

1 **Title: Obscure soil microbes and where to find them**

2 **Running title: Unclassified bacteria and fungi globally**

3 Manuel Delgado-Baquerizo

4 Departamento de Biología y Geología, Física y Química Inorgánica, Escuela Superior de Ciencias Experi-
5 mentales y Tecnología, Universidad Rey Juan Carlos, Calle Tulipán Sin Número, Móstoles 28933, Spain.

6 *Author for correspondence:

7 Manuel Delgado-Baquerizo. M.delgadobaquerizo@gmail.com. Departamento de Biología y Geología,
8 Física y Química Inorgánica, Escuela Superior de Ciencias Experimentales y Tecnología, Universidad
9 Rey Juan Carlos, Calle Tulipán Sin Número, Móstoles 28933, Spain.

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34 **Abstract**

35 Many soil bacteria and fungi remain unclassified at the highest taxonomic ranks (e.g., phyla level), which
36 hampers our ability to assess the ecology and functional capabilities of these soil organisms in terrestrial
37 ecosystems globally. The first logical step toward the classification of these unknown soil taxa is to identify
38 potential locations on Earth where these unclassified bacteria and fungi are feasibly most prevalent. To do
39 this, here I used data from a global soil survey across 235 locations, including amplicon sequencing infor-
40 mation for fungal and bacterial communities, and generated global atlases highlighting those soils where
41 the percentages of taxa of bacteria and fungi with an unknown phyla are expected to be more prevalent.
42 Results indicate that soil samples with the largest percentage of fungi with an unknown phyla can be found
43 in dry forests and grasslands, while those with the largest percentage of bacteria with an unknown phyla
44 are found in boreal and tropical forests. This information can be used by taxonomists and microbiologists
45 to target these potentially new soil taxa.

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69 **Text**

70 Soil microbial communities play an essential role in maintaining important soil processes such as nutrient
71 cycling, waste decomposition, climate regulation, and pollution degradation (Bardgett and van der Putten
72 2014; Delgado-Baquerizo et al. 2016). Today, sequencing technologies are well established and broadly
73 used (Caporaso et al. 2010). As such, producing large amounts of data on the composition and diversity of
74 bacterial and fungal communities is no longer so challenging. Moreover, the major ecological drivers of
75 the variation in these microbial communities are becoming increasingly visible (Tedersoo et al. 2014; Fierer
76 2017). The spotlight is now on the soil taxonomists. Although progress has been made in the past few years
77 (Marx 2017; York 2018), culturing, isolating, and classifying soil microbes is still a difficult task. For most
78 soil bacterial and fungal species, we know very little about their identity or the tasks performed even by the
79 most dominant microbial taxa (Delgado-Baquerizo et al. 2018). More concerning, in some cases, we lack
80 the most basic taxonomic information to classify these bacterial and fungal taxa as they do not match the
81 latest data within taxonomic databases (e.g., Greengenes and UNITE) even at the highest taxonomic ranks
82 (e.g., phyla level).

83 The first logical step toward the classification of these unknown microbial taxa is to identify potential lo-
84 cations where they could be found across the globe. This information can then be used by taxonomists and
85 microbiologists to target these new soil taxa. Here, I used data from a global soil survey (Delgado-Baquerizo
86 et al. 2018) across 235 locations (Fig. S1), and including amplicon sequencing information on fungal (ITS
87 gene) and bacterial (16S rRNA gene) communities from around the world, to highlight those locations on
88 Earth where taxa of bacteria and fungi with an unknown phyla are feasibly most prevalent. The database in
89 Delgado-Baquerizo et al. (2018) has been used previously to identify the dominant taxa of bacteria globally,
90 and more recently, the major ecological predictors of bacterial diversity (Delgado-Baquerizo and Eldridge
91 2019). I used the bioinformatics pipeline described in Delgado-Baquerizo et al. (2018), and two of the most
92 commonly used microbial databases for taxonomic identification (Greengenes and UNITE), to estimate, at
93 the global scale, the percentage of phylotypes of bacteria and fungi with an unknown phyla in soils across
94 the globe. These taxa are classified as fungi or bacteria using taxonomic databases, but do not match any
95 known phyla. As such, they are expected to be potential new phyla of fungi or bacteria.

96 As expected, the taxonomic information at the species level could not be found for 99% of bacterial and
97 63% of fungal phylotypes (clustered at 97% similarity). Notably, up to 1.36% and 9.37% of the retrieved
98 phylotypes classified as bacteria or fungi remained unclassified at the phyla level in soils across the globe.
99 For these microbes, we do not know the phylum to which they belong. In other words, for some soils,
100 almost 10% of taxa within bacteria and fungi are totally unknown to us. These taxa represent between 0.01-
101 1.86% (average of 0.12%) of all 16sRNA sequences, and between 0.00-22.11% (average of 3.98%) of all
102 ITS retrieved sequences. On average, soil samples with the largest percentage of phylotypes of bacteria
103 with an unknown phyla can be found in boreal and tropical forests (Fig. 1), while those with the largest
104 percentage of phylotypes of fungi with an unknown phyla are found in dry forests and grasslands (Fig. 1).

105 I then generated a global atlas highlighting those global soils where bacteria and fungi with an unknown
106 phyla are expected to be more prevalent. Building these global maps is possible for three main reasons;
107 firstly, the percentages of phylotypes of bacteria and fungi with an unknown phyla are highly correlated
108 with key environmental factors at the global scale (Table 1). This result suggests that environmental data
109 can be used to predict the distribution of phylotypes of fungi and bacteria unclassified at the phyla level.
110 Secondly, the database used here covers a wide gradient of environmental conditions and soil properties
111 found on Earth, being highly representative for globally distributed terrestrial ecosystems. For example,
112 mean annual precipitation and temperature in these locations ranged from 67 to 3085mm and -11.4° to
113 26.5°C, respectively. Moreover, soil pH ranged from 4.04 to 9.21; soil C from 0.15 to 34.77%; and, fine
114 texture fraction (% clay+silt) from 1.40 to 92.00%. Finally, high resolution maps for key environmental
115 factors predicting the percentage of unclassified taxa (Table 1) are available at the global scale. Therefore,
116 globally available information on environmental factors can potentially be used to predict global hotspots

117 for phylotypes of bacteria and fungi with an unknown phyla. These three important points allowed me to
118 generate global atlases for the potential distribution of percentages of phylotypes of bacteria and fungi with
119 an unknown phyla (Fig. 2). These global atlases were cross-validated as explained in Appendix 1 (Supple-
120 mentary Materials).

121 The global maps included in this study indicate the potential distribution of unclassified taxa within bacteria
122 and fungi. Interestingly, locations where bacteria with an unknown phyla are more prevalent are distinct
123 from those of fungi. This global atlas suggests that soils from Brazil, Chile, Russia, Indonesia, Iceland,
124 Northern Europe, and the coastlines of North America contain a relatively high percentage of bacteria with
125 an unknown phyla. On the other hand, deserts from Peru, China, Australia, South Africa, the Middle East,
126 the Saharan region, and the western coast of North America contain a relatively high percentage of unclas-
127 sified taxa within fungi. Soil taxonomists and microbiologists should target soils from these environments
128 and global locations to increase our chances of isolating and classifying these elusive yet significant soil
129 taxa, and thus, increase our knowledge of who they are and what they are doing in our soils.

130 **Methods**

131 **Soil sampling**

132 Soils were collected from 235 locations across eighteen countries and six continents. Soil samples (top
133 ~7.5cm depth) were collected under the most common vegetation across a wide range of ecosystem (forests,
134 grasslands, and shrublands) and climatic (arid, temperate, tropical, continental, and polar ecosystems) types.
135 The locations sampled represent wide gradients in environmental factors, which is critical for mapping
136 predictions. Detailed information about this survey can be found in Delgado-Baquerizo et al. (2018).

137 **Molecular analyses**

138 Soil DNA was extracted using the Powersoil® DNA Isolation Kit (MoBio Laboratories, Carlsbad, CA,
139 USA) according to the manufacturer's instructions. Amplicons targeting the bacterial 16S rRNA gene
140 (341F-805R; Herlemann et al. 2011) and the fungal ITS region (FITS7-ITS4R; Ihrmark et al. 2012) were
141 sequenced at Western Sydney University's NGS facility (Sydney, Australia) using the Illumina MiSeq plat-
142 form. Bioinformatic processing was performed using a combination of QIIME (Caporaso et al. 2010),
143 USEARCH (Edgar 2010) and UPARSE (Edgar 2013). Operational taxonomic units –OTUs– (phylotypes
144 hereafter), were identified at the $\geq 97\%$ identity level. Taxonomy for bacteria and fungi was assigned using
145 the Greengenes and UNITE databases, respectively. OTU abundance tables were constructed from these
146 analyses. 16s rRNA reads classified as Archaea, chloroplasts or mitochondria were removed. The percent-
147 age of phylotypes of bacteria and fungi with an unknown phyla for each sample were calculated from these
148 OTU tables. These phylotypes are classified as fungi or bacteria, but do not match data within taxonomic
149 databases at the phyla level (unclassified bacteria and fungi hereafter). Given that soil and DNA samples
150 were collected, extracted, and analysed following the same standardised protocol and within the same la-
151 boratory, any biases (e.g., sequencing error) would be consistent across analyses.

152 **Environmental factors**

153 For each location, information for twelve environmental factors was obtained: climate (maximum and min-
154 imum temperatures, precipitation seasonality; mean diurnal temperature range and Aridity Index); soil
155 properties (pH, texture and total organic carbon); dominant ecosystem type (forest and grasslands); plant
156 productivity, and UV light intensity. Information on soil pH, texture and total organic carbon (soil C) was
157 obtained using standard laboratory methods (Anderson 1993; Kettler et al. 2001) in the laboratories from
158 the Universidad Rey Juan Carlos (Spain). Climatic information (1km resolution) for all sampling locations
159 was obtained from the Worldclim database (www.worldclim.org; Hijmans et al. 2005; Zomer et al. 2018).
160 The dominant ecosystem types (forest and grasslands) were determined in the field. Plant productivity (net

161 primary productivity) data was obtained using the Normalized Difference Vegetation Index (NDVI) from
162 the Moderate Resolution Imaging Spectroradiometer (MODIS) aboard NASA's Terra satellites
163 (<http://neo.sci.gsfc.nasa.gov/>). The monthly average value for this variable was calculated between 2003-
164 2015 (~10km resolution), when all soil samplings were conducted. Information on the annual ultraviolet
165 index (UV index) was obtained from the NASA's Aura satellite (<https://neo.sci.gsfc.nasa.gov>).

166 **Mapping the global distribution of unclassified soil taxa**

167 The prediction-oriented regression model Cubist (Quinlan 1993) was used to predict the percentage of phy-
168 lotypes of bacteria and fungi with an unknown phyla across the globe. Mapping analyses were inde-
169 pendently done to find the percentage of unclassified taxa within bacteria and fungi. The Cubist algorithm
170 uses a regression tree analysis to generate a set of hierarchical rules using information on environmental
171 covariates, based on real data (235 locations), which are later used for spatial prediction (Kuhn et al. 2016).
172 Covariates in our models include the above described twelve environmental factors as well as space (lati-
173 tude and longitude). Global predictions on the distribution of the percentage of unclassified taxa within
174 bacteria and fungi were done on a 25km resolution grid, which resulted in a grid including 225530 locations.
175 Environmental information for each of these locations, including soil properties, climatic information, plant
176 production, ecosystem types and UV light, was obtained from global databases available online. Global
177 information on soil properties for this grid was obtained using the ISRIC (global gridded soil information)
178 Soil Grids (https://soilgrids.org/#!//?layer=geonode:taxnwr_b_250m). Global information on the major veg-
179 etation types in this study (grasslands and forests) was obtained using the Globcover2009 map from the
180 European Space Agency (http://due.esrin.esa.int/page_globcover.php). Global information on climate, UV
181 radiation and net primary productivity were obtained from the WorldClim database (www.worldclim.org)
182 and NASA satellites (<https://neo.sci.gsfc.nasa.gov>), as explained above. The R package Cubist was used to
183 conduct these analyses (Kuhn et al. 2016).

184 **References**

- 185
186 Anderson JM (1993) *JSI, Ingrassia, Tropical Soil Biology and Fertility: A Handbook of Methods* (CABI,
187 Wallingford, UK, ed. 2).
188 Bardgett RD, van der Putten WH (2014) Belowground biodiversity and ecosystem functioning. *Nature*
189 515, 505–511.
190 Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. (2010) QIIME
191 allows analysis of high-throughput community sequencing data. *Nat Method* 7, 335.
192 Delgado-Baquerizo M, Maestre FT, Reich PB, Jeffries TC, Gaitan JJ, Encinar D. et al. (2016) Microbial
193 diversity drives multifunctionality in terrestrial ecosystems. *Nat Commun* 28, 10541.
194 Delgado-Baquerizo M, Oliverio AM, Brewer TE, Benavent-González A, Eldridge DJ, Bardgett RD et al.
195 (2018) A global atlas of the dominant bacteria found in soil. *Science* 19, 320-325.
196 Delgado-Baquerizo, M. & Eldridge, D.J. *Ecosystems* (2019). <https://doi.org/10.1007/s10021-018-0333-2>.
197 Edgar R.C. (2010) Search and clustering orders of magnitude faster than BLAST *Bioinformatics* 26,
198 2460.
199 Edgar R.G. (2013) UPARSE: highly accurate OTU sequences from microbial amplicon reads *Nature*
200 *Methods* 10, 996-998.
201 Fierer, N (2017) Embracing the unknown: disentangling the complexities of the soil microbiome. *Nature*
202 *Reviews Microbiology* 15, 579-590.
203 Greengenes. <http://greengenes.secondgenome.com>
204 Herlemann, D.P., Labrenz M, Jürgens K, Bertilsson S, Waniek JJ, Andersson AF (2011) Transitions in
205 bacterial communities along the 2000 km salinity gradient of the Baltic Sea. *ISME Journal* 5,
206 1571–1579.

- 207 Ihrmark, K., Bödeker IT, Cruz-Martinez K, Friberg H, Kubartova A, Schenck J et al. (2012) New primers
208 to amplify the fungal ITS2 region-evaluation by 454-sequencing of artificial and natural commu-
209 nities. *FEMS Microbiol Ecol.* 82, 666–677.
- 210 J.R. Quinlan, *C4.5: Programs for Machine Learning* (Morgan Kaufmann Publishers, San Mateo, Califor-
211 nia, 1993).
- 212 Kettler TA, Doran JW, Gilbert TL (2001) Simplified method for soil particle-size determination to ac-
213 company soil-quality analyses. *Soil Sci Soc Am J* 65, 849.
- 214 M. Kuhn, S. Weston, C. Keefer, N. Coulter (2016) Cubist: Rule- And Instance-Based Regression Mode-
215 ling. R package version 0.0.19.
- 216 Marx V (2017) Microbiology: the return of culture. *Nature Methods* 14, 37–40
- 217 R.J. Hijmans, Cameron SE, Parra JL, Jones PG, Jarvis A (2005) Very high resolution interpolated climate
218 surfaces for global land areas *Int J Climatol* 25, 1965-1978.
- 219 Tedersoo L, Bahram M, Pölme S, Kõljalg U, Yorou NS, Wijesundera R et al. (2014) Fungal bio-
220 geography. Global diversity and geography of soil fungi. *Science* 346, 1256688.
- 221 UNITE. <https://unite.ut.ee>
- 222 York A. *Nature Reviews Microbiology* 16, 583 (2018).
- 223 Zomer RJ, Trabucco A, Bossio DA, Verchot LV (2008) Climate change mitigation: A spatial analysis of
224 global land suitability for clean development mechanism afforestation and reforestation. *Agric*
225 *Ecosyst Envir* 126, 67-80.

226 **Acknowledgements**

227 This project has received funding from the European Union’s Horizon 2020 research and innovation pro-
228 gramme under the Marie Skłodowska-Curie grant agreement No 702057. I would like to thank Melissa S.
229 Martín, David J. Eldridge and Fernando T. Maestre for their comments and suggestions, which have helped
230 to improve this piece. I would also like to thank Brajesh K. Singh, Noah Fierer, Richard Bardgett, Alberto
231 Benavent-González, David J. Eldridge and Fernando T. Maestre for their original contribution to the data-
232 bases included in this study.

233 **Competing financial interests.**

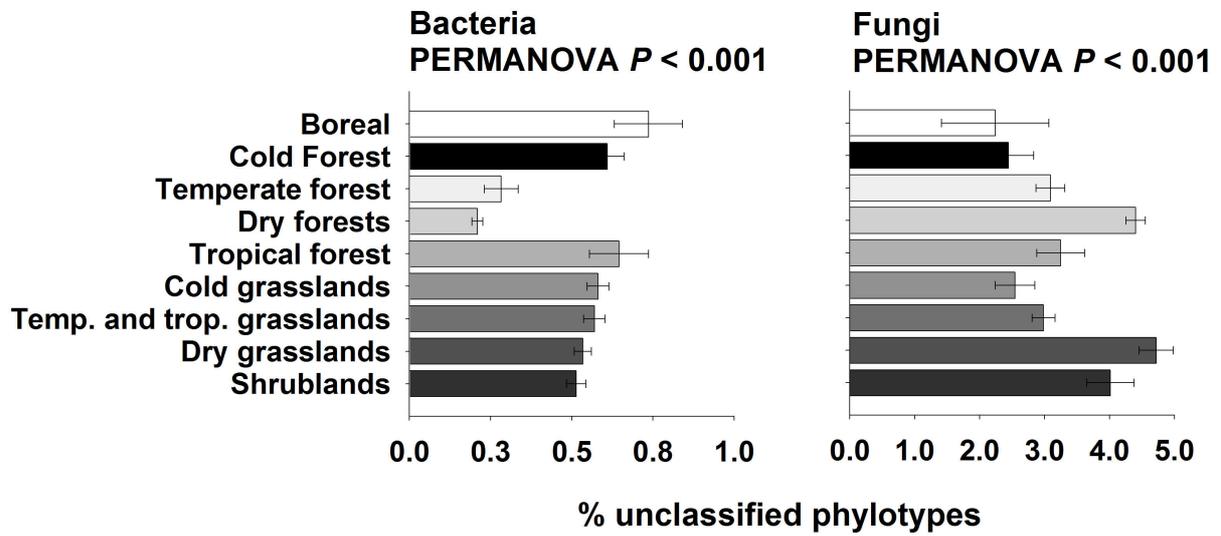
234 The authors declare no conflict of interest.

235 **Data accessibility**

236 The data used in this article will be made publicly available in a public repository (Figshare) upon publica-
237 tion.

238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255

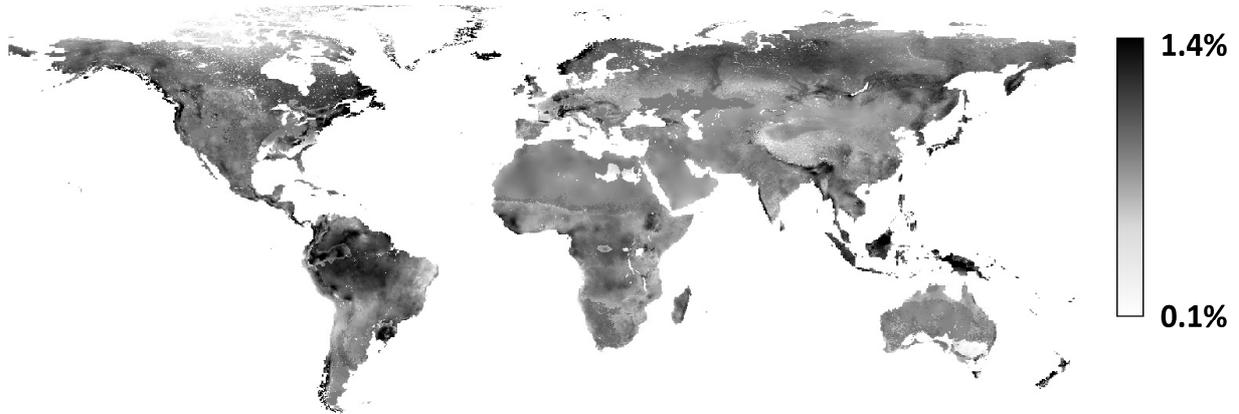
256 **Figure Captions**
257



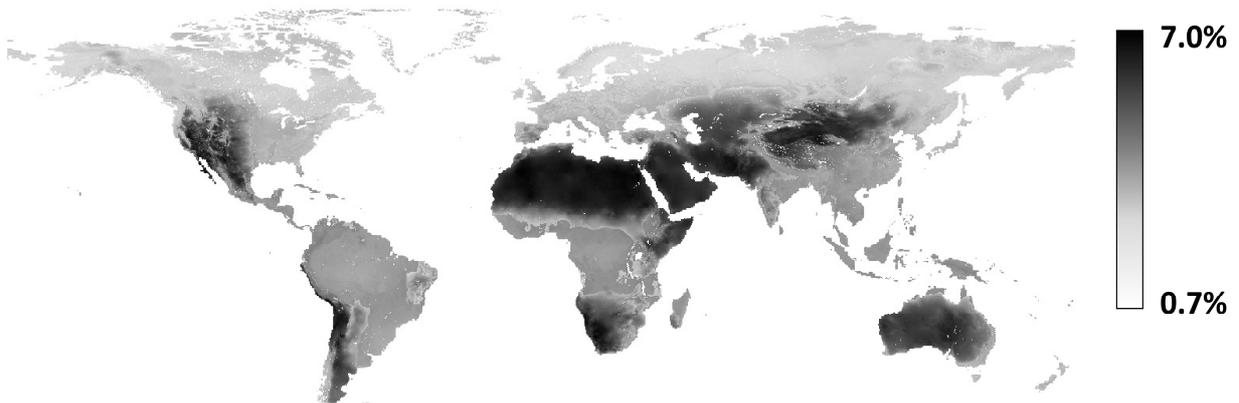
258 **Figure 1.** Mean values (\pm SE) for % phylotypes of bacteria and fungi with an unknown phyla across major
259 terrestrial biomes in 235 locations.
260

261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277

Unclassified bacteria



Unclassified fungi



278

279 **Figure 2.** Global atlas including the potential distribution of % of phylotypes of bacteria and fungi with an
280 unknown phyla (unclassified bacteria and fungi) based on their natural co-occurrence with climatic (aridity
281 index, maximum and minimum temperature, precipitation seasonality and mean diurnal temperature range),
282 primary productivity, dominant ecosystem type (forest and grasslands), soil properties (total organic car-
283 bon, pH and texture) and UV light in 235 locations. See Fig. S1 for the locations of the 235 in this study.
284 See Appendix S1 for a cross-validation of these maps.

285

286

287

288

289

290

291

292 **Table 1.** Correlation (Spearman) between the % phylotypes of bacteria and fungi with an unknown phyla
 293 (unclassified bacteria and fungi) with climate (aridity index, maximum and minimum temperature, precip-
 294 itation seasonality and mean diurnal temperature range), primary productivity, dominant ecosystem type
 295 (forest and grasslands), soil properties (total organic carbon, pH and texture) and UV light in 235 locations
 296 ($P < 0.05$). MAXT = maximum temperature. MINT = minimum temperature. Aridity Index = precipitation
 297 / potential Evapotranspiration. MDR = Mean diurnal temperature range. NPP = Net primary productivity.
 298
 299

	Longitude	Latitude	Aridity Index	MAXT	MINT	PSEA	MDR	NPP	For-ests	Grasslands	Texture	Soil C	pH	UV light
Unclassified bacteria	-0.66	0.59	0.30	-0.29	-0.33	0.45	-0.17		-0.51	0.44		-0.15		-0.25
Unclassified fungi	0.27	-0.25	-0.66	0.56	0.30		0.39	-0.51		-0.24		-0.41	0.57	0.43

300

301

302