distinguishing between communities, it is within this long tail

208     that we might identify taxa that could indicate ecological or functional differences among soil

209     communities[33,55,56].

210

211     To be ecological indicators[57,58], taxa need to vary in abundance in response to environmental

212     factors and have high occurrence across studies, as is the case for the phylum Acidobacteria[43].

213     Acidobacteria, however, are typically abundant and our analysis suggests that the most abundant

214     taxa are *not* the most important in determining community structure. While dominant taxa like

215     Acidobacteria do change with environmental factors such as pH (Figure 1D), those changes are

216     of lesser importance for the 'non-randomness' of community structure, and more confounded

217     with technical effects, than changes in less dominant, pH responsive taxa (Supplementary Figure

218  3A). Therefore, we assessed which taxonomic ranks are more or less distinguished from the

219  randomly permutated data. Although differences among domains and phyla are strongly

220  associated with differences among studies (Figure 4B) only taxa at a rank lower than phyla are

221  consistently better than random at identifying community structure (Figure 4A).

222

223  A very similar pattern was found for the sequence-matched community, emphasizing the

224  importance of taxa at the level of Class and below (Supplementary Figure 7A and 7B). However,

225  this was not apparent in individual studies (Supplementary Figure 4C-D), where phyla were

226  relatively important. A subset of the taxonomy-matched studies showed a pattern intermediate

227  between the single studies and the full dataset (phyla with some importance, but less than Class,

228  Order or Family, Supplementary Figure 7C). This, along with abundance analyses (Figure 3 and

229  Supplementary Figure 5), suggests that our name matching approach is consistent with, but less

230  powerful than a full sequence-matched analysis. At the same time, the taxonomy-matching is

231  worthwhile because, as with the findings on abundance (Figure 3), macroecological patterns (the

232  importance of taxa below phyla and of relatively low abundance in community structure) are

233  evident when we consider thousands of samples from tens of studies, that are not apparent from

234  hundreds of samples from one or a handful of studies.

235

236  To be a good ecological indicator a taxon should occur in most studies; we therefore looked

237  explicitly at the relationship between a taxon's importance in community structure and its

238  occurrence across studies. Low abundance taxa and taxa of lower taxonomic rank are

239  consistently important in determining community structure, but tend to be detected in fewer

240  studies ($\rho = 0.59$ and 0.31 respectively Supplementary Figure 3B and 3C). We discovered a

241 novel relationship between taxon occurrence across studies and importance for structuring

242 communities for all taxa (Figure 5, Supplementary Table 4). Comparison with the null

243 expectation reveals a range of taxa, occurring in multiple samples from most studies, which are

244 much more important in determining community structure than expected by chance. A similar

245 pattern is apparent in the sequence-matched dataset (Supplementary Figure 8A) and the same

246 subset of studies when taxonomy-matched (Supplementary Figure 8B). Altogether, the analyses

247 clearly illustrate the significance of taxonomic rank, for example *class* Gemmatimonadetes is

248 relatively unimportant for community structure but *genus* Gemmatimonadetes is relatively

249 important. The result also shows rarer taxa being more important in structuring communities and

250 suggests rarer bacterial taxa play overlooked ecologically important roles for bacterial

251 community dynamics[56]. This result is robust to artifacts caused by the rarest taxa (e.g.

252 differences between 0 and 1 reads in a sample could be significant for a model, without being

253 biologically significant) – a very similar pattern is seen when only taxa present at above 0.003%

254 in any given sample were included in this analysis (typically removing the rarest 10% of taxa

255 from any given sample, Supplementary Figure 9). Conversely, many taxa of high taxonomic rank

256 with high occurrence across samples, such as the phyla Actinobacteria, Acidobacteria,

257 Proteobacteria, and Bacteroidetes, were much less important for community structure than the

258 null expectation. These taxa have been reported elsewhere as 'core' members of the soil

259 community[43,59,60], and even been included in source-tracking of microbial communities due to

260 their ubiquitous presence in soil[61]. Yet, it is the consistent presence of the core taxa across

261 samples and studies that makes them inadequate for assessing community structure.

262

263 **Conclusions**

11

264  Our results demonstrate the power of combining global bacterial HTS data from multiple

265  independent sources for the detection of biogeographical patterns and for identifying community

266  patterns that can be used to generate hypotheses on the roles of certain taxa. Though our

267  assessment was on soil communities, our methods can be applied to broadly to other microbial

268  datasets and disciplines. Taxonomy-based merging gives results that are consistent with raw

269  sequence data, and expands opportunities for extracting information about microbial

270  communities from the wealth of existing and future studies. Moreover, we find that rarer

271  bacterial taxa are more important in differentiating communities than previously assumed, and

272  hold potential as overlooked soil indicators or keystone species. Still, there are considerable

273  challenges associated with merging large sequence datasets beyond the well-known biases that

274  accompany any molecular HTS study. Perhaps the most concerning was that so few raw

275  sequence datasets for publically deposited analyses could be retrieved. This highlights the need

276  for wider community adoption of open and accessible short read sequence databases[62], open

277  reference clustering[63], standardized databases[64] and—as always—that metadata should be

278  consistent and accessible. Regardless of these challenges, as HTS methods rapidly advance we

279  must find ways to simultaneously curate and carry our research knowledge forward[32]. Only then,

280  in combination with the many novel and classical approaches, can we uncover the full breadth of

281  soil diversity and the roles soil microbes play for ecosystem processes.

282

283   **Methods:**

284   *Description of datasets*:

285   Metadata from the 30 studies and 1998 samples were collected and compiled into a summary

286   data file. To do so, we standardized the metadata of each study using the dplyr package

287   (Wickham & Francois, 2016) of the R statistical platform (R Core Team, 2016). Samples were

288   collected from 21 counties representing all continents except Antarctica. In addition to location

289   and pH data (median = 6.1, quartile range=5.3-7.0), which were available from all studies,

290   information on altitude (10 m, 10-860 m), soil moisture (19.5%, 14.1-27.4%), and total soil

291   nitrogen (0.36 mg kg$^{-1}$, 0.23-0.51 mg kg$^{-1}$), carbon (4.7%, 1.9-7.5%) and phosphorus (20.7 mg

292   kg$^{-1}$, 7.0-223.0 mg kg$^{-1}$) was noted where available. Depth of sample collection was also noted

293   and ranged from surface collections to a maximum depth of 70 cm, with 83% of samples

294   originating from 0-10 cm below the soil surface. Samples represented anthropogenically

295   managed (59%) and natural (40%; remaining samples undefined) systems, and were taken from

296   arable, grassland, peatland, forest, scrub (including tundra) and urban habitats. The majority of

297   samples (71%) were described as non-experimental, meaning no treatments were applied, with

298   the remainder described as experimental. Sequencing data were either produced using Roche 454

299   technology (22%) or one of the Illumina platforms (78%). Primer pairs were defined for 92% of

300   the samples and nine different pairs were identified from the study meta data (27F:338R;

301   341F:518R; 341F:806R; 341F:907R; 357F:926R; 515F:806R; 577F:926R; 799F:1193R and

302   341F:805R) with the majority of samples (66%) using 515F and 806R to produce amplicons.

303   Post sequencing processing varied, but 81% of samples were run through the QIIME workflow

304   at some point. An OTU table for 1 study comprising 43 samples was programmatically retrieved

305   from the MG-RAST public metagenome repository[65]. Taxonomy for the different studies was

306     mainly assigned using the Greengenes database (84 %), but RDP (6 %;[46] and the Silva database

307     (9 %)[66] were also used.

308

309     *Merging OTU tables:*

310     For the OTU tables from the 30 individual studies to be merged, extensive data cleaning was

311     carried out on the OTU and taxonomy files to maximize the possibility of matching taxa across

312     datasets. This comprised several steps: (1) Most datasets contained a seven-level taxonomy,

313     recorded in a variety of ways, which was converted to a standardized format. (2) Individual

314     taxon names were cleaned, to give a single name at each taxonomic level (e.g. removing special

315     characters and extra annotations, such as 'candidate division' or details of containing taxa). (3)

316     For the many cases where a taxon was not assigned at a particular taxonomic level, a unified

317     'unassigned' label was created. Repeating analyses with all these taxa removed made no

318     qualitative difference to the results (Supplementary Figure 10). Merging at the taxonomy-based

319     level has the added benefit of lessening the impacts of hypervariable regions. For example, the

320     identification of an organism at a specific level in one sample also contributes to the

321     identification of the containing genus for that sample, allowing direct comparison with a sample

322     where, because a different region was sequenced, that same organism is only resolved to the

323     genus level. Next, relative abundance data were, where necessary, re-scaled to sum to 1 for a

324     sample, using original OTU count files where possible. These values were then manipulated to

325     give data tables usable for modeling using custom R scripts. For some analyses (Figures 3-5), a

326     dataset without community structure was created by randomly permuting the relative abundance

327     of each taxon across all samples. Unless otherwise stated, the analyses performed on the

328     permuted dataset was identical to that performed on the observed data.

329

330 *Merging raw sequence data and other validation datasets:*

331 While no dataset can currently provide a "ground truth" against which to judge our approach, we

332 can at least validate it. The primary validation of our taxonomy-matching approach was to merge

333 raw sequence data ('sequence-matched') from five studies (Supplementary Table 1). Per sample

334 fastq files were obtained for each individual dataset. Read files were quality filtered with sickle

335 [67] for single end reads trimming bases below phred score 36 and shorter than 100bp. These

336 stringent filtering criteria were applied to keep only high quality reads and to make sure it is

337 possible to map reads to full length 16S rRNA gene sequences. Full length 16S rRNA gene

338 sequences from the Silva 119 release [66] were obtained in Qiime compatible format from the Silva

339 Download Archive For each dataset, all reads were mapped to the full length 16S rRNA gene

340 sequences using the usearch global algorithm implemented in VSEARCH version 1.9.6 [68]. The

341 alignment results in usearch table format (uc) were directly converted to BIOM format using

342 biom version 2.1.5 [69]. Consensus/majority taxonomy was added as metadata to the biom file.

343 Finally, all BIOM files of each dataset were merged using Qiime version 1.9.1 [70]. All steps were

344 implemented in a workflow made with Snakemake version 3.5.4 [71] available: (De Hollander

345 2016). See Supplementary Fig 1 for workflow.

346

347 To use this sequence-matched dataset to validate our taxonomy-matching approach across

348 studies using different taxonomy databases (Supplementary Figures 5, 7 & 8) we created an

349 equivalent taxonomy-matched dataset from the same 5 studies. As with the full dataset, only taxa

350 occurring in at least two studies were included in either this or the sequence-matched dataset. To

351 test what is gained or lost by considering different numbers of studies simultaneously, we

352 considered, not only the full dataset (30 studies) and the subset of 5 studies used in the sequence-

353 matched dataset, but two of the largest individual studies: from Central Park, NYC

354 encompassing 594 samples (study #24) and a global dataset encompassing 103 samples (study

355 #30). In each case a simple subset of the full dataset was analyzed (Supplementary Figure 4). To

356 address PCR biases (Supplementary Table 2) and biases associated with rare taxa, we created a

357 filtered subset of the data where only taxa present at above 0.003% in any given sample were

358 considered, meaning that all taxa deemed present are represented by multiple sequence reads

359 (Supplementary Figure 9). To address the issue of differential 16S copy numbers skewing

360 abundance estimates, we created a binary dataset of the presence/absence of all taxa. The results

361 for a model separating studies using this dataset were very similar to the main dataset using

362 relative abundance, however, there was insufficient power to identify taxa important for

363 community structure (Supplementary Figure 6). Nonetheless, this analysis did agree with the

364 main analysis that phyla were the most stable taxonomic level, with lower importance than on

365 the permuted data (Supplementary Figure 6). Finally, to test the effect of 'unknown' or

366 unclassified bacterial taxa we created a reduced dataset where all taxa classified as 'unassigned'

367 at any level were removed (Supplementary Figure 10).

368

369 *Random forest models.*

370 To test for the importance of different taxa in the structuring of the data we used Random Forest

371 models [45,72] with the relative abundances of the taxa as explanatory variables. Random Forest

372 models have two principal advantages in this context: 1) they can deal easily with thousands of

373 explanatory variables and quantify their relative importance, and 2) they can run equivalently in

374 both supervised and un-supervised modes. In the latter, the importance of a variable describes

375    how effective it is at separating the observed data from randomized synthetic data[44]. In both

376    cases, a proximity matrix may be generated, which can be used for ordination (Supplementary

377    Figure 2). The importance of individual taxa in a Random Forest relate to traditional ecological

378    measures. For instance, the importance in a supervised model, such as that used separating

379    studies (x-axis in Figure 2) is closely correlated with the sensitivity component of the indicator

380    value of each taxon ($\rho = 0.89$, Supplementary Figure 3D)[58]. There are two key parameters that

381    may be adjusted in a Random Forest model, *mtry*, the number of variables randomly sampled as

382    candidates for a split in the constituent trees and *ntree*, the number of trees in the forest. *mtry* was

383    set at its default value (square root of the number of variables) *ntree* was set to 100,000 for each

384    forest. Such a large number of trees was found to be necessary to achieve stable importance

385    across taxa and was achieved by combining several forests run in parallel without normalizing

386    votes. Other parameters were left at default values, in particular, trees were grown to completion

387    (i.e. a minimum node size of 1). The un-scaled permutation importance of variables is used

388    throughout: Each variable importance is the difference between the classification error rate of a

389    tree on data not used to construct it (the 'out of bag' data) and the same error following random

390    permutation of the variable in question, averaged over all trees.

391

392    We used permuted data (see above) to create null distributions for taxon importance. For

393    unsupervised Random Forests analyses, such as the community structure model, this amounts to

394    calculating how important a taxon with a particular abundance distribution is for separating two

395    randomized distributions. This can then be compared to its importance for separating the

396    observed from a randomized distribution. This clarifies the fact that, even in null data without

397    community structure (Supplementary Figure 2), variable importance correlates with ecologically

398  important factors, such as abundance. This makes intuitive sense in as much as, even with

399  randomized samples, is easier to separate them on the basis of taxa that occur in only some of

400  them than on the basis of ubiquitous taxa. This, for instance, results in the negative slope of the

401  orange (permuted, null, data) line in Figure 5.

402

403  All analyses were completed with RandomForest package for R version 4.6.

404

417

418  **References**

419  1.    Prosser, J. I. Dispersing misconceptions and identifying opportunities for the use of

420        'omics' in soil microbial ecology. *Nat. Rev. Microbiol.* **13,** 439–46 (2015).

18

421  2.  Bardgett, R. D. & van der Putten, W. H. Belowground biodiversity and ecosystem

422       functioning. *Nature* **515,** 505–511 (2014).

423  3.  Caporaso, J. G. *et al.* Ultra-high-throughput microbial community analysis on the Illumina

424       HiSeq and MiSeq platforms. *ISME J.* **6,** 1621–4 (2012).

425  4.  Tedersoo, L. *et al.* Global diversity and geography of soil fungi. *Science.* **346,** (2014).

426  5.  Bik, H. M. *et al.* Sequencing our way towards understanding global eukaryotic

427       biodiversity. *Trends Ecol. Evol.* **27,** 233–243 (2012).

428  6.  Davison, J. *et al.* Global assessment of arbuscular mycorrhizal fungus diversity reveals

429       very low endemism. *Science.* **349,** (2015).

430  7.  Wieder, W. R., Bonan, G. B. & Allison, S. D. Global soil carbon projections are improved

431       by modelling microbial processes. *Nat. Clim. Chang.* **3,** 909–912 (2013).

432  8.  Allison, S. D., Wallenstein, M. D. & Bradford, M. A. Soil-carbon response to warming

433       dependent on microbial physiology. *Nat. Geosci.* **3,** 336–340 (2010).

434  9.  Karhu, K. *et al.* Temperature sensitivity of soil respiration rates enhanced by microbial

435       community response. *Nature* **513,** 81–84 (2014).

436  10.  Wieder, W. R., Cleveland, C. C., Smith, W. K. & Todd-Brown, K. Future productivity

437       and carbon storage limited by terrestrial nutrient availability. *Nat. Geosci.* **8,** 441–444

438       (2015).

439  11.  Barberán, A., Casamayor, E. O. & Fierer, N. The microbial contribution to macroecology.

440       *Front. Microbiol.* **5,** 203 (2014).

441  12.  Ramirez, K. S. *et al.* Biogeographic patterns in below-ground diversity in New York

442       City's Central Park are similar to those observed globally. *Proc. Biol. Sci.* **281,** 20141988-

443       (2014).

19

444    13.    O'Brien, S. L. *et al.* Spatial scale drives patterns in soil bacterial diversity. *Environ.*

445            *Microbiol.* **18,** 2039–2051 (2016).

446    14.    Evans, S., Martiny, J. B. H. & Allison, S. D. Effects of dispersal and selection on

447            stochastic assembly in microbial communities. *ISME J.* **11,** 176–185 (2017).

448    15.    Talbot, J. M. *et al.* Endemism and functional convergence across the North American soil

449            mycobiome. *Proc. Natl. Acad. Sci. U. S. A.* **111,** 6341–6 (2014).

450    16.    Barber, A. *et al.* Why are some microbes more ubiquitous than others? Predicting the

451            habitat breadth of soil bacteria. *Ecol. Lett.* **17,** 794–802 (2014).

452    17.    Fierer, N. *et al.* Reconstructing the microbial diversity and function of pre-agricultural

453            tallgrass prairie soils in the United States. *Science* **342,** 621–4 (2013).

454    18.    Ma, B. *et al.* Geographic patterns of co-occurrence network topological features for soil

455            microbiota at continental scale in eastern China. *ISME J.* **10**, 1891-1901 (2016).

456    19.    Ranjard, L. *et al.* Turnover of soil bacterial diversity driven by wide-scale environmental

457            heterogeneity. *Nat. Commun.* **4,** 1434 (2013).

458    20.    Jetz, W., McPherson, J. M. & Guralnick, R. P. Integrating biodiversity distribution

459            knowledge: toward a global map of life. *Trends Ecol. Evol.* **27,** 151–9 (2012).

460    21.    Ricketts, T. H. *et al.* Disaggregating the evidence linking biodiversity and ecosystem

461            services. *Nat. Commun.* **7,** 13106 (2016).

462    22.    Dirzo, R. *et al.* Defaunation in the Anthropocene. *Science..* **345,** 401–406 (2014).

463    23.    Gerstner, K., Dormann, C. F., Stein, A., Manceur, A. M. & Seppelt, R. Effects of land use

464            on plant diversity - A global meta-analysis. *J. Appl. Ecol.* **51,** 1690–1700 (2014).

465    24.    Patterson, D. J., Cooper, J., Kirk, P. M., Pyle, R. L. & Remsen, D. P. Names are key to the

466            big new biology. *Trends Ecol. Evol.* **25,** 686–691 (2010).

467  25.  Santos, A. M. & Branco, M. The quality of taxonomy-based species records in databases.

468       *Trends Ecol. Evol.* **27,** 6-7-8 (2012).

469  26.  Beiko, R. G. Microbial Malaise: How Can We Classify the Microbiome? *Trends*

470       *Microbiol.* **23,** 671–679 (2015).

471  27.  Tedersoo, L. *et al.* Standardizing metadata and taxonomic identification in metabarcoding

472       studies. *Gigascience* **4,** 34 (2015).

473  28.  Ramirez, K. S. *et al.* Toward a global platform for linking soil biodiversity data. *Front.*

474       *Ecol. Evol.* **3,** (2015).

475  29.  Turner, W. *et al.* Free and open-access satellite data are key to biodiversity conservation.

476       *Biol. Conserv.* **182,** 173–176 (2015).

477  30.  Gilbert, J. A., Jansson, J. K. & Knight, R. The Earth Microbiome project: successes and

478       aspirations. *BMC Biol.* **12,** 69 (2014).

479  31.  Joppa, L. N. *et al.* Filling in biodiversity threat gaps. *Science.* **352,** (2016).

480  32.  Sinha, R., Abnet, C. C., White, O., Knight, R. & Huttenhower, C. The microbiome quality

481       control project: baseline study design and future directions. *Genome Biol.* **16,** 276 (2015).

482  33.  Sogin, M. L. *et al.* Microbial diversity in the deep sea and the underexplored "rare

483       biosphere". *Proc. Natl. Acad. Sci. U. S. A.* **103,** 12115–20 (2006).

484  34.  García-Palacios, P. *et al.* Are there links between responses of soil microbes and

485       ecosystem functioning to elevated CO2, N deposition and warming? A global perspective.

486       *Glob. Chang. Biol.* **21,** 1590–1600 (2015).

487  35.  Hermans, S. M. *et al.* Bacteria as Emerging Indicators of Soil Condition. *Appl. Environ.*

488       *Microbiol.* **83,** AEM.02826-16 (2017).

489  36.  Philippot, L. *et al.* The ecological coherence of high bacterial taxonomic ranks. *Nat. Rev.*

490        *Microbiol.* **8,** 523–529 (2010).

491    37.    Shade, A., Caporaso, J. G., Handelsman, J., Knight, R. & Fierer, N. A meta-analysis of

492        changes in bacterial and archaeal communities with time. *ISME J.* **7,** 1493–506 (2013).

493    38.    Delgado-Baquerizo, M., Grinyer, J., Reich, P. B. & Singh, B. K. Relative importance of

494        soil properties and microbial community for soil functionality: insights from a microbial

495        swap experiment. *Functional Ecology*. **30**, 1862-1873 (2016).

496    39.    Hendershot, J. N., Read, Q. D., Henning, J. A., Sanders, N. J. & Classen, A. T.

497        Consistently inconsistent drivers of microbial diversity and abundance at macroecological

498        scales. *Ecology* **98,** 1757–1763 (2017).

499    40.    Bier, R. L. *et al.* Linking microbial community structure and microbial processes: an

500        empirical and conceptual overview. *FEMS Microbiol. Ecol.* **91,** (2015).

501    41.    Walters, W. A., Xu, Z. & Knight, R. Meta-analyses of human gut microbes associated

502        with obesity and IBD. *FEBS Lett.* **588,** 4223–4233 (2014).

503    42.    Kõljalg, U. *et al.* Towards a unified paradigm for sequence-based identification of fungi.

504        *Mol. Ecol.* **22,** 5271–5277 (2013).

505    43.    Lauber, C. L., Hamady, M., Knight, R. & Fierer, N. Pyrosequencing-Based Assessment of

506        Soil pH as a Predictor of Soil Bacterial Community Structure at the Continental Scale.

507        *Appl. Environ. Microbiol.* **75,** 5111–5120 (2009).

508    44.    Shi, T. & Horvath, S. Unsupervised Learning With Random Forest Predictors. *J. Comput.*

509        *Graph. Stat.* **15,** 118–138 (2006).

510    45.    Breiman, L. & Cutler, A. Random Forests Manual v4.0. *Technical report, UC Berkeley*

511        (2003).

512    46.    McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological

513    and evolutionary analyses of bacteria and archaea. *ISME J.* **6,** 610–8 (2012).

514    47.    Cole, J. R. *et al.* Ribosomal Database Project: data and tools for high throughput rRNA

515    analysis. *Nucleic Acids Res.* **42,** D633-42 (2014).

516    48.    Werner, J. J., Zhou, D., Caporaso, J. G., Knight, R. & Angenent, L. T. Comparison of

517    Illumina paired-end and single-direction sequencing for microbial 16S rRNA gene

518    amplicon surveys. *ISME J.* **6,** 1273–1276 (2011).

519    49.    Lozupone, C. & Stombaugh, J. Meta-analyses of studies of the human microbiota.

520    *Genome Res.* **10**, 1704-14 (2013).

521    50.    Pawluczyk, M. *et al.* Quantitative evaluation of bias in PCR amplification and next-

522    generation sequencing derived from metabarcoding samples. *Anal. Bioanal. Chem.* **407,**

523    1841–1848 (2015).

524    51.    Smith, D. P., Peay, K. G., Palmer, M., Gillikin, C. & Keefe, D. Sequence Depth, Not PCR

525    Replication, Improves Ecological Inference from Next Generation DNA Sequencing.

526    *PLoS One* **9,** e90234 (2014).

527    52.    Zhbannikov, I. Y. & Foster, J. A. MetAmp: combining amplicon data from multiple

528    markers for OTU analysis. *Bioinformatics* **31,** 1830–1832 (2015).

529    53.    Lu, X., Seuradge, B. J. & Neufeld, J. D. Biogeography of soil Thaumarchaeota in relation

530    to soil depth and land usage. *FEMS Microbiol. Ecol.* **93,** (2017).

531    54.    Jung, S. P. & Kang, H. Assessment of microbial diversity bias associated with soil

532    heterogeneity and sequencing resolution in pyrosequencing analyses. *J. Microbiol.* **52,**

533    574–580 (2014).

534    55.    Langille, M., Zaneveld, J. & Caporaso, J. Predictive functional profiling of microbial

535    communities using 16S rRNA marker gene sequences. *Nature* (2013).

536   56.   Jousset, A. *et al.* Where less may be more: how the rare biosphere pulls ecosystems

537        strings. *ISME J.* **11,** 853–862 (2017).

538   57.   Hermans, S. M. *et al.* Bacteria as emerging indicators of soil condition. *Appl. Environ.*

539        *Microbiol.* AEM.02826-16 (2016).

540   58.   Cáceres, M. De & Legendre, P. Associations between species and groups of sites: indices

541        and statistical inference. *Ecology* **90,** 3566–3574 (2009).

542   59.   Maestre, F. T. *et al.* Increasing aridity reduces soil microbial diversity and abundance in

543        global drylands. *Proc. Natl. Acad. Sci. U. S. A.* **112,** 15684–9 (2015).

544   60.   Fierer, N., Bradford, M. a & Jackson, R. B. Toward an ecological classification of soil

545        bacteria. *Ecology* **88,** 1354–64 (2007).

546   61.   Knights, D. *et al.* Bayesian community-wide culture-independent microbial source

547        tracking. *Nat. Methods* **8,** 761–763 (2011).

548   62.   Muir, P. *et al.* The real cost of sequencing: scaling computation to keep pace with data

549        generation. *Genome Biol.* **17,** 53 (2016).

550   63.   Rideout, J. R. *et al.* Subsampled open-reference clustering creates consistent,

551        comprehensive OTU definitions and scales to billions of sequences. *PeerJ* **2,** e545 (2014).

552   64.   Yilmaz, P. *et al.* The genomic standards consortium: bringing standards to life for

553        microbial ecology. *ISME J.* **5,** 1565–7 (2011).

554   65.   Wilke, A. *et al.* The MG-RAST metagenomics database and portal in 2015. *Nucleic Acids*

555        *Res.* **44,** D590–D594 (2016).

556   66.   Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data

557        processing and web-based tools. *Nucleic Acids Res.* **41,** D590-6 (2013).

558   67.   Joshi & Fass, J. Sickle: A sliding-window, adaptive, quality-based trimming tool for

559        FastQ files. (2011).

560   68.    Rognes, T. *et al.* vsearch: VSEARCH 1.9.6. (2016).

561   69.    McDonald, D. *et al.* The Biological Observation Matrix (BIOM) format or: how I learned

562        to stop worrying and love the ome-ome. *Gigascience* **1,** 7 (2012).

563   70.    Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing

564        data. *Nat Meth* **7,** 335–336 (2010).

565   71.    Koster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine.

566        *Bioinformatics* **28,** 2520–2522 (2012).

567   72.    Breiman, L. Random Forests. *Mach. Learn.* **45,** 5–32 (2001).

568

569   Hadley Wickham and Romain Francois (2016). dplyr: A Grammar of Data

570    Manipulation. R package version 0.5.0.

571    https://CRAN.R-project.org/package=dplyr

572

573   R Core Team (2016). R: A language and environment for statistical computing.

574    R Foundation for Statistical Computing, Vienna, Austria. URL

575    https://www.R-project.org/.

576

577

**Figure 1. Merging of data from 32 independent studies demonstrates wide geographic breadth, community variation, and confirms the well-known importance of soil pH. A.** Map of locations from which samples were collected, with zoom panels on the United States (left) and western Europe (right). Points in blue were used in both the taxonomy-based and raw-unified analyses and red points were only used in taxonomy-based analyses. **B.** Average proportion of total prokaryotic abundance and **C.** eukaryotic abundance, represented by taxa shared among different numbers of datasets at different taxonomic levels. Level 1 indicates the complete data, levels 2-4 are subsets of the data containing only taxa present in a minimum of 2-4 separate datasets. **D.** Correlation plot of Acidobacteria relative abundance to soil pH where ach color represents a different study ($r = $ -0.42 $p$=8.6 x $10^{-87}$).

**Figure 2: Regardless of technical differences between studies, many bacterial taxa are still informative about bacterial community structure.** Machine learning models classify the study from which samples came (x-axis) based on the relative abundance of taxa within samples and distinguish the observed distribution of taxa among samples from random (y-axis). Plotted alongside bacterial taxa (black) are technical factors (red) and ecological factors (purple), including soil pH, minimum and maximum soil depth, longitude, latitude and degrees from the equator. All values are variable importance from Random Forest models (see *Methods*) – points further to the right on the x-axis have more importance in separating studies, while points higher up on the y-axis, have more importance for community structure. Note the non-linear axes.

601 **Figure 3: Rarer taxa are more important for structuring communities than abundant taxa.**

602 Here we show the thousand most important bacterial taxa in community structure (A) and in

603 separating studies (B) with respect to their average relative abundance across samples. Plotted

604 are the 'observed' points (green) and 'permuted' points (orange) which are a null distribution

605 from performing the same analysis on a permuted dataset (see *Methods*). The y-axis reports the

606 rank variable importance in the Random Forests model of community structure (see *Methods*),

607 i.e. the taxon with the greatest importance in this model is ranked 1, the second greatest 2, etc.

608

609 **Figure 4: The importance of bacterial taxa classified at different taxonomic ranks.** Lower

610 taxonomic rank is more important for community structure (A), while high taxonomic rank is

611 more important for separating studies (B). For each taxon, the difference was calculated between

612 the variable importance (see *Methods*) of that taxon in a Random Forests model of either

613 community structure or separating studies and the equivalent value from an analysis performed

614 on the permuted dataset (see *Methods*). The lines and grey ribbons show the mean and standard

615 error respectively of these values across taxa at each taxonomic r considered.

616

617 **Figure 5: Importance of bacterial taxa in community structure related to their occurrence**

618 **in different studies.** The y-axis reports the variable importance in the Random Forests model of

619 community structure (see *Methods*). Green 'observed' points correspond to those taxa shown in

620 Figure 1. Orange 'permuted' points correspond to the same analysis on a null distribution (see

621 *Methods*). Lines are general additive model (gam) smoothers. Each line is shown with a

622 confidence interval (grey); where this is not visible it is narrower than the line it surrounds.

623

a

b

c

d

e

Data available
• Names only
• Raw sequence

Number of samples
• 1
• 10
• 100

Latitude

Longitude

Median proportion of abundance

Bacteria

Fungi

Minimum number of studies
in which a taxon occurs
1
2
3
4

Domain  Kingdom  Phylum  Class  Order  Family  Genus  Species

Taxonomic level

Proportion of Acidobacteria

pH