The final typeset version of this paper is available from Nature at the following address: http://dx.doi.org/10.1038/s41586-018-0386-6

When citing this paper, use the following citation:

Bahram M, Hildebrand F, Forslund SK, Anderson JL, Soudzilovskaia NA, Bodegom PM, Bengtsson-Palme J, Anslan S, Coelho LP, Harend H, Huerta-Cepas J, Medema MH, Maltz MR, Mundra S, Olsson PA, Pent M, Põlme S, Sunagawa S, Ryberg M, Tedersoo L, Bork P: **Structure and function of the global topsoil microbiome**. *Nature*, 560, 7717, 233–237 (2018). doi: 10.1038/s41586-018-0386-6

1 Structure and function of the global topsoil microbiome

Mohammad Bahram^{1,2,3*}†, Falk Hildebrand^{4*}, Sofia K. Forslund^{4,5,6}, Jennifer L. Anderson², Nadejda A. Soudzilovskaia⁷, Peter M. Bodegom⁷, Johan Bengtsson-Palme^{8,9}, Sten Anslan^{1,10}, 2 3 Luis Pedro Coelho⁴, Helery Harend¹, Jaime Huerta-Cepas^{4,11}, Marnix H. Medema¹², Mia R. Maltz¹³, Sunil Mundra¹⁴, Pål Axel Olsson¹⁵, Mari Pent¹, Sergei Põlme¹, Shinichi Sunagawa^{4,16}, Martin Ryberg², Leho Tedersoo¹⁷[†], Peer Bork^{4, 18,19}[†] 4 5 6 7 ¹Department of Botany, Institute of Ecology and Earth Sciences, University of Tartu, 40 Lai St., 8 51005 Tartu, Estonia.²Department of Organismal Biology, Evolutionary Biology Centre, 9 Uppsala University, Norbyvägen 18D, 75236 Uppsala, Sweden. ³Department of Ecology, 10 Swedish University of Agricultural Sciences, Ulls väg 16, 756 51 Uppsala, Sweden. ⁴Structural 11 and Computational Biology Unit, European Molecular Biology Laboratory, 69117 Heidelberg, 12 Germany. ⁵Current address: Experimental and Clinical Research Center, a cooperation of 13 Charité-Universitätsmedizin Berlin and Max Delbruck Center for Molecular Medicine, 13125 14 Berlin, Germany. 6Current address: Max Delbrück Centre for Molecular Medicine, 13125 Berlin, 15 Germany. ⁷Conservation Biology Department, Institute of Environmental Sciences, CML, 16 Leiden University, Einsteinweg 2, 2333 CC Leiden, The Netherlands. ⁸Department of Infectious 17 Diseases, Institute of Biomedicine, The Sahlgrenska Academy, University of Gothenburg, 18 Guldhedsgatan 10, SE-413 46 Göteborg, Sweden. ⁹Centre for Antibiotic Resistance research 19 (CARe) at University of Gothenburg, Göteborg, Sweden.¹⁰Braunschweig University of 20 Technology, Zoological Institute, Mendelssohnstr. 4, 38106 Braunschweig, Germany. ¹¹Current 21 address: Centro de Biotecnología y Genómica de Plantas, Universidad Politécnica de Madrid 22 (UPM) - Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria (INIA), 23 Campus de Montegancedo-UPM, 28223-Pozuelo de Alarcón (Madrid) Spain. ¹²Bioinformatics 24 Group, Wageningen University, Droevendaalsesteeg 1, Wageningen, The Netherlands. ¹³Center 25 for Conservation Biology, University of California Riverside, USA Germany.¹⁴Section for 26 Genetics and Evolutionary Biology (Evogene), Department of Biosciences, University of Oslo, 27 P.O. Box 1066 Blindern, 0316 Oslo, Norway.¹⁵Biodiversity Unit, Department of Biology, 28 Ecology building, Lund University, SE-223 62 Lund, Sweden. ¹⁶Current address: Department of 29 Biology, Institute of Microbiology, ETH Zurich, 8092 Zurich, Switzerland. ¹⁷Natural History 30 Museum, University of Tartu, 14A Ravila, 50411 Tartu, Estonia.¹⁸Max Delbrück Centre for 31 Molecular Medicine, 13125 Berlin, Germany.¹⁹Department of Bioinformatics, University of 32 33 Würzburg, 97074 Würzburg, Germany. 34

- 35 *: These authors contributed equally.
- 36 *†*: Corresponding author.
- 37
- 38
- 39 Summary
- 40 Soils harbour some of Earth's most diverse microbiomes and are essential for both nutrient
- 41 cycling and carbon storage. To understand soil functioning, it is necessary to model the
- 42 global distribution patterns, biotic and environmental associations of the diversity and
- 43 structure of both bacterial and fungal communities, and their functional gene repertoires¹⁻

⁴. By leveraging metagenomics and metabarcoding of global topsoil samples (189 sites, 7560 44 subsamples), we show that bacterial, but not fungal, genetic diversity is highest in 45 temperate habitats and that microbial gene composition varies more strongly with 46 environmental variables than geographic distance. We demonstrate that fungi and bacteria 47 show global niche differentiation associated with contrasting diversity responses to 48 precipitation and soil pH. Furthermore, we provide evidence for strong bacterial-fungal 49 antagonism, inferred from antibiotics resistance genes, in topsoil and ocean habitats, 50 indicating a substantial role of biotic interactions in shaping microbial communities. Our 51 results suggest that both competition and environmental filtering affect bacterial and 52 fungal abundance, composition and their encoded gene functions, implying spatially 53 54 different relative contributions of these microbes to global nutrient cycling.

55

56 Bacteria and fungi dominate terrestrial soil habitats in terms of biodiversity, biomass, and their influence over essential soil processes⁵. Specific roles of microbial communities in 57 biogeochemical processes are reflected by their taxonomic composition, biotic interactions and 58 gene- functional potential¹⁻⁴. While microbial biogeography studies have focused largely on 59 single taxonomic groups, and on how their diversity and composition respond to local abiotic 60 soil factors (e.g. pH^{6,7}), both global patterns and the impact of biotic interactions on microbial 61 biogeography remain relatively unexplored. In addition to constraints imposed by environmental 62 factors, biotic interactions may strongly influence bacterial communities. For example, to 63 outcompete bacteria, many fungal taxa secrete substantial amounts of antimicrobial compounds⁸, 64 which select for antibiotic resistant (AR) bacteria and effectively increase relative antibiotic 65 resistance gene (ARG) abundance. Here we employed metagenomics and DNA metabarcoding 66 (16S, 18S, ITS rRNA gene markers), soil chemistry and biomass assessments (phospholipid fatty 67 acids analyses, PLFAs) to determine the relationships among genetic (functional potential), 68 phylogenetic, and taxonomic diversity and abundance in response to biotic and abiotic factors in 69 189 topsoil samples, covering all terrestrial regions and biomes of the world⁹ (Extended Data 70 Figure 1a; Supplementary Table 1). Altogether 58,000 topsoil subsamples were collected from 71 0.25-ha plots from 1450 sites (40 subsamples per site), harbouring homogeneous vegetation that 72 were minimally affected by humans. We minimized biases and shortcomings in sampling¹⁰ as 73 well as technical variation including batch effects¹¹ by using highly standardized-collection and 74

75 processing protocols. From the total collection, 189 representative sites were selected for this

analysis. We validated our main findings in external datasets, including an independent soil

dataset (145 topsoil samples; Supplementary Table 1) that followed the same sampling and

78 sequencing protocol.

79

Using metagenomics, we constructed a gene catalogue for soils, by combining our newly 80 generated data with published soil metagenomes (n=859, Supplementary Table 1) and identified 81 159,907,547 unique genes (or fragments thereof). Only 0.51% of these 160 million genes 82 overlapped with those from published genomes and large gut¹² and ocean¹³ gene catalogues that 83 are much closer to saturation (Supplementary Table 2), indicating that the functional potential of 84 soil microbiomes is enormously vast and undersampled. For functional analysis, we annotated 85 genes and functional modules via Orthologous Groups (OGs) using the eggNOG database¹⁴. For 86 each sample, we also constructed taxonomic profiles at the class and phylum levels for both 87 bacteria and fungi from relative abundance of rRNA genes in metagenomic datasets (miTags¹⁵), 88 complemented by operational taxonomic units (OTUs) based on clustering 18S rRNA and 89 internal transcribed spacer (ITS)¹⁶ genes for soil fungi and 16S rRNA genes for soil bacteria at 90 97% similarity threshold (see Methods). In total, 34,522 16S-based bacterial, 2,086 18S-based 91 and 33,476 ITS-based fungal OTUs were analysed in the context of geographic space and 16 92 93 edaphic and climatic parameters determined for each sampling site (see Methods). Archaea were poorly represented in our metabarcoding (<1% of OTUs) and metagenomics data (<1% miTags) 94 95 and hence are excluded from most analyses.

96

We examined whether the latitudinal diversity gradient (LDG), a trend of increasing diversity 97 from the poles to the tropics seen in many macroscopic organisms, especially plants¹⁷, applies to 98 microbial global distribution patterns¹⁰. We found that contrary to the typical LDG, both 99 taxonomic and gene functional diversity of bacteria peaked at mid-latitudes and declined towards 100 the poles and the equator, as is also seen in the global ocean¹³, although the pattern was relatively 101 weak for taxonomic diversity herein (Figure 1a, c; Extended Data Figure 1b,2). The deviation of 102 several bacterial phyla (5 of 20) from the general trends may be explained by responses to 103 104 edaphic and climate factors weakly related to latitude (Extended Data Figure 1b) or contrasting 105 effects at lower taxonomic levels (Supplementary discussion). In contrast, the LDG does apply to

106 overall fungal taxonomic diversity, and to 3 of 5 fungal phyla when examined separately, but not to fungal functional diversity, which was lowest in temperate biomes and exhibited an inverse 107 unimodal relationship with latitude (Figure 1b,d; Extended Data Figure 2c). The LDG was 108 negligible for oceanic fungi $(p>0.05)^{13}$, possibly due to their lower dispersal limitation and 109 paucity of plant associations. While fungal taxonomic diversity decreased poleward, the total 110 fungal biomass (inferred from PLFA markers) and the fungi-to-bacteria biomass ratio increased 111 poleward, partly due to decline of bacterial biomass decreased with latitude (Extended Data 112 Figure 3a-c). 113

114

We tested the extent to which deterministic processes (such as competition and environmental 115 filtering; i.e. the niche theory) *versus* neutral processes (dispersal and drift; the neutral theory) 116 explain distributions of fungal and bacterial taxa and functions¹⁸. In bacteria, environmental 117 variation correlated strongly with taxonomic composition (partial Mantel test accounting for 118 119 geographic distance between samples: r_{EnviGeo}=0.729, p=0.001) and moderately with gene functional composition ($r_{EnviGeo}$ =0.100, p=0.001), whereas the overall effect of geographic 120 distance among samples was negligible (p>0.05). The weak correlation between geographic and 121 taxonomic as well as functional composition suggests that environmental variables are more 122 important than dispersal capacity in determining global distributions of soil bacteria and their 123 encoded functions, as suggested by Baas Becking¹⁹ and observed for oceanic prokaryotes¹³. 124 125 For fungi, both geographic distance and environmental parameters were correlated with 126 taxonomic composition (ITS data: $r_{Geo|Env}=0.307$, p=0.001; $r_{Env|Geo}=0.208$, p=0.001; 18S data: 127 r_{Geo|Env}=0.193, p=0.001; r_{Env|Geo}=0.333, p=0.001). Environmental distance (but not geographic 128 distance) correlated with composition of fungal functional genes (r_{Env/Geo}=0.197, p=0.001), as 129 also observed for bacteria. The relatively weaker correlation of fungi with environmental 130 variation is consistent with results from local scales⁷. Thus, at both global and local scales, 131 different processes appear to underlie community assembly of fungi and bacteria. 132 133

134 To more specifically investigate the association of environmental parameters with the

135 distribution of taxa and gene functions on a global scale, we used multiple regression modelling

136 (see Methods). We found that bacterial taxonomic diversity, composition, richness and biomass

137 as well as relative abundance of major bacterial phyla can be explained by soil pH and nutrient concentration, and to a lesser extent by climatic variables (Extended Data Figures 4,5; 138 Supplementary Table 4). Bacterial community composition responded most strongly to soil pH, 139 followed by climatic variables, particularly mean annual precipitation (MAP; Extended Data 140 Figures 4,5). This predominant role of pH agrees with studies from local to continental scales⁶, 141 and may be ascribed to the direct effect of pH or confounded variables such as concentration of 142 calcium and other cations⁶. The relative abundance of genes encoding several metabolic and 143 transport pathways were strongly increased with pH (Extended Data Figure 4c), suggesting that 144 there may be greater metabolic demand for these functions for bacteria in high-nutrient and 145 alkaline conditions. 146

147

148 Compared to temperate biomes, tropical and boreal habitats contained more closely related taxa at the tip of phylogenetic trees, but from more distantly related clades (Extended Data Figure 149 2d), indicating a deeper evolutionary niche specialization in bacteria²⁰. Together with global 150 biomass patterns (Extended Data Figure 2a), these results suggest that soil bacterial communities 151 in the tropics and at high latitudes are subjected to stronger environmental filtering and include a 152 relatively greater proportion of edaphic niche specialists, possibly rendering these communities 153 more vulnerable to global change. In contrast, phylogenetic overdispersion in temperate bacterial 154 communities, may result from greater competitive pressure²⁰ or nutrient availability as predicted 155 by the niche theory 21 . 156

157

158 In contrast to the strong association between bacterial taxonomic diversity and soil pH, diversity of bacterial gene functions was more strongly correlated with MAP (Extended Data Figure 5a-h). 159 The steeper LDG in gene functions than in taxa (Figure 1a,c) may thus relate to the stronger 160 161 association of specific metabolic functions to climate than to local soil conditions. While soil and 162 climate variables exhibited comparable correlations with fungal taxa, soil carbon-to-nitrogen (C/N) ratio was the major predictor for fungal biomass and relative abundance and composition 163 of gene functions (Extended Data Figures 3g,4b,d; Supplementary Table 4). We hypothesize that 164 compared to bacteria, global distribution of fungi is more limited by resource availability due to 165 specialization for the use of specific compounds as substrates and greater energy demand. 166 167

168 We interpret opposing biogeographic trends for bacteria and fungi as niche segregation, driven by differential responses of bacteria and fungi to environmental factors⁷ and their direct 169 competition. Gene functional diversity of both bacteria and fungi responded to MAP and soil pH, 170 albeit in opposite directions (Extended Data Figure 5c.d.g.h; Supplementary Table 3). This may 171 partly explain the observed inverse pattern of gene functional diversity across the latitudinal 172 gradient, i.e. niche differentiation, between bacteria and fungi (Figure 1; Extended Data Figure 173 2). While increasing precipitation seems to favour higher fungal diversity, it is associated with 174 higher B/F biomass and abundance ratios (Extended Data Figure 3d,g; Extended Data Figure 175 5f,h). The increasing proportion of fungi towards higher latitudes may be explained by 176 competitive advantages perhaps due to a greater tolerance to nutrient and water limitation 177 associated with potential long-distance transport by hyphae. 178

179

A role of inter-kingdom biotic interactions in determining the distributions of functional diversity 180 and biomass in fungi and bacteria has been suggested previously²². As competition for resources 181 affect the biomass of fungi and bacteria^{22,23}, we hypothesized that B/F biomass ratio is related to 182 the prevalence of fungi and bacterial AR capacity because of broader activities of fungi than 183 bacteria in utilizing complex carbon substrates²⁴ as well as increased antibiotic production of 184 fungi in high C/N environments²⁵. Consistent with this hypothesis, we found that both fungal 185 biomass and the B/F biomass ratio correlated with ARG relative abundance (Extended Data 186 Figure 6) and that most fungal OG subcategories, particularly those involved in biosynthesis of 187 antibiotic and reactive oxygen species, increased with soil C/N ratio (Supplementary Table 4; 188 Supplementary results). We also found that ARG relative abundance in topsoil is more strongly 189 related to fungal relative abundance (r=0.435, p $<10^{-9}$) and B/F abundance ratio (r=-0.445, p $<10^{-1}$ 190 ¹²; Figure 2b) than to bacterial relative abundance (r=0.232, p=0.002, based on miTags), which is 191 supported by our external validation dataset (fungal relative abundance r=0.637, p< 10^{-15} ; B/F 192 abundance ratio r=-0.621, p< 10^{-15} ; bacterial relative abundance r=0.174, p=0.036). Also, topsoil 193 ARG relative abundance was significantly negatively correlated with bacterial phylogenetic 194 diversity and OTU richness based on 16S rRNA gene (Extended Data Figures 7a,c,8a), further 195 supporting a role for biotic interactions in shaping microbial communities. 196 197

198 We also tested possible direct and indirect relationships between ARGs and 16 environmental

199 predictors using structural equation modelling (SEM; Supplementary Table 5). The optimized model suggests that soil C/N ratio and moisture, rather than pH – the predominant driver of 200 bacterial diversity (Extended Data Figure 3g, Supplementary results) – affect B/F abundance 201 ratio that in turn affects ARG relative abundance at the global scale (Figure 2c). In line increased 202 antibiotics production in high competition environments, soil C/N ratio was the best predictor for 203 richness of fungal functional genes ($r^2=0.331$, $p<10^{-15}$; Supplementary Table 3) and bacterial 204 CAZyme genes involved in degrading fungal carbohydrates (r=0.501, p<10⁻¹²). ARG relative 205 abundance was also strongly correlated with C/N ratio in the external validation dataset (r=0.505, 206 p<10⁻¹⁰). 207

208

While the concomitant increase in AR potential and relative abundance of bacteria (as potential 209 210 ARG carriers) was expected, the strong correlation of fungal relative abundance with ARG 211 relative abundance and in turn bacterial phylogenetic diversity may be explained by selection 212 against bacteria that lack ARGs, such that bacteria surviving fungal antagonism are enriched for ARGs. Among all studied phyla, the relative abundance of Chloroflexi, Nitrospirae, and 213 Gemmatimonadetes bacteria (based on miTags), taxa with relatively low genomic ARG content 214 (Supplementary Table 6) were most strongly negatively correlated with ARG relative abundance 215 (Figure 3a). In contrast, ARGs were strongly positively correlated with the relative abundance of 216 Proteobacteria, which have the greatest average number of ARGs per genome²⁶ among bacteria 217 (Supplementary Table 6), and the fungal phyla Ascomycota and Zygomycota s.lat. (including 218 219 Zoopagomycota and Mucoromycota) in both the global soil and the external validation sets (Figure 3a,b; Extended Data Figure 9a,c; Supplementary Table 7). More specifically, ITS 220 metabarcoding revealed increasing relative abundances of ARGs with numerous fungal OTUs 221 222 (Supplementary Table 8), particularly those belonging to *Oidiodendron* (Myxotrichaceae, Ascomycota) and *Penicillium* (Aspergillaceae, Ascomycota), which are known antibiotic 223 producers^{27,28} (Supplementary Results). Among bacterial ARGs, the relative abundance of efflux 224 pumps and beta-lactamases, which act specifically on fungal-derived antibiotics, were 225 significantly correlated to the relative abundance of Ascomycota (Extended Data Figure 10a; 226 Supplementary Table 7). Actinobacteria, encompassing antibiotics-producing Streptomyces, also 227 significantly correlated to ARG diversity in topsoil (Supplementary Table 6). Together these 228 229 results suggest that relationships between organismal and ARG abundances are likely the result

230 of selective and/or suppressive actions of antibiotics on bacteria.

231

Consistent with our observations in topsoil, we found evidence for antagonism between fungi 232 and bacteria in oceans by reanalysing ARG distribution in 139 water samples from the global 233 Tara Oceans project¹³ (see Methods; Supplementary Table 1; Extended Data Figure 8a): the 234 fungi-like stramenopile class Oomycetes (water moulds) and the fungal phylum Chytridiomycota 235 constituted the groups most strongly associated with bacterial ARG relative abundance (Figure 236 3a,c, Extended Data Figures 9b,d,10b,d). Although there is little direct evidence that oomycetes 237 produce antibiotics, their high antagonistic activity can trigger bacteria²⁹ and other organisms 238 including fungi³⁰ to produce antibiotics (Supplementary Discussion). As in topsoil, bacterial 239 phylogenetic diversity was significantly negatively correlated with ARG relative abundance in 240 241 ocean samples (Extended Data Figure 7b,c). In addition, the ARG relative abundance declined 242 with increasing distance from the nearest coast in ocean samples (Extended Data Figure 8b), 243 which may reflect the effect of a decreasing nutrient gradient along distance from the coast on the pattern of bacteria and fungi abundance and in turn ARG abundance. The agreement of 244 results from these disparate habitats suggests that competition for resources related to nutrient 245 availability and climate factors drive a eukaryotic-bacterial antagonism in both terrestrial and 246 oceanic ecosystems. 247

248

Our results indicate that both environmental filtering and niche differentiation determine global 249 soil microbial composition, with a minor role of dispersal limitation at this scale (for limitations, 250 see Methods). In particular, global distribution of soil bacteria and fungi was most strongly 251 associated with soil pH and precipitation, respectively. Our data further indicate that inter-252 kingdom antagonism, as reflected in the association of bacterial ARGs with fungal relative 253 254 abundance, is also important in structuring microbial communities. Although further studies are 255 needed to explicitly address the interplay of B/F abundance ratio and ARG abundance, our data suggest that environmental variables that impact B/F abundance ratio may have consequences for 256 257 microbial interactions and favouring fungi- or bacteria-driven soil nutrient cycling. This unprecedented view of global patterns of microbial distributions implies that global climate 258 259 change may differentially affect bacterial and fungal composition and their functional potential, 260 because acidification, nitrogen pollution and shifts in precipitation all have contrasting effects on

261 topsoil bacterial and fungal abundance, diversity and functioning.

262

263 **Acknowledgments** The authors thank Ingrid Liiv for technical and laboratory assistance; Sebastian Waszak for constructive comments on the manuscript; Yan P. Yuan and Anna Glazek 264 for bioinformatics support; and Alexander Holm Viborg for help in retrieving the CAZY 265 database. We also thank Vladimir Benes, Raina Hercog and other members of the EMBL 266 GeneCore (Heidelberg), who provided assistance and facilities for sequencing. This study was 267 funded by the Estonian Research Council (grants PUT171, PUT1317, PUT1399, IUT20-30, 268 MOBERC, KIK, RMK, ECOLCHANGE), the Swedish Research Council (VR grant 2017-269 05019), Royal Swedish Academy of Sciences, Helge Axson Johnsons Stiftelse, EU COST 270 Action FP1305 Biolink (STSM grant), Netherlands Organization for Scientific research (vidi 271 grant 016.161.318), EMBL European Union's Horizon 2020 Research and Innovation 272 Programme (#686070; DD-DeDaF) and Marie Skłodowska-Curie (600375). 273 274 Author Contributions Author Contributions M.B., L.T. and P.B. conceived the project. L.T. 275 supervised DNA extraction and sequencing. M.B., F.H., S.F., J.L.A., M.R. and P.B. designed 276 and supervised the data analyses. F.H. designed and performed bioinformatics analysis. N.A.S 277 and P.A.O. performed biomass analysis. S.F., S.M., M.P., S.A., H.H., S.P., M.R.M., S.S., and 278 L.T. contributed data. M.B., F.H., S.F., J.L.A., P.M.B., S.A., J.B.P., M.H.M., L.P.C. and J.H.C. 279 performed the data analyses. M.B. wrote the first draft of the manuscript with significant input 280 from F.H., S.F., J.L.A., L.T. and P.B. All authors contributed to data interpretation and editing of 281 282 the paper.

Author Information The authors declare no competing financial interests. Correspondence
should be addressed to P.B. (bork@embl.de), L.T. (leho.tedersoo@ut.ee) or M.B.
(bahram@ut.ee). Requests for materials should be addressed to M.B. (bahram@ut.ee) and F.H.
(falk.hildebrand@embl.de).

288 289

291

283

290 **References**

- Green, J. L., Bohannan, B. J. & Whitaker, R. J. Microbial biogeography: from taxonomy to traits. *Science* 320, 1039-1043 (2008).
- 294 2 Reed, D. C., Algar, C. K., Huber, J. A. & Dick, G. J. Gene-centric approach to integrating
 295 environmental genomics and biogeochemical models. *Proc. Natl. Acad. Sci. USA* 111,
 296 1879-1884 (2014).
- Maynard, D. S., Crowther, T. W. & Bradford, M. A. Fungal interactions reduce carbon
 use efficiency. *Ecol. Lett.* 20, 1034-1042 (2017).
- 299 4 de Menezes, A. B., Richardson, A. E. & Thrall, P. H. Linking fungal-bacterial co-
- 300 occurrences to soil ecosystem function. *Curr. Opin. Microbiol.* **37**, 135-141 (2017).
- Bardgett, R. D. & van der Putten, W. H. Belowground biodiversity and ecosystem
 functioning. *Nature* 515, 505-511 (2014).
- Lauber, C. L., Hamady, M., Knight, R. & Fierer, N. Pyrosequencing-based assessment of
 soil pH as a predictor of soil bacterial community structure at the continental scale. *Appl. Environ. Microbiol.* 75, 5111-5120 (2009).

306 307	7	Rousk, J. <i>et al</i> . Soil bacterial and fungal communities across a pH gradient in an arable soil. <i>ISME J</i> . 4 , 1340 (2010).
308	8	de Boer, W., Folman, L. B., Summerbell, R. C. & Boddy, L. Living in a fungal world:
309		impact of fungi on soil bacterial niche development. FEMS Microbiol. Rev. 29, 795-811
310	0	(2005).
311	9	Olson, D. M. <i>et al.</i> Terrestrial Ecoregions of the World: A New Map of Life on Earth: A
312 313		new global map of terrestrial ecoregions provides an innovative tool for conserving biodiversity. <i>Bioscience</i> 51 , 933-938 (2001).
314	10	Green, J. & Bohannan, B. J. Spatial scaling of microbial biodiversity. <i>Trends Ecol. Evol.</i>
315	-	21 , 501-507 (2006).
316	11	Yuan S Cohen D B Ravel J Abdo Z & Forney L J Evaluation of methods for the
317		extraction and purification of DNA from the human microbiome <i>PLoS ONE</i> 7 e33865
318		(2012)
319	12	Li L et al. An integrated catalog of reference genes in the human gut microbiome Nat
320	12	Riotechnol 32 , 834-841 (2014)
321	13	Sunagawa S et al. Structure and function of the global ocean microbiome. Science 348
321	15	1261359 (2015)
322	14	Huerta-Cepas I et al. $eggNOG 4.5$: a hierarchical orthology framework with improved
323	17	functional annotations for eukarvotic prokarvotic and viral sequences Nucleic Acids Res
324		44 D286-D293 (2015)
325	15	Logares R <i>et al.</i> Metagenomic 16S rDNA Illumina tags are a nowerful alternative to
320	15	amplicon sequencing to explore diversity and structure of microbial communities
220		Environ Microbiol 16 2650 2671 (2014)
320	16	Tederson I <i>et al</i> Global diversity and geography of soil fungi <i>Science</i> 346 1078_+
330	10	(2014)
331	17	Willig, M. R., Kaufman, D. & Stevens, R. Latitudinal gradients of biodiversity: pattern.
332		process, scale, and synthesis, Annu, Rev. Ecol. Evol. Syst. 34 , 273-309 (2003).
333	18	Martiny, J. B. H. <i>et al.</i> Microbial biogeography: putting microorganisms on the map. <i>Nat.</i>
334		<i>Rev. Microbiol.</i> 4 . 102-112 (2006).
335	19	Baas-Becking, L. G. M. <i>Geobiologie: of inleiding tot de milieukunde</i> . (WP Van Stockum
336		& Zoon NV. 1934).
337	20	Webb, C. O., Ackerly, D. D., McPeek, M. A. & Donoghue, M. J. Phylogenies and
338		community ecology. Annu. Rev. Ecol. Syst. 33, 475-505 (2002).
339	21	Bryant, J. A., Stewart, F. J., Eppley, J. M. & DeLong, E. F. Microbial community
340		phylogenetic and trait diversity declines with depth in a marine oxygen minimum zone.
341		<i>Ecology</i> 93 . 1659-1673 (2012).
342	22	Frev-Klett, P. <i>et al.</i> Bacterial-fungal interactions: hyphens between agricultural, clinical.
343		environmental, and food microbiologists. <i>Microbiol. Mol. Biol. Rev.</i> 75 , 583-609 (2011).
344	23	Mille Lindblom, C., Fischer, H. & J Tranvik, L. Antagonism between bacteria and
345	-	fungi: substrate competition and a possible tradeoff between fungal growth and tolerance
346		towards bacteria. <i>Oikos</i> 113 . 233-242 (2006).
347	24	Koranda M <i>et al.</i> Fungal and bacterial utilization of organic substrates depends on
348		substrate complexity and N availability. FEMS Microbiol. Ecol. 87, 142-152 (2014)
349	25	Platas G Pelaez F Collado J Villuendas G & Diez M Screening of antimicrobial
350		activities by aquatic hyphomycetes cultivated on various nutrient sources. Cryptogamie:
351		<i>Mycol.</i> 19 , 33-43 (1998).

35226Mende, D. R. *et al.* proGenomes: a resource for consistent functional and taxonomic353annotations of prokaryotic genomes. *Nucleic Acids Res.*, D529-D534 (2017).

- Bérdy, J. Thoughts and facts about antibiotics: where we are now and where we are
 heading. J. Antibiot. 65, 385-395 (2012).
- Andersen, N. R. & Rasmussen, P. The constitution of clerocidin a new antibiotic isolated
 from Oidiodendron truncatum. *Tetrahedron Lett.* 25, 465-468 (1984).
- Zhao, Y., Qian, G., Chen, Y., Du, L. & Liu, F. Transcriptional and antagonistic responses
 of biocontrol strain Lysobacter enzymogenes OH11 to the plant pathogenic oomycete
 Pythium aphanidermatum. *Front. Microbiol.* 8, 1025 (2017).
- 36130Takahashi, K. *et al.* Cladomarine, a new anti-saprolegniasis compound isolated from the362deep-sea fungus, Penicillium coralligerum YK-247. J. Antibiot. 70, 911 (2017).

365 Figure Legends

363 364

366

367 Figure 1 | Fungal and bacterial diversity exhibit contrasting patterns across the latitudinal gradient. Latitudinal distributions of bacterial (left columns) and fungal (right columns) 368 taxonomic (a and b; n=188 biologically independent samples) and gene functional (c and d; 369 n=189 biologically independent samples) diversity in the global soil samples. The order of 370 polynomial fit was chosen based on the corrected Akaike Information Criterion (AICc; see 371 Methods) of first and second order polynomial models (ANOVA: **a**: F=34.28; $p<10^{-7}$; **b**: 372 F=3.84, p=0.052; c: F= 50.48, p<10⁻¹⁰; d: F= 18.55, p= p<10⁻⁴). Grey dashed and black solid 373 lines are the first and second order polynomial regression lines, respectively. Diversity was 374 measured using Inverse Simpson Index (these trends were robust to choice of index, see 375 Extended Data Figure 2b, c). The latitudinal distribution of the high-level biome (tropical, 376 temperate and boreal-arctic) is given at the top of **a**) and **b**). 377 378 Figure 2 | Global relative abundance of antibiotic resistance genes (ARGs) can be explained 379 by a combination of biotic and abiotic factors. a, Pairwise Spearman correlation matrix of 380 main biotic and abiotic determinants of ARG relative abundance. b, B/F abundance ratio 381 significantly correlated with ARG relative abundance on a global scale. **c**, Structural equation 382 modelling (SEM) of ARG relative abundance of soil (green) and ocean (blue) datasets 383 (explaining 44% and 51% of variation, respectively; Supplementary Table 5). The goodness of 384

- fit was acceptable (Soil: RMSEA=0.00, PCLOSE=0.989, n=189 biologically independent
- samples; Ocean: RMSEA=0.059, PCLOSE=0.302, n=139 biologically independent samples).
- Abbreviations: C/N, carbon to nitrogen ratio; N, nitrates; Bacteria/Fungi (B/F), the ratio of bacterial to fungal abundance/biomass; Bacterial richness, bacterial OTU (>97% similarity)
- richness based on metabarcoding dataset; Abundance, relative abundance of miTags determined as fungi or bacteria; Biomass (nmol/g), absolute biomass based on PLFA analysis; MAP: mean
- annual precipitation; MAT: mean annual temperature; n.a.: not applicable; n.s.: not significant (p>0.05, q>0.1).
- 393
- 394Figure 3 | Fungi are the main determinants of antibiotic resistance gene (ARG) relative
- **abundance in soils and oceans. a**, The association between ARG relative abundance and major

bacterial and fungal (incl. fungal-like protist) phyla in metagenomic samples from soil and
 ocean. Outer circle colour corresponds to the Pearson correlation coefficient. Circle fill colour

corresponds to significance after adjustment for multiple testing (q-value), as indicated in the 398 legend. b-c, Relationships (non-parametric correlations) between the relative abundances of most 399 correlated fungal groups with ARGs in soil metagenomes from this study (b) and ocean 400 metagenomes (c). For statistical details and significance, see Supplementary Table 8. Asterisks 401 denote significance after Benjamini-Hochberg correction for multiple testing (*, q < 0.1). See also 402 supplementary analysis and Supplementary Table 8 for analogous results as in (a) but at the class 403 level and in other habitats besides soil and ocean including published non-forest and agricultural 404 soil as well as human skin and gut samples. 405 406

407

408 METHODS

409 Soil sample preparation

Composite soil samples from 1450 sites worldwide were collected using highly standardized 410 protocols¹⁶. The sampling was conducted broadly across the most influential known 411 environmental gradient – the latitude - taking advantage of a global "natural laboratory" to study 412 the impact of climate on diversity across vegetation, biome and soil types and to enable testing 413 414 the effects of environmental parameters, spatial distance, and biotic interactions in structuring microbial communities. We carefully selected representative sites for different vegetation types 415 416 separated by spatial distances sufficient to minimize spatial autocorrelation and to cover most 417 areas of the globe. Total DNA was extracted from 2.0 g of soil from each sample using the PowerMax Soil DNA Isolation kit (MoBio, Carlsbad, CA, USA). A subset of 189 high-quality 418 DNA samples representing different ecoregions spanning multiple forest, grassland and tundra 419 biomes (Supplementary Table 1) were chosen for prokaryote and eukaryote metabarcoding 420 (ribosomal rRNA genes) and whole metagenome analysis. Samples from desert (n=8; G4010, 421 G4034, S357, S359, S411, S414, S418 and S421) and mangrove (n=1: G4023) biomes yielded 422 sufficient DNA for metabarcoding, but not for metagenomics sequencing, thus these samples 423 were used for global mapping of taxonomic diversity but excluded from all comparisons between 424 425 functional and taxonomic diversity. One sample (S017) contained no 16S sequences; thus, altogether 189 and 197 samples were used for metagenomics and metabarcoding analyses, 426 427 respectively.

428

To determine the functional gene composition of each sample, 5 µg total soil DNA (300-400 bp 429 fragments) was ligated to Illumina adaptors using the TruSeq Nano DNA HT Library Prep Kit 430 (Illumina Inc., San Diego, CA, USA) and shotgun-sequenced in three runs of the Illumina HiSeq 431 2500 platform $(2 \times 250 \text{ bp paired-end chemistry, rapid run mode})^{31}$ in the Estonian Genomics 432 Center (Tartu, Estonia). Taxonomic composition was estimated from the same DNA samples 433 434 using ribosomal DNA metabarcoding for bacteria (16S V4 subregion) and eukaryotes (18S V9 subregion). For amplification of prokaryotes and eukaryotes, universal prokaryote primers 515F 435 and 806RB³² (although this pair may discriminate against certain groups of Archaea and Bacteria 436 such as Crenarchaeota/Thaumarchaeota (and SAR11, see ref. ³³) and eukaryote primers 1389f 437

and $1510r^{34}$ were used. While the resolution of 16s rRNA sequencing is limited to genus (and

439 higher) -level assignments, it is currently a standard approach in profiling bacterial communities

and thus enabled us at least to explore patterns at coarse phylogenetic resolution.

441

Each primer was tagged with a 10-12-base identifier barcode¹⁶. DNA samples were amplified using the following PCR conditions: 95 °C for 15 min, followed by 30 cycles of 95 °C for 30 s,

50 °C 45 s and 72 °C for 1 min with a final extension step at 72 °C for 10 min. The 25 μ I PCR 444 mix consisted of 16 μ l sterilized H₂O, 5 μ l 5× HOT FIREPol Blend MasterMix (Solis Biodyne, 445 Tartu, Estonia), 0.5 µl each primer (200nM) and 3 µl template DNA. PCR products from three 446 technical replicates were pooled and their relative quantity was evaluated after electrophoresis on 447 an agarose gel. DNA samples producing no visible band or an overly strong band were amplified 448 using 35 and 25 cycles, respectively. The amplicons were purified (FavorPrep[™] Gel/PCR 449 Purification Kit; Favorgen), checked for quality (ND 1000 spectrophotometer; NanoDrop 450 Technologies), and quantified (Oubit dsDNA HS Assay Kit; Life Technologies). Quality and 451 concentration of 16S amplicon pools were verified using Bioanalyzer HS DNA Analysis Kit 452 (Agilent) and Oubit 2.0 Fluorometer with dsDNA HS Assay Kit (Thermo Fisher Scientific), 453 respectively. Sequencing was performed on an Illumina MiSeq at the EMBL GeneCore facility 454 455 (Heidelberg, Germany) using a v2 500 cycle kit, adjusting the read length to 300 and 200 bp for read1 and read2, respectively. 18S amplicon pools were quality checked using Bioanalyzer HS 456 457 DNA Analysis Kit (Agilent), quantified using Qubit 2.0 Fluorometer with dsDNA HS Assay Kit (Thermo Fisher Scientific) and sequenced on an Illumina HiSeq at Estonian Genomics Center 458 459 (Tartu, Estonia). Sequences resulting from potential contamination and tag-switching were identified and discarded based on two negative and positive control samples per sequencing run. 460

461

462 Soil chemical analysis and biomass analysis

463 All topsoil samples were subjected to chemical analysis of pH_{KCl} , P_{total} , K, Ca and Mg; the 464 content of ${}^{12}C$, ${}^{13}C$, ${}^{14}N$ and ${}^{15}N$ were determined using an elemental analyzer (Eurovector, 465 Milan, Italy) coupled with an isotope ratio mass spectrometer⁵⁵.

466

To calculate the absolute abundance of bacteria and fungi using an independent approach, 467 bacterial and fungal biomass were estimated from Phospholipid Fatty Acids (PLFAs)³⁵ in nmol/g 468 as follows. Lipids were extracted from 2 g freeze dried soil in a one-phase solution of 469 chloroform, methanol and citrate buffer³⁶. Chloroform and citrate buffer was added to split the 470 collected extract into one lipophilic phase, and one hydrophilic phase. The lipid phase was 471 collected and applied on a pre-packed silica column³⁶. The lipids were separated into neutral 472 lipids, intermediate lipids and polar lipids (containing the phospholipids) by subsequent elution 473 with chloroform, acetone and methanol. The neutral and phospholipids were dried using a speed 474 vac. Methyl nonadecanoic acid (Me19:0) was added as an internal standard. The lipids were 475 subjected to a mild alkaline methanolysis, in which fatty acids were derivatised to fatty acid 476 methyl esters (FAMEs). The FAMEs from neutral (NLFAs) and phospholipids (PLFAs) were 477 478 dried, using speed vac, and then dissolved in hexane before analysis on a gas-chromatograph as described by ref.³⁷. Fungal biomass was estimated as the concentration of PLFA 18:206.9 and 479 bacterial biomass from the sum of nine PLFAs (i15:0, i16:0, i17:0, a15:0, a17:0, cy17:0, cy19:0, 480 10Me17:0 and 10Me18:0)³⁶. The nomenclature of fatty acids follows Frostegård et al.³⁷. 481

482

483 Acquisition of metadata from public databases

Climate data including monthly temperature and precipitation were obtained from the WorldClim database (www.worldclim.org). In addition, estimates of soil carbon, moisture, pH, potential evapotranspiration (PET) and net primary productivity (NPP) at 30 arc minute resolution were obtained from the Atlas of the Biosphere (www.sage.wisc.edu/atlas/maps.php). Samples were categorized into 11 biomes⁹, with all grassland biomes being categorized as "grasslands". Thus, the following biomes were considered and summarized to three global levels: moist tropical forests, tropical montane forests and dry tropical forests, savannas as
tropical; Mediterranean, grasslands and shrublands, southern temperate forests, coniferous
temperate forests and deciduous temperate forests as temperate; and boreal forests and arctic
tundra as boreal-arctic. The time from the last fire disturbance was estimated based on inquiry
from local authorities or collaborators and evidence from the field.

495

496 Metagenome analysis

Most soil microbes are uncultured, making their identification difficult. Metagenomics analysis
has emerged as a way around this to capture both genetic and phylogenetic diversity. As such it
can only directly reveal the potential for functions through determining and tracing gene family
abundances (as opposed to realized protein activity), which may be involved in various
functional pathways³⁸, but we can safely assume a strong correspondence between gene

functional potential and the resulting ecosystem functioning³⁹ or enzyme activities⁴⁰.

503

Reads obtained from the shotgun metagenome sequencing of topsoil samples were qualityfiltered, if the estimated accumulated error exceeded 2.5 with a probability of $\ge 0.01^{41}$, or >1ambiguous position. Reads were trimmed if base quality dropped below 20 in a window of 15 bases at the 3' end, or if the accumulated error exceeded 2 using the sdm read filtering software⁴². After this, all reads shorter than 70% of the maximum expected read length (250 bp

software⁴². After this, all reads shorter than 70% of the maximum expected read length (250 bp unless noted otherwise for external datasets) were removed. This resulted in retention of

510 894,017,558 out of 1,307,037,136 reads in total (Supplementary Table 1). We implemented a

511 direct mapping approach to estimate the functional gene composition of each sample. First, the

512 quality-filtered read pairs were merged using FLASH⁴³. The merged and unmerged reads were

513 mapped against functional reference sequence databases (see below) using DIAMOND 0.8.10 in

blastx mode44 using "-k 5 -e 1e-4 --sensitive" options. The mapping scores of two unmerged

515 query reads that mapped to the same target were combined to avoid double counting. In this case, 516 the hit scores were combined by selecting the lower of the two e-values and the sum of the bit

517 scores from the two hits. The best hit for a given query was based on the highest bit score,

518 longest alignment length and highest percent identity to the subject sequence. Finally, aligned

reads were filtered to those, having an alignment %identity >50% and matching with an e-value

- 520 <1e-9 (see below for parameter choice).
- 521

The functional databases to which metagenomic reads were mapped included gene categories related to ROS sources (peroxidases genes databases^{45,46}, KEGG⁴⁷ (Kyoto Encyclopedia of Genes and Genomes) and CAZyme genes (<u>www.cazy.org</u>, accessed 22.11.2015)⁴⁸. To facilitate 522 523 524 interpretation of the results, the relative abundance of CAZyme genes were summed based on the 525 substrates for each gene family. Substrate utilization information for CAZyme families was 526 obtained from ref. well CAZvpedia 527 as as (http://www.cazypedia.org/index.php?title=Carbohydrate-binding modules&oldid=9411). Based 528 on the KEGG Ortholog (KO) abundance matrices we calculated SEED functional module 529 abundances. For functional annotations of metagenomic reads, we used in silico annotation based 530 on a curated database of the orthologous gene family resource $eggNOG 4.5^{14}$. 531

532

For all databases that included taxonomic information (eggNOG, KEGG, CAZy), reads were mapped competitively against all kingdoms and assigned into prokaryotic and eukaryotic groups, based on the best bit score in the alignment and the taxonomic annotation provided with the

database at kingdom level. All functional abundance matrices were normalized by the total 536 number of reads used for mapping in the statistical analysis, unless mentioned otherwise (e.g. 537 rarefied in the case of diversity analysis, see below). This normalization better takes into account 538 differences in library size as it has the advantage of including the fraction of unmapped (that is 539 functionally unclassified) reads. Although there are limitations in using relative abundance of 540 genes, our analysis shows, which potential functions are relatively more important. Without any 541 normalisation, such analyses cannot be performed. It is currently difficult to test the absolute 542 numbers, due to limitations to reliably quantify soil DNA resulting from differences in extraction 543 efficiency and level of degradation. 544

545

To identify ARGs in our metagenome samples, the merged and unmerged reads were mapped to 546 a homology expansion (see ref.⁵¹) of the Antibiotic Resistance gene Data Base (ARDB). Only 547 hits surpassing the minimum sequence identity values as listed in the ARDB for each family 548 549 were taken further into account. While there exist newer ARG databases, only the ARDB presently have curated family inclusion thresholds that directly allow application to our topsoil 550 551 dataset: as soil microbial diversity is so large, unlike for gut datasets, high-fidelity gene catalogue construction will not be possible until many more samples are available. Therefore, 552 direct mapping of reads to the gene family databases becomes necessary for our analysis, in turn 553 necessitating ARG inclusion thresholds that are well-defined also for single reads, not merely for 554 full-length genes. Thus, the cut-offs curated for e.g. ResFams⁵² or CARD⁵³ are inappropriate, 555 since they are defined in the length-dependent bit score space. The ARDB cut-offs, however, are 556 defined as sequence identities, thus in principle applicable also to shorter than full-length 557 558 sequences. Because of these technical limitations, we used a soil gene catalogue to determine CARD based ARG abundance matrices (see further on). 559 560

- 561 It is important to note that functional gene including ARG measurements represent relative 562 proportions of different gene families, because the absolute amount of DNA differs among 563 samples. This necessitates, as we have done, to choose statistical tests that do not assume 564 absolute measurements, and centres analysis of this type on comparisons across the set of 565 samples.
- 566

567 miTag taxa abundance estimation

We used a miTag approach¹⁵ to determine bacterial and fungal community composition from 568 metagenome sequence data. First, SortMeRNA⁵⁴ was used to extract and blast search rRNA 569 genes against the SILVA LSU/SSU database. Reads approximately matching these databases 570 with e-values $<10^{-1}$ were further filtered with custom Perl and C++ scripts, using FLASH to 571 572 attempt merging all matched read pairs. In case read pairs could not be merged, as happens if the overlap between them is too small, the reads were interleaved such that the second read pair was 573 reverse complemented and then sequentially added to the first read. To fine-match candidate 574 interleaved or merged reads to Silva LSU/SSU databases, lambda⁵⁵ was used. Using the lowest common ancestor (LCA) algorithm adapted from LotuS (version 1.462)⁴², we determined the 575 576 identity of filtered reads based on lambda hits. This included a filtering step, where queries were 577 only assigned to phyla and classes if they had at least 88% and 91% similarity to the best 578 database hit, respectively. The taxon by sample matrices were normalized by the total number of 579 reads per sample to minimize the effects of uneven sequencing depth. The average of SSU and 580 581 LSU matrices was used for calculating the relative abundance of phyla/classes. The abundance of

miTag sequences matching bacteria and fungi was used to determine B/F abundance ratio. While 582 LSU/SSU assessments refer to number of fungal cells rather than number of discrete 583 multicellular fungi, since this can apply to all samples equally, it is not systematically biased for 584 comparing the trends of bacterial to fungal abundance across samples. 585

586

587 **External metagenomic datasets**

We validate and compare the global trends with those on a smaller scale, we used a regional 588 scale dataset of 145 topsoil generated and processed using the same protocol as our global 589 dataset (Supplementary Table 1). 590

591

In addition, to compare patterns of ARG diversity in soils and oceans on a global scale, we re-592

analysed the metagenomics datasets of the Tara Oceans¹³, including all size fractions 593

(Supplementary Table 1). After quality filtering, 41,790,928,650 out of 43,076,016,494 reads 594 595 were retained from the Tara Oceans dataset.

- 596

597 The quality-filtered reads from all datasets were mapped to the corresponding databases using Diamond, with the exception that no merging of read pairs was attempted, because the chances 598

of finding overlapping reads were too low (with a read length of 100 bp and insert size of 300 bp 599

(Tara Oceans). Sequences for SSU/LSU miTags were extracted from these metagenomics 600

datasets as described above. ARG abundance matrices were also obtained from the Tara Oceans 601

- project based on the published gene catalogues annotated using a similar approach as in the 602 current study. 603
- 604

Gene catalogue construction 605

To create a gene catalogue, we first searched for complete reference genes that matched to read 606 pairs in our collection using bowtie2⁵⁶ with the options "--no-unal --end-to-end". The resulting 607 bam files were sorted and indexed using samtools 1.3.1⁵⁷ and the jgi_summarize_bam_contig_depths provided with MetaBat⁵⁸ was used to create a depth profile 608

609

of genes from the reference databases that were covered with \geq 95% nucleotide identity. This cut-off is commonly used in constructing gene catalogues^{13,59} and chosen to delineate genes 610

611

belonging to the same species. Using the coverage information, we extracted all genes that had at 612 least 200bp with $>1\times$ coverage by reads from our topsoil metagenomes. The reference databases

613 614

included an ocean microbial gene catalogue¹³, a gut microbial gene catalogue¹², as well as all genes extracted from 25,038 published bacterial genomes²⁶. Altogether 273,723 and 2,376 and 615

8,642 genes from proGenomes, IGC and Tara database, respectively, could be matched to soil 616

reads and were used in the gene catalogue. 617

618

The majority of genes in our catalogue were assembled from the topsoil samples presented here. 619 To reduce the likelihood of chimeric reads, each sample was assembled separately using Spades

620 3.7-0 (development version obtained from the authors)⁶⁰ in metagenomic mode with the 621

parameters "--only-assembler -m 500 --meta -k 21,33,67,111,127". Only sdm⁴² filtered paired 622

reads were used in the assembly, with the same read filtering parameters as described above. 623

Resulting assemblies had an average N50 of 469 bases (total of all assemblies 21,538 MBp). The 624

low N50 reflects difficulties in the assembly of soil metagenomes, most likely reflecting the vast 625

microbial genetic diversity of these ecosystems. We further de novo assembled reads from two 626

other deep sequencing soil⁶¹ and sediment studies⁶², using the same procedure and parameters, 627

except that the Spades parameter "-k 21,33,67,77" was adjusted to a shorter read length. 628 Furthermore, we included publicly available data from the European Nucleotide Archive (ENA). 629 ENA was gueried to identify all projects with publicly available metagenomes and whose 630 metadata contained the keyword "soil". The initial set of hits was then manually curated to select 631 relevant project/samples that were assembled as described above. Additionally, we integrated 632 gene predictions from soil metagenomes downloaded from MG-RAST⁶³ (Supplementary Table 633 1). Assembly was not attempted for these samples due to the absence of paired end reads, and 634 relatively low read depth; rather, only long reads or assemblies directly uploaded to MG-RAST 635 with \geq 400bp length were retained. Therefore, only scaffolds and long reads, with at least 400 bp 636 length, were used for analysis. On these filtered sequences genes were de novo predicted using 637 prodigal 2.6.1⁶⁴ in metagenomic mode. Finally, we merged the predicted genes from assemblies, 638 long reads, gene catalogues and references genomes to construct a comprehensive soil gene 639 catalogue. 640

641

Thus, 53,294,555,100 reads were processed, of which 31,015,827,636 (58,20%) passed our 642 stringent quality control. The initial gene set predicted on the soil assemblies and long reads was 643 separated into 17,114,295 complete genes and 111,875,596 incomplete genes. A non-redundant 644 gene catalogue was built by comparing all genes to each other. This operation was performed 645 initially in amino-acid space using DIAMOND⁴⁴. Subsequently, any reported hits were checked 646 in nucleotide space. Any gene that covered at least 90% of another one (with at least 95% 647 identity over the covered area) was considered to be a potential representative of it (genes are 648 also potential representatives of themselves). The final set was chosen by greedily picking the 649 650 genes which are representative of the highest number of input genes until all genes in the original input have at least one representative in the output. This resulted in a gene catalogue with a total 651 of 159,907,547 non-redundant genes at 95% nucleotide identity cut-off. We mapped reads from 652 our experiment on the gene catalogue with bwa^{65} , requiring >45 nt overlap and >95% identity. 653 The average mapping rate was $26.2 \pm 7.4\%$. Although the gene catalogue is an invaluable 654 resource for future explorations of the soil microbiome, we decided to rely on using the direct 655 mapping approach to gene functional composition, due to the low overall mapping rate. Further, 656 using minimap2⁶⁶ to find genes at 95% similarity threshold, we compared the soil gene catalogue with the Tara Oceans gene catalogue¹³, human gut gene catalogue¹² and the proGenomes 657 658 prokaryotic database²⁶. The gene catalogue nucleotide and amino acid sequences and abundance 659

660 matrix estimates from rtk^{67} have been deposited at <u>http://vm-</u>

661 <u>lux.embl.de/~hildebra/Soil_gene_cat/</u>.

662

663 CARD ARG abundance estimation

 $\begin{array}{rcl} 664 & CARD \ abundances in topsoil samples were estimated by annotating the soil gene catalogue using \\ 665 & a \ DIAMOND \ search \ of \ the \ predicted \ amino \ acid \ sequences \ against \ the \ CARD \ database \ and \\ 666 & filtering \ hits \ to \ the \ specified \ bit-score \ cut-offs \ in \ the \ CARD \ database. \ Based \ on \ the \ gene \\ 667 & abundances \ in \ each \ sample, \ we \ estimated \ the \ abundance \ of \ different \ CARD \ categories \ per \\ 668 & metagenomic \ sample. \ Despite \ qualitative \ similarities \ in \ overall \ trends \ of \ ARDB \ and \ CARD \\ 669 & abundance \ matrices, \ CARD \ abundance \ estimation \ is \ limited \ by \ being \ based \ on \ the \ gene \\ 670 & catalogue \ (only \ a \ 26.2\pm7.4\% \ of \ all \ metagenomic \ reads \ could \ be \ mapped \ to \ the \ gene \ catalogue). \end{array}$

671

672 **Processing of metabarcoding sequence data**

The LotuS pipeline⁴² was used for bacterial 16S rRNA amplicon sequence processing. Reads 673 were demultiplexed with modified quality-filtering settings for MiSeq reads, increasing strictness 674 to avoid false positive OTUs. These modified options were the requirement of correctly detected 675 forward 16S primer, trimming of reads after an accumulated error of 1 and rejecting reads below 676 28 average quality or, exceeding an estimated accumulated error >2.5 with a probability of 677 $\geq 0.01^{41}$. Further, we required each unique read (reads preclustered at 100% identity) to be 678 present 8 or more times in at least one sample, 4 or more times in at least two samples, or three 679 or more times in at least three samples. In total 27,883,607 read pairs were quality-filtered and 680 clustered with uparse⁶⁸ at 97% identity. Chimeric OTUs were detected and removed based on 681 both reference-based and *de novo* chimera checking algorithms, using the RDP reference 682 database (http://drive5.com/uchime/rdp_gold.fa) in uchime⁶⁸, resulting in 13,070,436 high-683 quality read pairs to generate and estimate the abundance of bacterial OTUs The seed sequence 684 for each OTU cluster was selected from all read pairs assigned to that OTU, selecting the read 685 pair with the highest overall quality and closest to the OTU centroid. Selected OTU seed read 686 pairs were merged with FLASH⁴³ and a taxonomic identity was assigned to each OTU by 687 aligning full-length sequences with lambda⁵⁵ to the SILVA v123 database⁶⁹ and the LotuS least 688 common ancestor (LCA) algorithm. This was performed using the following LotuS command 689 line options: "-p miSeq -derepMin 8:1,4:2,3:3 -simBasedTaxo 2 -refDB SLV -thr 8". OTU 690 abundances per sample were summed to class and phylum level per sample, according to their 691 taxonomic classification, to obtain taxa abundance matrices. However, the choice of clustering 692 693 method (e.g. Swarm) and identity threshold had little effect on retrieved OTU richness (comparison with 99% threshold: r=0.977, p<10⁻¹⁵; comparison with Swarm clustering: r=0.979 694 $p < 10^{-15}$). 695

696

For eukaryotic 18S rRNA genes, we used the same options in LotuS, except that reads were 697 698 rejected if they did not occur at least six times each in a minimum of two samples or at least four times each in a minimum of three samples. This was done to account for lower sequencing depth 699 in 18S rRNA compared to 16S rRNA dataset. Further, the database to annotate fungal taxonomy 700 was extended to include general annotations of SILVA and information from unicellular 701 eukaryotes (PR2 database⁷⁰). Of 7,462,813 reads, 2,890,093 passed quality filtering. The fungal ITS metabarcoding dataset¹⁶ was downloaded and used in addition to 18S data in specific 702 703 analyses, such as finding associated fungal OTUs with ARG relative abundance. The resulting 704 taxon abundance matrix was further filtered to remove sequences of chloroplast origin for all 705 three metabarcoding experiments. 706

707

Full-length sequences representing OTUs were aligned using the SILVA reference alignment as
 a template in mothur⁷¹. A phylogenetic tree was constructed using FastTree2⁷² with the
 maximum-likelihood method using default settings. This program uses the Jukes-Cantor models
 to correct for multiple substitutions.

712

713 Parametrization and validation of metagenomics approach

Although we used state-of-art molecular approaches, there are several potential limitations

- regarding our analyses related to the used technologies. All metagenomics and amplicon-based
- analysis are affected by taxonomic biases in sequence databases, while (PCR-free) miTag as well
- as amplicon sequencing are biased due to differential ribosomal gene copy number across
- taxonomic groups. Amplicon-based metabarcoding, specifically, is affected by both primer PCR

artefacts and PCR biases that may affect estimates of absolute organism abundance. These biases

- are inherent to all metagenomics and metabarcoding studies. However, all these biases affect
- different samples equally (same rRNA gene copy numbers, same PCR biases per species, same

database bias per taxa) and thus we estimate that our results are robust to these methodological
 shortcomings. Shotgun-based metagenomics is affected by reference bias, in which human

723 shoreonings: shoreonings: shoreonings is affected by reference bias, in which human 724 pathogens or Proteobacteria are overrepresented. The necessity for lenient thresholds becomes

obvious from annotating phylogenetic profiles with MetaPhlAn2⁷³ using standard parameters:

while we observed that most fungal phyla are present abundantly in our samples. MetaPhlAn2

detected Ascomycota only in 2 out of 189 samples. In 48 out of 189 samples, no organism

728 (bacteria/archaea/eukaryotes) was detected, and the most abundant phylum was Proteobacteria

(55%). Since these results are clearly deviating from our miTag, 16S, 18S and ITS based
 analysis, specific database cut-off thresholds were required for this project.

731

732 To optimize the analysis pipeline and identify suitable e-values for filtering blastx results, we used metagenomic simulations of four reference genomes where CAZy assignments in the CAZy 733 database were available. Simulated reads were created as 250 bp paired reads with 400 bp insert 734 at differing sequence abundances from the four reference genomes in each simulated 735 metagenome, using iMessi⁷⁴. For this simulated dataset, we used the pipeline described above to 736 derive CAZy functional profiles. We found that querying short reads processed as above against 737 databases results in the retrieval of most genes at relative abundances consistent with 738 expectations based on the reference genomes at e-value $< 1e^{-9}$ (r=0.95±0.01, p<0.001). Further, 739 we simulated 200 metagenomes from 18 bacterial genomes, five bacterial plasmids, one fungal 740 741 mitochondrion and two fungal genomes at differing relative proportions in each of these simulated metagenomes (Supplementary Table 11). We subsequently simulated 1,000,000 reads 742 of 250 bp and 400 bp insert size using iMessi, and mapped these against reference databases and 743 retained hits that fulfilled the following arbitrary criteria (used in all subsequent analyses): e-744 value cut-off of e^{-9} , alignment length ≥ 20 amino acids, and similarity $\ge 50\%$ amino acids to the 745 target sequence. From these, we generated functional profiles and found a strong correlation of 746 simulated to expected functional metagenomic composition based on mixed fungal and bacterial 747 genomes (r=0.94±0.05, p<0.001). 748

749

750 Estimating fungal antibiotics production

751 We also specifically screened for fungal gene clusters directly associated with antibiotic activity,

based on a compiled database of MIBiG (Minimum Information about a Biosynthetic Gene

cluster, https://mibig.secondarymetabolites.org) repository entries that describe gene clusters for

- which the products have been shown experimentally to display antimicrobial activities T_{1}
- 755 (Supplementary Table 12). To extend the range of genes that can be associated with the
- validated, antibiotics producing, MiBIG protein domains, we downloaded all published non-
- redundant fungal genomes deposited in JGI (Supplementary Table 13) as well as all non-
- redundant fungal genes deposited in NCBI. The set of MiBIG, and fungal derived genes was
 screened with custom HMMs for domains from secondary metabolite production (specifically)
- these were dmat, AMP-binding, Condensation, PKS KS and Terpene synthesis domains). All
- ⁷⁶¹ identified domains were aligned together with the MiBIG domains using Clustal Omega⁷⁵ and a
- 762 tree was constructed with FastTree2. Phylogenetic trees were rooted to midpoint and
- automatically scanned to identify highly supported clades (aLRT branch support ≥ 0.99) where
- antibiotic producing MiBIG domains were monophyletically grouped. The average nucleotide

⁷⁶⁵ identity within each such group was subsequently used as identity cut-off in the mapping step.

All metagenomic reads were mapped with diamond in blastx mode to the newly created

database, using before-mentioned sequence identity cut-offs and rejecting domains of reads thatwere mapping to bacterial NOGs.

769

770 Statistical analyses

771 Data normalization and diversity estimates

All statistical analyses were performed using specific packages in R (version 3.3.2) unless otherwise noted. Diversity parameters were estimated from OTU and functional gene matrices that were rarefied to an equal number per sample to reduce the effect of variation in sequencing depth using the function *rrarefy* in vegan (version 2.2.1)⁷⁶. ARG matrices were normalized by the total number of merged and singleton reads. Total abundance of ARGs per sample was estimated by summing the abundance of all individual ARGs per sample. ARG diversity

- measures indicate the variety and their proportions produced.
- 779

780 From the rarefied matrices we calculated OTU, OG and CAZyme gene richness (function

specnumber) and diversity (function *diversity*, based on the Inverse Simpson index). The latter

782 measure accounts for both richness and evenness, and it gives more weight to abundant groups 783 compared to Shannon Index. Our results were robust to choice of index, and the various diversity

783 compared to shallfor index. Our results were robust to choice of index, and the various diversity 784 indices highly correlated in the present dataset (e.g. bacterial taxonomic diversities calculated

using Inverse Simpson versus using Shannon diversity were highly correlated: r=0.888, $p<10^{-15}$;

for a comparison of richness and diversity trends, see Extended Data Figure 2b,c). Since

evenness and richness were highly correlated in all datasets, we report the results based on

diversity index that represent both richness and evenness. The rarefaction process was repeated

for calculating taxonomic and gene functional diversity and richness based on the average of 100rarefied datasets.

791

Phylogenetic diversity was calculated based on Faith's Phylogenetic Diversity (PD) metric in
Picante package of R⁷⁷. In addition, to assess phylogenetic clustering and overdispersion, Nearest
Relative Index (NRI) and Nearest Taxon Index (NTI) were calculated in Picante. Although both
measures are closely related, NRI is more sensitive to phylogenetic diversity at deep nodes,

taxon labels (100 times) was used to randomize phylogenetic relationships among OTUs.

whereas NTI is more sensitive to phylogenetic clustering towards tips. A null model of shuffling

796 797

798

799 Correlating environmental parameters to taxa and functions

To identify the main determinants of taxonomic and gene functional composition or diversity 800 and relative abundance of phyla/classes, we used a series of statistical tests. We included all 801 prominent environmental variables that we expected to have a significant effect on microbial 802 diversity based on previous studies, and which were feasible to collect. These included soil pH, 803 carbon and nutrient levels and factors that can affect these, such as fire, assuming soil as the 804 major resource for microbial nutrition. We also included isotope ratios of nitrogen ($\partial^{15}N$) and 805 carbon $(\partial^{13}C)$ as these provide principal components for carbon and nitrogen cycling. To avoid 806 overfitting and to ensure model simplicity, we excluded the variables that had no significant 807 impact on fungal or bacterial diversity, such as altitude, age of vegetation, plant diversity and 808 community (the first two PCA axes of Plant community variation at both genus and family level) 809

and basal areas of trees. Thus, for univariate regression modelling, 16 variables (Supplementary
 Table 14) were included.

812

To understand, which factors explain the OG- and OTU-based community composition, variable 813 selection was performed in the *Forward.sel* function of Packfor (version 0.0-8/r109)⁷⁸ according 814 to the coefficient of determination (threshold, $r^2=0.01$). All functional and taxonomic 815 compositional matrices were transformed using Hellinger transformation prior to statistical 816 analysis. Further, Mantel tests and partial Mantel tests were used to test the effects of 817 geographical vs. environmental distances on OTU and OG compositional similarity as 818 implemented in vegan. Mantel tests allow testing the correlation of two distance matrices, 819 whereas partial Mantel tests are similar but also control for variation in a third distance matrix. In 820 821 our analysis, we controlled for the effect of geographic distance while testing the correlation of environmental variation and functional or taxonomic composition variation. The importance of 822 823 biome type in explaining functional gene and taxonomic composition was tested in Permutational Multivariate Analysis of Variance (PERMANOVA) using the Adonis function of 824 vegan (using 10^3 permutation for calculating pseudo-F test statistic and its statistical 825 significance). For constructing OG and OTU distance matrices, the Bray-Curtis dissimilarity was 826 calculated between each pair of samples. Great-circle distance was used to calculate a geographic 827 distance matrix between samples based on geographical coordinates. This test compares the 828 intragroup distances to intergroup distances in a permutation scheme and from this assesses 829 significance. PERMANOVA post-hoc p-values were corrected for multiple testing using the 830 Benjamini–Hochberg correction. We visualized taxonomic (OTU) and functional (OG) 831 832 composition of bacteria using global nonmetric multidimensional scaling (GNMDS) in vegan with the following options: two dimensions, initial configurations = 100, maximum iterations = 833 200, and minimum stress improvement in each iteration $=10^{-7}$. The main environmental drivers 834 835 of the relative abundance of major taxonomic groups and main functional categories were recovered by random forest (RF) analysis⁷⁹ using the R-package randomForest (version 4.6-10). 836 837 To examine latitudinal gradients of diversity at phylum level (Figure 2), the diversity of OTUs 838 assigned to each phylum was calculated based on Inverse Simpson index. Diversity values were 839 modelled in response to environmental variables and predicted values were extracted, which 840 were used in a clustering and bootstrapping analysis to depict the similarities of phyla 841 environmental associations using pyclust (version 1.3-2)⁸⁰ with 1000 iterations. To model 842 latitudinal gradients and environmental associations of diversity and biomass (Figure 1, 843 844 Extended Data Figure 3), we compared the goodness of fit estimates between first and second order polynomial models based on the corrected Akaike information criterion (AICc) using 845 analysis of variance (ANOVA). AICc reflects both goodness of fit and parsimony of the models. 846 847 848 For univariate regression modelling of diversity and biomass measures, ordinary least squares (OLS) or generalized least squares (GLS) regression models were employed depending on the 849 importance of the spatial component. The model variance structure (Gaussian, exponential. 850 spherical and linear) was evaluated based on AICc. Following selection of variance structure, 851 variables were combined in a set of models with specified variance structure (i.e. number of 852 tested models: 2^{number of variables}). The resulting models were sorted according to AICc values to 853

reveal the best model. Lists of the 5 best-fitting models for each response variable are given in

855 Supplementary Table 4. Prior to model selection, all variables were evaluated for linearity,

normality, and multicollinearity (excluded if the variance inflation factor was >5). The degree of

polynomial functions (linear, quadratic, cubic) was chosen based on the lowest AIC values.

Because of non-linear relationships with response variables, a quadratic term for pH was also

859 included in the model selection procedure. The accuracy of the final models was evaluated using

10-fold 'leave-one-out' cross-validation. For this, we used 1000 randomly sampled 90%-data

- subsets for model training and predicting the withheld data. To minimize biases due to the
- partitioning of the data and potential overfitting, the average of 1000 resulting determination

coefficients are reported as cross-validated $r^2(r^2cv)$ for each regression model.

864

865 Correlating biotic interactions to taxa and functions

To test the associations of biotic variables on ARG relative abundance, we used a sparse partial 866 least squares (sPLS) analysis, which reduces dimensionality by projecting predictor variables 867 onto latent components to identify the 16S/18S lineages (phyla/classes) and the ITS OTUs most 868 strongly associated with ARG relative abundance, as implemented in the mixOmics (version 5.0-869 4)⁸¹ package. ARG composition and taxonomic community matrices (miTags classes/phyla and 870 ITS OTUs) were normalized by library size using Hellinger transformation. Significance of 871 associations was examined by bootstrap tests of subsets of each dataset. We subsequently used 872 partial least squares (PLS) analysis to predict ARG relative abundance based on significantly 873 correlated lineages, which allows the dimensionality of multivariate data to be reduced into PLS 874 components. Optimal numbers of PLS components for prediction of the relative ARG abundance 875 were selected based on leave-one-out cross-validation. To confirm the results of PLS analysis, 876 we further used a cross-validated LASSO model to simultaneously perform variable selection 877 and model fitting, as implemented in glmnet (version 2.0-2)⁸². First the lambda shrinkage 878 parameter was determined from a cross-validated lasso-penalized logistic regression classifier. 879 Using this shrinkage parameter, a new logistic regression classifier was fit to the data to predict 880

881 ARG relative abundance.

882

To further test direct and indirect effects of geographic and environmental variables on microbial 883 distributions, we built SEM models in the AMOS software (SPSS, Chicago, IL) by including 884 predictors of the best GLS model. In *a priori* models, all indirect and direct links between 885 variables were established based on their pairwise correlations. We subsequently removed non-886 significant links and variables or created new links between error terms until a significant model 887 fit was achieved. Goodness of fit was assessed based on Chi-square test to evaluate the 888 difference between observed and estimated by model covariance matrices (non-significant value 889 indicates that the model fits the observed data). We also used Root Mean Square Error of 890 Approximation (RMSEA) and *PCLOSE* (p-value for test of close fit) to assess the discrepancy 891 between the observed data and model per degree of freedom, which is less sensitive to sample 892 size compared to chi-square test (RMSEA < 0.08 and PCLOSE > 0.05 show a good fit). 893 Observed correlations between diversity and environmental values can serve as the first step 894 towards understanding the structure and function of global topsoil microbiome; however, they 895 are not proof of causations and mechanism. Despite the fact that we used SEM modelling to infer 896 indirect links, we cannot preclude the possibility of other biotic or soil variables confounded with 897 climate variables that we did not include in our models. Further laboratory experiments may 898 enable to address causality of relationships reported in this study. 899 900

901 Differences between univariate variables such as taxonomic and functional richness were tested using a non-parametric Wilcoxon rank-sum test, with Benjamini-Hochberg multiple testing 902 correction. Post-hoc statistical testing for significant differences between all combinations of two 903 groups was conducted only for taxa with p<0.2 in the Kruskal-Wallis test. For this, wilcoxon 904 rank-sum tests were calculated for all possible group combinations and corrected for multiple 905

- testing using Benjamini-Hochberg multiple testing correction. 906
- 907

Geographic coordinates were plotted on a world map transformed to a Winkler2 projection, 908 using the maptools (version 0.8-36) package⁸³. 909

910

Limitations of statistical modelling on a global scale 911

Although we performed cross-validations to test the accuracy of most of our statistical models, 912 predictions might be limited by the vast diversity in soil microbiomes. For example, strong local 913 914 variation in soil pH may lead to deviation from general patterns, which is a common limitation in environmental sciences. Given the large spatial scale and strong environmental gradient in our 915 sampling design, and long-term persistence of DNA in soil⁸⁴, seasonal variation in soils is 916 expected to have a minor impact⁸⁵ (in contrast to ocean). In addition, the vast majority of our 917 samples were collected during growing season, further reducing possible seasonal biases. We 918 nevertheless tested the effect of sampling month and seasons and found no significant effect of 919 seasonality on diversity indices (P>0.05). We also compared the effect of seasons and years in a 920 time series study in two of our sites, which revealed no seasonal effects on richness and 921 composition (unpublished data). In particular, the relationship between bacterial phylogenetic 922 diversity and pH, are strongly consistent with studies performed at the local to continental scales 923 and within a single season 6,7,86 , which indicates the robustness of our results. Nonetheless, 924 validation of the proposed models needs to be performed by other researchers with extended data 925 or an independent dataset, particularly by including samples from under-sampled regions 926 (Extended Data Figure 1a) and from different seasons (to account for seasonality). For example, 927 there were some under-sampled regions in our dataset (e.g. North Asia) lowering precision of our 928 models for those regions. Unfortunately, there are no published global datasets with comparable 929 sampling protocols used that could be directly compared and used for model validation, and we 930 encourage future studies that will make this possible. 931

932 933

934 Data availability All metagenomics and metabarcoding sequences have been deposited in the European Bioinformatics Institute-Sequence Read Archive database, under accession number 935 936 PRJEB24121 (ERP105926): Estonian forest and grassland topsoil samples; PRJEB19856 (ERP021922): 16S metabarcoding data of global soil samples; PRJEB19855 (ERP021921): 18S 937 metabarcoding data of global soil samples; PRJEB18701 (ERP020652): Global analysis of soil 938 microbiomes. The soil gene catalogue and dataset are available at http://vm-939

lux.embl.de/~hildebra/Soil gene cat/. The Tara Oceans data are available at http://ocean-940

microbiome.embl.de/companion.html. All other data that support the findings of this study are 941

available from the corresponding authors upon request. 942

943

 http://psbweb05.psb.ugent.be/lotus/. The pipeline to process shotgun metagenomic samples is available under https://github.com/hildebra/MATAFILER and https://github.com/hildebra/Rarefaction. 31 Tedersoo, L. <i>et al.</i> Shotgun metagenomes and multiple primer pair-barcode combinatio of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). 32 Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). 					
 available under https://github.com/hildebra/MATAFILER and https://github.com/hildebra/Rarefaction. 31 Tedersoo, L. <i>et al.</i> Shotgun metagenomes and multiple primer pair-barcode combinatio of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). 32 Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). 33 Weltern West of human disperied bacterial 16C. DNA and 1144 (2) and for a bit in 	http://psbweb05.psb.ugent.be/lotus/. The pipeline to process shotgun metagenomic samples is				
 947 https://github.com/hildebra/Rarefaction. 948 949 950 31 Tedersoo, L. <i>et al.</i> Shotgun metagenomes and multiple primer pair-barcode combinatio of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). 953 32 Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). 955 32 Weltern West of human dispersivel 16C and the diversity of the second sec	available under <u>https://github.com/hildebra/MATAFILER</u> and				
 948 949 950 31 Tedersoo, L. <i>et al.</i> Shotgun metagenomes and multiple primer pair-barcode combination of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). 953 32 Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). 955 22 Wettern West of human distribution (14 - 11/4 f) and for the formation of the sequences of the sequen					
 949 950 31 Tedersoo, L. <i>et al.</i> Shotgun metagenomes and multiple primer pair-barcode combinatio 951 of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). 953 32 Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of 954 sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). 955 22 					
 Tedersoo, L. <i>et al.</i> Shotgun metagenomes and multiple primer pair-barcode combinatio of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). 					
 of amplicons reveal biases in metabarcoding analyses of fungi. <i>MycoKeys</i> 10, 1-43 (2015). 32 Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011). Weltern W. et al. Improved hetericial 16C, DNA (14, 11, 14, 14, 14, 14, 14, 14, 14, 14,	ns				
 (2015). <					
95332Caporaso, J. G. <i>et al.</i> Global patterns of 16S rRNA diversity at a depth of millions of954sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108, 4516-4522 (2011).95522					
954 sequences per sample. <i>Proc. Natl. Acad. Sci. USA</i> 108 , 4516-4522 (2011).					
0.55 22 Welters West at Language the stand 14 (C DNA (374 1374 5) 16 1)					
955 55 waiters, w. <i>et al.</i> Improved bacterial 16S rKNA gene (V4 and V4-5) and fungal interna	ıl				
956 transcribed spacer marker gene primers for microbial community surveys. <i>mSystems</i> 1 ,					
957 e00009-00015 (2016).					
958 34 Amaral-Zettler, L. A., McCliment, E. A., Ducklow, H. W. & Huse, S. M. A method for					
studying protistan diversity using massively parallel sequencing of V9 hypervariable					
regions of small-subunit ribosomal RNA genes. <i>PLoS ONE</i> 4 , e6372 (2009).					
961 35 Frostegård, Å., Tunlid, A. & Bååth, E. Use and misuse of PLFA measurements in soils.					
962 Soil Biol. Biochem. 43 , 1621-1625 (2011).					
963 36 van Aarle, I. M. & Olsson, P. A. Fungal lipid accumulation and development of myceli	al				
964 structures by two arbuscular mycorrhizal fungi. <i>Appl. Environ. Microbiol.</i> 69 , 6762-676	57				
965 (2003).					
966 37 Frostegård, A., Tunlid, A. & Bååth, E. Phospholipid fatty acid composition, biomass, at	ıd				
967 activity of microbial communities from two soil types experimentally exposed to					
968 different heavy metals. Appl. Environ. Microbiol. 59 , 3605-3617 (1993).					
969 38 Prosser, J. I. Dispersing misconceptions and identifying opportunities for the use					
970 of omics' in soil microbial ecology. <i>Nat. Rev. Microbiol.</i> 13 , 439-446 (2015).					
971 39 Salles, J. F., Le Roux, X. & Poly, F. Relating phylogenetic and functional diversity					
9/2 among denitrifiers and quantifying their capacity to predict community functioning.					
9/3 Front. Microbiol. 3 (2012).					
40 Invedi, P. <i>et al.</i> Microbial regulation of the soil carbon cycle: evidence from gene-					
975 enzyme relationships. <i>ISME J.</i> 10, 2595-2004 (2010).					
9/6 41 Fuence-Sanchez, F., Aguirre, J. & Fairlo, V. A novel conceptual approach to read-intern in high throughout amplicon sequencing studiog. Nucleic Acids Res. 44, e40, e40 (2016)	.ig				
9/7 In high-unoughput amplicon sequencing studies. Nucleic Actus Res. 44, e40-e40 (2010) 079 42 Hildebrand E. Tadeo, P. Voiet, A. V. Bork, D. & Paes, I. Lotus: an afficient and use	1. r				
⁹⁷⁸ 42 Indebiand, F., Tadeo, K., Volgi, A. T., Dork, T. & Kaes, J. Lotus. an efficient and use	1-				
Magoc T & Salzberg S I ELASH: fast length adjustment of short reads to improve					
agenome assemblies <i>Bioinformatics</i> 27 2057 2063 (2011)					
981 genome assemblies. <i>Diologormatics</i> 27, 2951-2905 (2011).					
DIAMOND Nat Methods 12, 59, 60 (2015)					
12, 59-00 (2013).					
Microbiol 14 117 (2014)					
986 46 Fawal N <i>et al</i> PeroxiBase: a database for large-scale evolutionary analysis of					
987 neroxidases Nucleic Acids Res 41 D441-D444 (2012)					
988 47 Kanehisa M & Goto S KEGG kyoto encyclonedia of genes and genomes <i>Nucleic</i>					
989 Acids Res. 28, 27-30 (2000).					

990	48	Cantarel, B. L. et al. The Carbohydrate-Active EnZymes database (CAZy): an expert
991		resource for glycogenomics. Nucleic Acids Res. 37, D233-D238 (2009).
992	49	Cantarel, B. L., Lombard, V. & Henrissat, B. Complex carbohydrate utilization by the
993		healthy human microbiome. PLoS ONE 7, e28742 (2012).
994	50	Cardenas, E. et al. Forest harvesting reduces the soil metagenomic potential for biomass
995		decomposition. ISME J. 9, 2465-2476 (2015).
996	51	Forslund, K. et al. Country-specific antibiotic use practices impact the human gut
997		resistome. Genome Res. 23, 1163-1169 (2013).
998	52	Gibson, M. K., Forsberg, K. J. & Dantas, G. Improved annotation of antibiotic resistance
999		determinants reveals microbial resistomes cluster by ecology. ISME J. 9, 207-216 (2015).
1000	53	McArthur, A. G. et al. The comprehensive antibiotic resistance database. Antimicrob.
1001		Agents Chemother. 57, 3348-3357 (2013).
1002	54	Kopylova, E., Noé, L. & Touzet, H. SortMeRNA: fast and accurate filtering of ribosomal
1003		RNAs in metatranscriptomic data. <i>Bioinformatics</i> 28, 3211-3217 (2012).
1004	55	Hauswedell, H., Singer, J. & Reinert, K. Lambda: the local aligner for massive biological
1005		data. Bioinformatics 30, i349-i355 (2014).
1006	56	Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nat. Methods
1007		9 , 357-359 (2012).
1008	57	Li, H. et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25,
1009		2078-2079 (2009).
1010	58	Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately
1011		reconstructing single genomes from complex microbial communities. PeerJ 3, e1165
1012		(2015).
1013	59	Qin, J. et al. A human gut microbial gene catalogue established by metagenomic
1014		sequencing. Nature 464, 59 (2010).
1015	60	Bankevich, A. et al. SPAdes: a new genome assembly algorithm and its applications to
1016		single-cell sequencing. J. Comput. Biol. 19, 455-477 (2012).
1017	61	Howe, A. C. et al. Tackling soil diversity with the assembly of large, complex
1018		metagenomes. Proc. Natl. Acad. Sci. USA 111, 4904-4909 (2014).
1019	62	Sharon, I. et al. Accurate, multi-kb reads resolve complex populations and detect rare
1020		microorganisms. Genome Res. 25, 534-543 (2015).
1021	63	Meyer, F. et al. The metagenomics RAST server-a public resource for the automatic
1022		phylogenetic and functional analysis of metagenomes. BMC Bioinformatics 9, 386
1023		(2008).
1024	64	Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site
1025		identification. BMC Bioinformatics 11, 119 (2010).
1026	65	Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler
1027		transform. Bioinformatics 25, 1754-1760 (2009).
1028	66	Li, H. Minimap2: fast pairwise alignment for long DNA sequences. arXiv:1708.01492
1029		(2017).
1030	67	Saary, P., Forslund, K., Bork, P. & Hildebrand, F. RTK: efficient rarefaction analysis of
1031		large datasets. Bioinformatics 33, 2594-2595 (2017).
1032	68	Edgar, R. C. UPARSE: highly accurate OTU sequences from microbial amplicon reads.
1033		<i>Nat. Methods</i> 10 (2013).

1034	69	Pruesse, E. et al. SILVA: a comprehensive online resource for quality checked and
1035		aligned ribosomal RNA sequence data compatible with ARB. Nucleic Acids Res. 35
1036		(2007).
1037	70	Guillou, L. et al. The Protist Ribosomal Reference database (PR2): a catalog of
1038		unicellular eukaryote small sub-unit rRNA sequences with curated taxonomy. <i>Nucleic</i>
1039		Acids Res., gks1160 (2012).
1040	71	Schloss, P. D. et al. Introducing mothur: Open-Source, Platform-Independent,
1041		Community-Supported Software for Describing and Comparing Microbial Communities.
1042		Appl. Environ. Microbiol. 75, 7537-7541 (2009).
1043	72	Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-
1044		likelihood trees for large alignments. PLoS ONE 5 (2010).
1045	73	Truong, D. T. et al. MetaPhlAn2 for enhanced metagenomic taxonomic profiling. Nat.
1046		Methods 12, 902 (2015).
1047	74	Mende, D. R. et al. Assessment of metagenomic assembly using simulated next
1048		generation sequencing data. PLoS ONE 7, e31386 (2012).
1049	75	Sievers, F. et al. Fast, scalable generation of high-quality protein multiple sequence
1050		alignments using Clustal Omega. Mol. Syst. Biol. 7 (2011).
1051	76	Oksanen, J. et al. Vegan: community ecology package. R Package (2015).
1052	77	Faith, D. P. Conservation evaluation and phylogenetic diversity. <i>Biol. Conserv.</i> 61
1053		(1992).
1054	78	Dray, S., Blanchet, F. & Legendre, P. packfor: Forward selection with permutation. <i>R</i>
1055		<i>Package</i> (2013).
1056	79	Breiman, L. Random forests. Mach. learn. 45, 5-32 (2001).
1057	80	Suzuki, R. & Shimodaira, H. pvclust: hierarchical clustering with p-values via multiscale
1058		bootstrap resampling. <i>R Package</i> (2009).
1059	81	Lê Cao, K. <i>et al.</i> mixOmics: Omics Data Integration Project. <i>R Package</i> (2017).
1060	82	Friedman, J. & Hastie, T. glmnet: lasso and elastic-net regularized generalized linear
1061	~~	models. R Package (2015).
1062	83	Bivand, R. <i>et al.</i> Maptools: Tools for reading and handling spatial objects. <i>R package</i>
1063	0.4	(2015).
1064	84	Carini, P. <i>et al.</i> Relic DNA is abundant in soil and obscures estimates of soil microbial
1065	05	diversity. Nal. Microbiol. 2, 10242 (2010). Žifšálvová L. Větrovolví, T. Howa, A. & Doldnion, D. Mionohiol activity in forest soil
1066	83	Zilcakova, L., Vetrovsky, I., Howe, A. & Baldrian, P. Microbial activity in lorest soli
1067		Microbiol 19, 288, 201 (2016)
1068	96	Micropiol. 18, 288-501 (2010). Eigrar N. & Laskson B. D. The diversity and his geography of soil hesterial
1069	80	communities Prog. Natl. Acad. Sci. USA 103, 626, 621 (2006)
1070		communues. <i>Froc. Natl. Acaa. Sci. USA</i> 105 , 626-651 (2006).
10/1		
1072		
1073	Fyton	dad Data laganda
1074	EAU	
1076	Exten	ded Data Figure 1 Distribution of tonsoil samples and diversity natterns of phyla a
1077	A mar	of samples used for metagenomic and metabarcoding analysis. Colours indicate biomes
1078	as indi	icated in the legend. Desert samples were only used in metabarcoding analysis and were
10,0	as ma	the second is the regent of the second of th

1079 excluded in comparative analysis of functional and taxonomic patterns. Black symbols refer to

samples from an independent soil dataset (145 topsoil samples; Supplementary Table 1) that

1081 were used for validation our results. b, Scatterplots showing the relationship between the
 1082 diversity of major microbial phyla (classes for Proteobacteria) and environmental variables

diversity of major microbial phyla (classes for Proteobacteria) and environmental variables
 across the global soil samples (n=197 biologically independent samples). Only regression lines

for significant relationships after Bonferroni correction are shown. Diversity was measured using

1085 Hellinger-transformed matrices based on Inverse Simpson Index. Latitude: absolute latitude;

1086 MAP: mean annual precipitation; MAT: mean annual temperature; C/N: carbon to nitrogen ratio. 1087

Extended Data Figure 2 | Contrasting microbial structure and function in major terrestrial 1088 **biomes.** a-d, The average total biomass (n=152 biologically independent samples) as well as 1089 richness, diversity and relative abundance (n=188 biologically independent samples) of fungi and 1090 bacteria across samples categorized into major terrestrial biomes, including tropical (moist and 1091 dry tropical forests and savannas), temperate (coniferous and deciduous forests, grasslands and 1092 1093 shrublands, and Mediterranean biomes) and boreal-arctic ecosystems: total biomass (a); richness (b); diversity (c); phylogenetic structure including Nearest relative index (NRI) and Nearest 1094 taxon index (NTI) (see Methods) (d). e-i, Relative abundance of major phyla (n=188 biologically 1095 independent samples) and functional categories (n=189 biologically independent samples) across 1096 biomes: bacterial phyla (classes for Proteobacteria) and archaea (e); fungal classes (f); functional 1097 categories of bacteria (g); functional categories of fungi (h); bacterial KEGG metabolic pathways 1098 (i). Biomass was measured based on phospholipid-derived fatty acids (PLFA) analysis (see 1099 1100 Methods). Different letters denote significant differences between groups (shown in the legend) at the 0.05 probability level based on Kruskal-Wallis test corrected for multiple testing. 1101 1102 Additional details for these comparisons are presented in Supplementary Table 14. Taxonomic and gene functional diversity indices were calculated based on Inverse Simpson Index. The 1103

1104 centre values and error bars represent mean and SD, respectively.

1105

1106 Extended Data Figure 3 | Significant decline of bacterial to fungal biomass ratio with

increasing latitude due to the joint effect of climate and soil fertility. a, The second order polynomial relationship of absolute latitude and the total biomass of bacteria (n=152 biologically

independent samples). **b**, The relationship of absolute latitude and the total biomass of bacteria $(n-1)^2$ biological

1110 The relationship of absolute latitude and the ratio of bacterial to fungal (B/F) biomass. **d-f**, The

relationship of B/F biomass ratio and mean annual precipitation (MAP), mean annual

temperature (MAT) and carbon to nitrogen ratio (C/N), as the main correlated environmental

1113 variables with B/F biomass ratio. Linear regression analysis (Pearson correlation) was used in **b-f**

1114 (n=152 biologically independent samples). **g**, Pairwise Spearman correlation matrix of biotic and

abiotic variables in soil. **h**, Direct and indirect relationships and directionality between variables determined from Best-fitting Structural Equation Model. Determination coefficients (R^2) are

given for biomass and diversity factors (see Supplementary Table 5 for more details). Goodness

of fit (see Methods): bacteria, Chi square=15.37, df=11, P=0.166; RMSEA=0.041,

1119 PCLOSE=0.573, n=189; fungi, Chi square=7.74, df=12, P=0.805; RMSEA=0.00,

1120 PCLOSE=0.970, n=189). Biomass (nmol/g) was measured based on phospholipid-derived fatty

acids (PLFA) analysis. pH, soil pH representing soil pH and its quadratic term; Ca, calcium; Mg,

1122 magnesium; P, phosphorous; K, potassium; C, carbon; N, nitrogen; d¹⁵N, nitrogen stable isotope

signature; d¹³C, carbon stable isotope signature; PET, potential of evapotranspiration; Fire, time

from the last fire disturbance; NPP, net primary productivity.

1125

1126 Extended Data Figure 4 | Environment has stronger effect on bacterial taxa and functions

1127 **than those of fungi**. Correlation and best random forest model for major taxonomic (**a** and **b**;

n=188 biologically independent samples) and functional (c and d; n=189 biologically

1129 independent samples) categories of bacteria (left column) and fungi (right column) in the global

- soil samples (n=189 biologically independent samples). **a**, Relative abundance of major 16S-
- based bacterial phyla (class for Proteobacteria). b, Relative abundance of ITS-based fungal
 classes. c-d, Major orthologous genes (OG) categories of bacteria (c) and fungi (d). For variable
- selection and estimating predictability, the random forest machine-learning algorithm was used.
- 1134 Circle size represents the variable importance, i.e. decrease in the prediction accuracy (estimated
- 1135 with out-of-bag cross-validation) as a result of permutation of a given variable. Colours represent
- 1136 Spearman correlations. pH, soil pH; Ca, calcium; Mg, magnesium; P, phosphorous; K,
- 1137 potassium; C, carbon; N, nitrogen; $d^{15}N$, nitrogen stable isotope signature; $d^{13}C$, carbon stable
- isotope signature; C/N, carbon to nitrogen ratio; Latitude, absolute latitude; MAP, mean annual
- 1139 precipitation; MAT, mean annual temperature; PET, potential of evapotranspiration; Fire, time 1140 from the last fire disturbance.
- 1140

1142 Extended Data Figure 5 | Niche differentiation between bacteria and fungi is likely related

1143to precipitation and soil pH. Contrasting effect of pH and mean annual precipitation (MAP) on1144bacterial (16S; left columns) and fungal (18S; right columns) taxonomic (n=188 biologically1145independent samples) and gene functional (n=189 biologically independent samples) diversity in

the global soil samples: **a**, **b**, Relationship of soil pH and taxonomic diversity of bacteria (**a**) and fungi (**b**); **c**, **d**, Relationship of soil pH and gene functional diversity of bacteria (**c**) and fungi (**d**);

- e, f, Relationship of MAP and taxonomic diversity of bacteria (e) and fungi (f); g, h,
- 1149 Relationship of MAP and gene functional diversity of bacteria (g) and fungi (h). Lines represent
- regression lines of best fit. The choice of degree of polynomial was determined by a goodness of fit (see Methods). Colours denote biomes as indicated in the legend. MAP: mean annual

precipitation. Taxonomic and gene functional diversity indices were calculated based on Inverse

1153 Simpson Index. **i-l**, Non-metric multidimensional scaling (NMDS) plots of trends in taxonomic

(16S and 18S-based datasets) and gene functional composition (OGs from metagenomes) of

bacteria (left column) and fungi (right column) based on Bray-Curtis dissimilarity. Taxonomic

1156 composition of bacteria (16S). **j**, Taxonomic composition of fungi (18S). **k**, Gene functional

- 1157 composition of bacteria. I, Gene functional composition of fungi. i, Colours denote biomes as
 1158 indicated in the legend. Vectors are the prominent environmental drivers fitted onto ordination.
- 1158 1159

1160 Extended Data Figure 6 | Fungal biomass is significantly related to the relative abundance

of antibiotic resistance genes (ARG). a, Increase in fungal biomass is related to ARG relative

- abundance. **b**, Bacterial biomass is unrelated to ARG relative abundance. **c**, ARG relative
- abundance: b, Bacterial biomass is unrelated to ARG relative abundance: c, ARG relative abundance is inversely correlated with Bacteria-to-Fungi biomass ratio. Biomass (nmol/g) was

measured based on Phospholipid Fatty Acids (PLFA) analysis (see Methods). Spearman

- 1165 correlation was used (n=152 biologically independent samples).
- 1166

1167 Extended Data Figure 7 | Topsoil and ocean bacterial phylogenetic diversity is negatively

- 1168 correlated with the abundance of antibiotic resistance genes. a, b, Spearman correlations
- between ARG relative abundance and bacterial phylogenetic diversity (Faith's index; see
- 1170 Methods) in soil (n=188 biologically independent samples). (a) and ocean (n=139 biologically
- 1171 independent samples). (b) at the global scale. Similar trends were observed for richness (r=-

1172 0.219, p=0.007 and r=-0.659, p<10⁻¹⁵) in soil and ocean, respectively). **c**, Global map of 1173 observed bacterial phylogenetic diversity (Faith's index; see Methods) at the sampled sites. Note 1174 that hotspots of bacterial diversity do not correspond to hotspots of ARG relative abundance (See

1175 Extended Data Figure 8).

1176

Extended Data Figure 8 | Antibiotic resistance gene (ARG) relative abundance within and 1177 1178 between terrestrial and oceanic ecosystems. a, Heat map of observed antibiotic resistance gene (ARG) relative abundance at the global scale. Squares and circles correspond to soil and to ocean 1179 samples, respectively. ARG abundance is given on three relative scales for these three datasets. 1180 **b.** ARG relative abundance in ocean samples (across depths) declines with distance from land 1181 (n=139 biologically independent samples), a pattern which was significant at two water depths, 1182 including surface (red) and deep chlorophyll maximum (DCM; green), but not at mesopelagic 1183 (blue). Spearman correlation statistics for specified comparisons are given in the legends. Dotted 1184 1185 lines display Spearman correlations across the whole dataset and within the three depth categories, respectively. n: number of biologically independent samples. 1186

1187

Extended Data Figure 9 | Antibiotic resistance gene (ARG) relative abundance in both 1188 ocean and topsoil samples can be modelled by the relative abundance of fungi and fungi-1189 like protists. a, b, Correlation circle indicating the relationships among fungal classes and ARG 1190 relative abundance as well as the first two partial least squares regression (PLS) components. 1191 1192 Length and direction of vectors indicate the strength and direction of correlations. Percentages show the variation explained by each PLS component. c, d, Linear (Pearson) correlations 1193 1194 between observed and modelled ARG relative abundance based on the relative abundance of fungal taxa in soil (c) and ocean (d). The two principal axes were chosen based on leave-one-out 1195 cross-validation (LOOCV) and explained 42% (LOOCV: R²=0.401) and 71% (LOOCV: 1196 $r^2=0.684$) of the variation of ARG relative abundance in soil and ocean, respectively. Only taxa 1197 significantly associated with ARG relative abundance are shown. Cross validation and Lasso 1198 regression confirmed this result: soil dataset: r=0.619, RMSE=10⁻⁹; n=189 biologically 1199 independent samples; Ocean dataset, r=0.832, RMSE=10⁻⁹; n=139 biologically independent 1200 samples. 1201

1202

1203 Extended Data Figure 10 | Fungal classes are among the main taxa associated with antibiotic resistance gene (ARG) relative abundance, diversity and richness in different 1204 habitats. a, b, Heat map derived from sPLS analysis showing correlation of total ARG relative 1205 abundance, richness and diversity to that of the main taxonomic classes in soil (a) and ocean (b) 1206 metagenomes (see also the supplementary results for analogous results in previously published 1207 soil (from grasslands, deserts agricultural soils) as well as human skin and gut samples). For 1208 statistical details and significance, see Supplementary Table 8. c, d, Heat maps showing 1209 1210 correlation of total ARG relative abundance to that of the main eukarvotic and prokarvotic taxa in soil (c) and ocean (d) based on sparse partial least square (sPLS) regression analysis. All 1211 matrices were normalized by library size and Hellinger transformation. Fungal and fungal-like 1212 classes are shown in **bold** text. See Supplementary Table 15 for ARG gene letter abbreviations. 1213 1214 1215

1216

1217







