# Does diversity beget diversity in microbiomes?

**Authors**: Naïma Madi[1], Michiel Vos[2], Carmen Lia Murall[1],

Pierre Legendre[1] and B. Jesse Shapiro[1,3,4*]


1. Département de sciences biologiques, Université de Montréal, Canada

2. European Centre for Environment and Human Health, University of Exeter, Penryn, UK

3. Department of Microbiology and Immunology, McGill University, Canada

4. McGill Genome Centre, McGill University, Canada


*correspondence: jesse.shapiro@mcgill.ca

18 **Abstract**

19 Microbes are embedded in complex communities where they engage in a wide array of

20 intra- and inter-specific interactions. The extent to which these interactions drive or

21 impede microbiome diversity is not well understood. Historically, two contrasting

22 hypotheses have been suggested to explain how species interactions could influence

23 diversity. 'Ecological Controls' (EC) predicts a negative relationship, where the evolution

24 or migration of novel types is constrained as niches become filled. In contrast, 'Diversity

25 Begets Diversity' (DBD) predicts a positive relationship, with existing diversity

26 promoting the accumulation of further diversity via niche construction and other

27 interactions. Using high-throughput amplicon sequencing data from the Earth

28 Microbiome Project, we provide evidence that DBD is strongest in low-diversity biomes,

29 but weaker in more diverse biomes, consistent with biotic interactions initially favoring

30 the accumulation of diversity (as predicted by DBD). However, as niches become

31 increasingly filled, diversity hits a plateau (as predicted by EC).

32

33

34 **Impact statement:**

35 Microbiome diversity favors further diversity in a positive feedback that is strongest in

36 lower-diversity biomes (*e.g.* guts) but which plateaus as niches are increasingly filled in

37 higher-diversity biomes (*e.g.* soils).

## Introduction

The majority of the genetic diversity on Earth is encoded by microbes (Hug et al., 2016; Lapierre & Gogarten, 2009; Sunagawa et al., 2015) and the functioning of all Earth's ecosystems is reliant on diverse microbial communities (Falkowski et al., 2008). High-throughput 16S rRNA gene amplicon sequencing studies continue to yield unprecedented insight into the taxonomic richness of microbiomes (e.g. (Louca et al., 2019; Sogin et al., 2006)), and abiotic drivers of community composition (e.g. pH; Lauber et al., 2009; Power et al., 2018) are increasingly characterized. Although it is known that biotic (microbe-microbe) interactions can also be important in determining community composition (Needham & Fuhrman, 2016), comparatively little is known about how such interactions, either positive (e.g. cross-feeding; Seth & Taga, 2014) or negative (e.g. toxin-mediated interference competition; Czárán et al., 2002; Hibbing et al., 2010), shape microbiome diversity as a whole.

The dearth of studies exploring how microbial interactions could influence diversity stands in marked contrast to a long research tradition on biotic controls of plant and animal diversity (Elton, 1946; Gause, 2003). In an early study of 49 animal (vertebrate and invertebrate) community samples, Elton plotted the number of species versus the number of genera and observed a ~1:1 ratio in each individual sample, but a ~4:1 ratio when all samples were pooled (Elton, 1946). He took this observation as evidence for competitive exclusion preventing related species, more likely to overlap in niche space, to co-exist. This concept, more recently referred to as niche filling or Ecological Controls (EC) (Schluter & Pennell, 2017), predicts speciation (or, more generally, diversification) rates to decrease with increasing standing species diversity
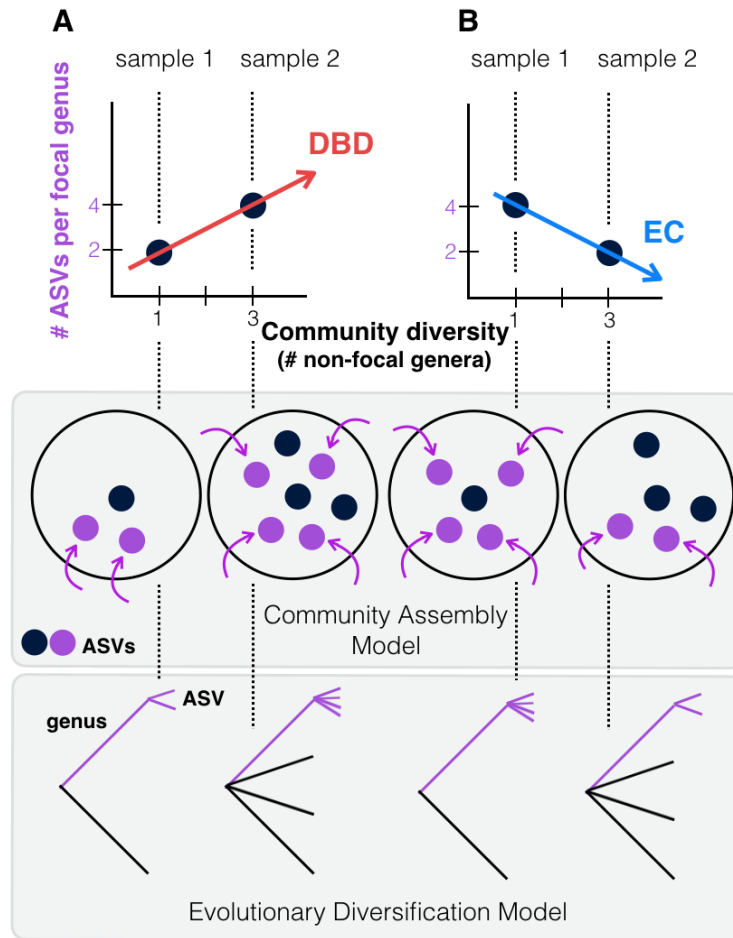
61    because less niche space is available (Rabosky & Hurlbert, 2015). In contrast, the

62    Diversity Begets Diversity (DBD) model predicts that when species interactions create

63    novel niches, standing biodiversity favors further diversification (Calcagno et al., 2017;

64    Whittaker, 1972). For example, niche construction (i.e. the physical, chemical or

65    biological alteration of the environment) could influence the evolution of the species

66    constructing the niche, as well as that of co-occurring species (Laland et al., 1999; San

67    Roman & Wagner, 2018). An alternative to either EC or DBD is The Neutral Theory of

68    Biodiversity and Biogeography, in which all species are functionally equivalent and

69    communities assemble via random sampling (Hubbell, 2001). Neutral Theory serves as a

70    null hypothesis of community assembly in macrobes (Azaele et al., 2016; N. J. Gotelli &

71    McGill, 2006), and more recently in microbiome research (Harris et al., 2017; Li & Ma,

72    2016).

73          Empirical evidence for the action of EC vs. DBD in natural plant and animal

74    communities has been mixed (Calcagno et al., 2017; Emerson & Kolm, 2005; Palmer &

75    Maurer, 1997; Price et al., 2014; Rabosky et al., 2018). Laboratory evolution experiments

76    tracking the diversification of a focal bacterial lineage in communities of varying

77    complexity have also yielded contradictory results, with support for EC, DBD, or

78    intermediate scenarios (Brockhurst et al., 2007; Meyer & Kassen, 2007). For example,

79    diversification of a focal *Pseudomonas* clone was favored by increasing community

80    diversity in the range of 0-20 other strains or species within the same genus (Calcagno et

81    al., 2017; Jousset et al., 2016) but diversification was inhibited in highly diverse

82    communities (*e.g.* hundreds or thousands of species in compost; (Gómez & Buckling,

83 2013)). These experiments are consistent with interspecific competition initially driving

84 (Bailey et al., 2013), but eventually inhibiting diversification as niches are filled.

85        Most laboratory experiments are restricted to relatively short evolutionary time

86 scales and include only a small number of taxa; it is therefore unclear if they can be

87 generalized to natural communities consisting of many more taxa evolving and

88 assembling over much longer periods, spanning more environmental change, greater

89 evolutionary diversification, and frequent migration events. Although the absence of a

90 substantial prokaryotic fossil record hinders deconvoluting speciation and extinction rates

91 (Louca & Pennell, 2020; Marshall, 2017), Louca et al. (Louca et al., 2018) recently

92 estimated that bacterial diversity has mostly increased over the past billion years, with

93 speciation rates slightly exceeding extinction rates. However, because many free-living

94 microbes have high migration rates ("everything is everywhere, but the environment

95 selects" (de Wit & Bouvier, 2006)), we expect that the majority of diversity present

96 within a typical microbiome sample is selected from a pool of migrants rather than

97 having evolved *in situ*. As such, here we broadly define "diversity begets diversity"

98 (DBD) to include the combined effects of community assembly from a migrant pool

99 ('ecological species sorting') and *in situ* evolutionary diversification (**Fig. 1**).

100

101

**Fig. 1. Contrasting the Diversity Begets Diversity (DBD) and Ecological Controls (EC) models. (A)** In this hypothetical scenario, microbiome sample 1 contains one non-focal genus, and two amplicon sequence variants (ASVs) within the focal genus (point at x=1, y=2 in the plot). Sample 2 contains three non-focal genera, and four ASVs within the focal genus (point at x=3, y=4). Tracing a line through these points yields a positive diversity slope, supporting the DBD model (red). **(B)** Alternatively, a negative slope would support the Ecological Controls (EC) model (blue line). In the middle panel, we consider a community assembly model to explain the hypothetical data of the top panel, in which standing diversity (black points) in a community selects (for or against) new types (referred to here as ASVs) which arrive via migration (purple points & arrows). In the bottom panel, we consider an evolutionary diversification model of a focal lineage (genus) into ASVs as a function of initial genus-level community diversity present at the time of diversification.

6

115   To test whether patterns of diversity in natural communities conform to EC or

116 DBD dynamics, we used 2,000 microbiome samples from the Earth Microbiome Project

117 (EMP), the largest available repository of biodiversity based on standardized sampling

118 and sequencing protocols, with 16S rRNA gene amplicon sequence variants (ASVs) as

119 the finest-grained taxonomic unit (Thompson et al., 2017). Following Elton (Elton,

120 1946), we use the equivalent of Species:Genus ratios, calculating a range of taxonomic

121 diversity ratios (up to the Class:Phylum level) as proxies for diversity within a focal

122 taxon, from shallow to deep evolutionary time. We then plot each ratio as a function of

123 the number of non-focal taxa (Genera, Families, Orders, Classes, and Phyla, respectively)

124 with which the focal taxon could interact. We refer to the slope of these plots as the

125 "diversity slope", with negative slopes supporting EC and positive slopes supporting

126 DBD (**Fig. 1**). As a null, we compare these slopes to the expectation under Neutral

127 Theory. To avoid a trivially positive diversity slope due to variation in sequencing effort,

128 all samples were rarefied to 5,000 observations (counts of 16S rRNA gene sequences), as

129 diversity estimates are highly sensitive to sampling effort (Nicholas J. Gotelli & Colwell,

130 2001). As 16S evolves at a rate of roughly 1-2 substitutions per million years (Kuo &

131 Ochman, 2009b), evolutionary diversification within individual EMP samples cannot be

132 uncovered using this marker; rather our data represent mainly a record of community

133 assembly.

134

135

**Results**

**Quantifying the DBD-EC continuum in prokaryote communities compared to neutral null models.** We used generalized linear mixed models (GLMMs) to estimate the diversity slope at each taxonomic level in the EMP data, which revealed a tendency toward positive slopes with significant variation explained by the random effects of lineage, environment, and their interaction (**Table 1, Figure 2, Figure 2 supplements 1-6, Supplementary Data file 1 Section 1**). All models reported here provide significantly better fits compared to models without the fixed effect of community diversity, and coefficients of determination ($R^2$) are higher with the inclusion of random effects, showing their importance (**Supplementary Data file 2**). Examples of how the diversity slope varies across lineages and environments are shown in **Figure 2** and **Figure 2 supplements 2-6**. To assess the significance of these slope estimates in light of potential sampling bias and data structure (Gotelli & Colwell, 2001; Jarvinen, 1982), we considered null models, all of which randomize the associations between ASVs within a sample, thus randomizing any true biotic interactions. Models 1 and 2 are based on draws from the zero-sum multinomial (ZSM) distribution, which arises from the standard Neutral Theory of Biodiversity (**Methods**). Model 1, in which each microbiome sample is drawn from the same ZSM distribution, produces a significantly negative diversity slope (**Figure 2 supplement 7; Table 2**). Model 2, in which each environment draws from a separate distribution, is effectively a composite of Model 1 in which different environments, each with a negative slope, are 'stacked' to yield an overall positive slope (**Figure 2 supplement 7**). However, the Model 2 slope is not significant in a GLMM
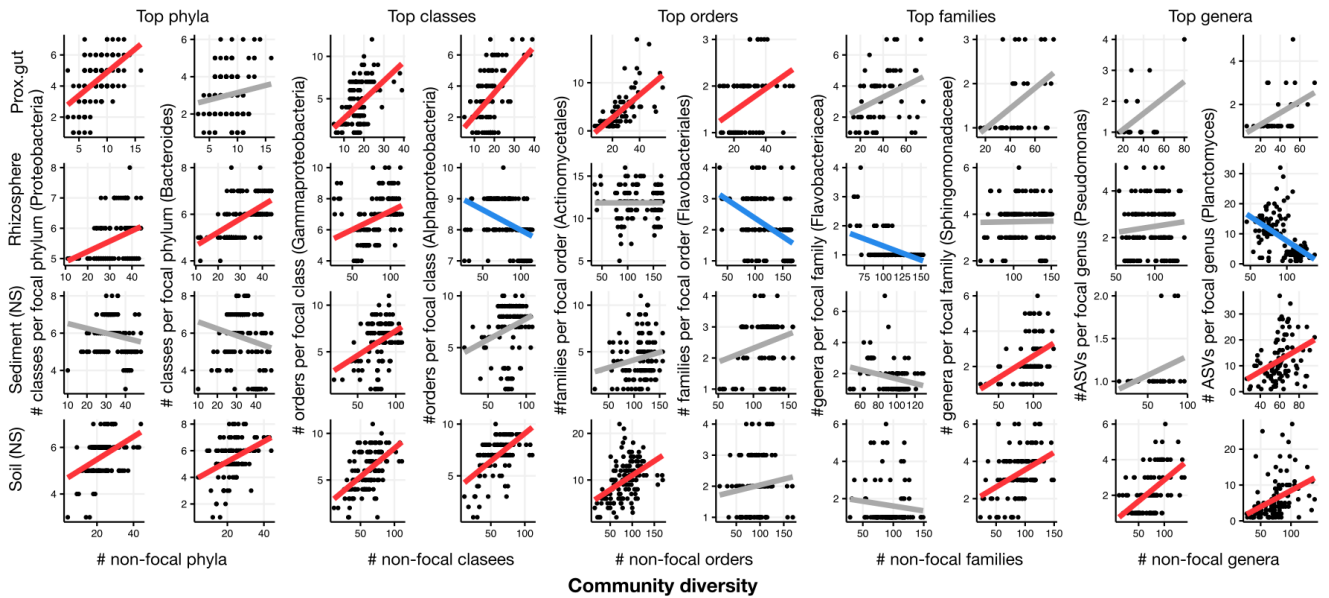
159    accounting for variation across environments (**Table 2, Supplementary Data file 3**

160    **Section 1.2**). In the real EMP data, most individual environments tend toward a positive

161    slope (**Figure 2 supplement 8**). The tendency toward positive diversity slopes in the

162    EMP is therefore not straightforwardly explained by neutral processes.

163         To estimate the power to detect either DBD or EC, we specifically added each of

164    these effects to data simulated under a null model. As expected, adding DBD reversed the

165    negative slope and rendered it positive (**Table 2; Figure 2 supplement 7,**

166    **Supplementary Data file 3 Section 2.1**), suggesting reasonable power to detect DBD

167    when truly present. In contrast, the addition of EC had little effect on the slope,

168    suggesting low power to detect EC under some null models. Taken together, these

169    modelling results suggest that positive diversity slopes observed in the EMP are more

170    readily explained by DBD than by Neutral Theory, whereas negative slopes could be

171    explained by EC, Neutral Theory, or some combination of the two.

172         Because taxonomic labels can be unavailable or inconsistent with phylogenetic

173    relationships (Parks et al., 2018; Vos, 2011) we repeated the analyses using nucleotide

174    sequence identity in the 16S rRNA gene instead of taxonomy, and again recovered

175    generally positive diversity slopes (**Methods**). As a final sensitivity analysis, we repeated

176    the GLMMs using unrarefied community Shannon diversity instead of richness

177    (**Methods**) and obtained similar results, with generally positive diversity slopes that

178    could in some cases be reversed depending on the lineage or environment (**Table 3,**

179    **Supplementary Data file 1 Section 2**). The Shannon diversity metric is robust to

180    sampling effort, suggesting that the results are not biased by undersampling in diverse

181    biomes. Even if undersampling could bias the diversity slope downward in more diverse

182 samples, the effect is unlikely to be large at a rarefaction to 5,000 sequences, and only to

183 occur at the extremes of diversity (*e.g.* very many genera and high ASV:genus ratios) and

184 not at higher taxonomic levels (*e.g.* Class:Phylum) (**Figure 2 supplement 9**).

185



186

**Fig. 2. Focal lineage diversity as a function of community diversity in the top two most prevalent taxa at each taxonomic level.** As in **Fig. 1**, the x-axes show community diversity in units of the number of non-focal taxa (*e.g.* the number of non-Proteobacteria phyla for the left-most column), and the y-axes show the taxonomic ratio within the focal taxon (*e.g.* the number of classes within Proteobacteria). Significant positive diversity slopes are shown in red, negative in blue (linear models, $P < 0.05$, Bonferroni corrected for 17 tests), and non-significant in grey. Note that linear models are distinct from GLMMs, and are for illustrative purposes only. Four representative environments are shown (see **Figure 2 supplements 2-6** for plots in all 17 environments).
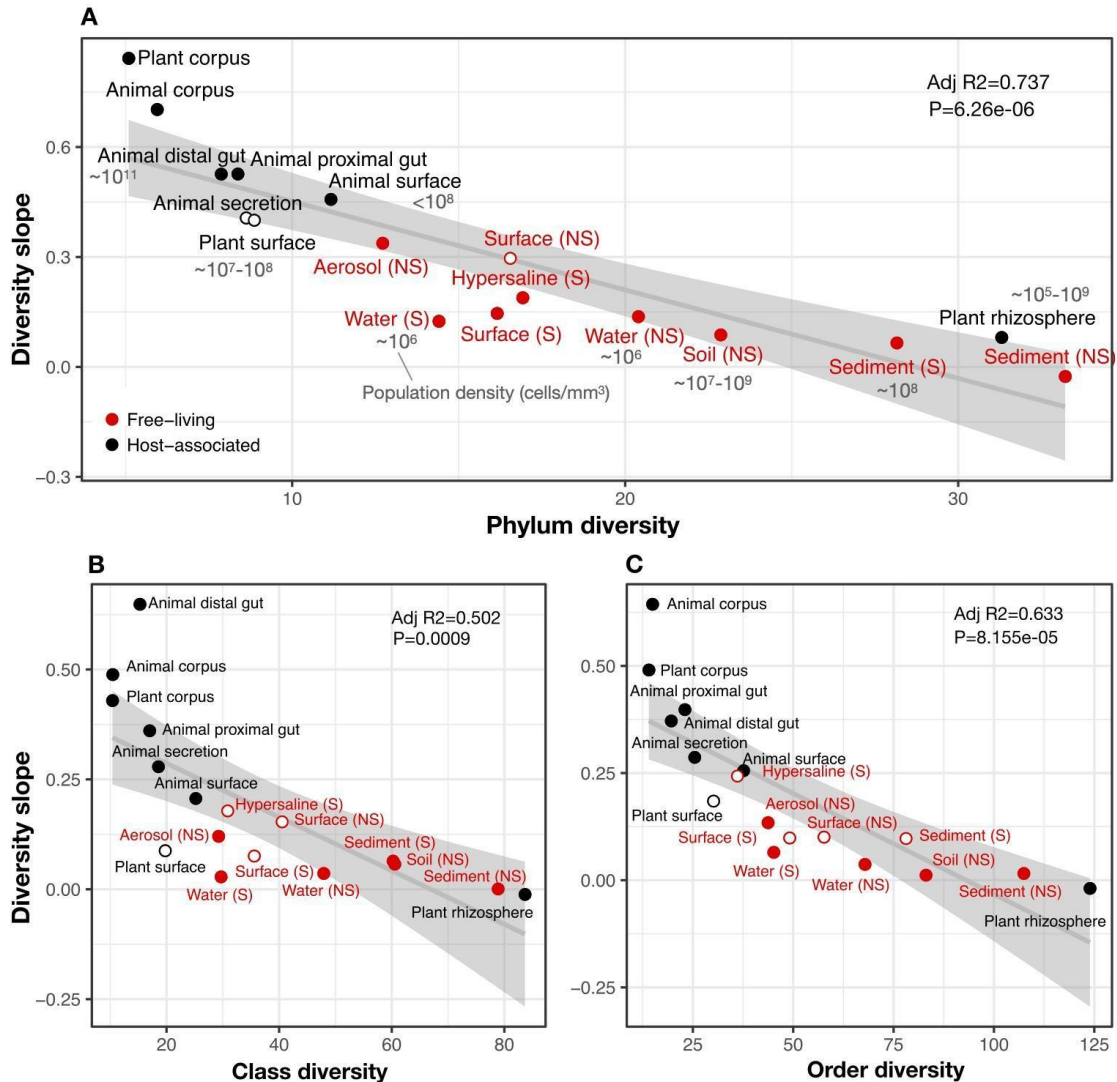
**DBD reaches a plateau at high diversity.** It is expected from theory and experimental studies that a positive DBD relationship should eventually reach a plateau, giving way to EC as niches become saturated (Brockhurst et al., 2007; Gómez & Buckling, 2013). This

200  expectation is borne out in our dataset, particularly in the nucleotide sequence-based

201  analyses which support quadratic or cubic relationships over linear diversity slopes

202  (**Figure 2 supplement 10**). For example, in the animal distal gut, a relatively low-

203  diversity biome, we observed a strong linear DBD relationship at most phylogenetic

204  depths; in contrast, the much more diverse soil biome clearly reaches a plateau (**Figure 2**

205  **supplement 11**).

206         To comprehensively test the hypothesis that more diverse microbiomes

207  experience weaker DBD due to saturated niche space, we used a GLMM including the

208  interaction between diversity and environment as a fixed effect. We considered this

209  model only for taxonomic ratios with significant diversity slope variation by environment

210  (**Table 1**): Family:Order, Order:Class, and Class:Phylum. Diversity slopes were

211  significantly higher in less diverse (often host-associated) biomes, suggesting that niche

212  filling leads to a plateau of DBD in more diverse biomes (**Fig. 3, Supplementary Data**

213  **file 1 Section 3**). The interaction observed in the real EMP data between community

214  diversity and biome type in shaping focal lineage diversity was not observed under a

215  neutral null (Model 2, in which each environment has its own characteristic level of

216  diversity) (**Supplementary Data file 3 Section 1.2**). The DBD plateau observed in more

217  diverse biomes is thus not readily explained by a neutral model, nor is rarefaction

218  expected to bias the diversity slope estimates, particularly at the Class:Phylum level

219  (**Figure 2 supplement 9**). This suggests that the plateau of DBD at higher levels of

220  community diversity is not an artefact of data structure or sampling effort. Finally, we

221  considered whether variation along the EC-DBD continuum could be explained by

222  differential cell density across environments, which could affect both the frequency of

223    cell-cell interactions (a biological effect) or the sampling depth (a technical artefact).

224    Although precise estimates of cell densities in all EMP biomes are not available, we

225    extracted plausible ranges for eight biomes from the literature (Kennedy & de Luna,

226    2005; Lindow & Brandl, 2003; Sender et al., 2016; Whitman et al., 1998) and annotated

227    these in **Figure 3**. It is clear from this figure that relatively high- and low-density samples

228    are found along the range of community taxonomic diversities, demonstrating that cell

229    density is unlikely to drive the trend of decreasing diversity slopes with increasing

230    community diversity.

231

**Fig. 3. The diversity slope of focal taxa is higher in low-diversity (often host-associated) microbiomes.** The x-axis shows the mean number of non-focal taxa: (A) phyla, B) classes, and C) orders in each biome. On the y-axis, the diversity slope was estimated by a GLMM predicting focal lineage diversity as a function of the interaction between community diversity and environment type at the level of A) Class:Phylum, B) Order:Class, and C) Family:Order ratios (**Supplementary Data file 1 Section 3**). The line represents a linear regression; the shaded area depicts 95% confidence limits of the fitted values. Adjusted $R^2$ and $P$-values from the linear fits are shown at the top right of each panel. See **Supplementary Data file 2** for model goodness of fit. Slopes not significantly different from zero are shown as empty circles. Estimates of bacterial cell

243  density from the literature are indicated in grey text, in units of bacteria/mm$^3$. For animal

244  (skin) and plant surface, units of bacteria/mm$^2$ were converted to mm$^3$ assuming layers of

245  bacteria 1 micron thick. For rhizosphere samples we assume a density of 1-2g/cm$^3$

246  (Kennedy & de Luna, 2005).

247

248  **Abiotic drivers of diversity**. Our results thus far suggest that community diversity is a

249  major determinant of the EC-DBD continuum, and by extension that biotic interactions

250  may override abiotic factors in determining where a community lies on the continuum.

251  To formally test for the additional role abiotic drivers might play in generating the

252  observed EC-DBD continuum, we analyzed two data sets in more detail.

253      First, we analyzed a subset of 192 EMP samples with measurements of four key

254  abiotic factors shown to affect microbial diversity (pH, temperature, latitude, and

255  elevation; (Delgado-Baquerizo et al., 2018; Lauber et al., 2009; Power et al., 2018;

256  Schluter & Pennell, 2017)). We fitted a GLMM with focal lineage-specific diversity as

257  the dependent variable, and with the number of non-focal lineages, the four abiotic

258  factors and their interactions as predictors (fixed effects). As in the full EMP dataset

259  (**Table 1**), focal lineage diversity was positively associated with community diversity at

260  all taxonomic ratios in the EMP subset (**Table 4**). As expected, certain abiotic factors,

261  alone or in combination with diversity, had significant effects on focal lineage diversity

262  (**Table 4**). However, the effects of abiotic factors were always weaker than the effect of

263  community diversity (**Table 4**; **Supplementary Data file 1 Section 4**).

264      Second, we used a global 16S sequencing dataset of 237 soil samples associated

265  with more detailed environmental metadata (Delgado-Baquerizo et al., 2018) which we

266  reprocessed to yield ASVs comparable to those in the EMP (**Methods**). This dataset

267  revealed weaker evidence for DBD and stronger effects of abiotic variables on diversity.

268  Community diversity generally had significant positive effects on focal-lineage diversity,

269  but the effect was weak and not detectable at all taxonomic ratios (**Table 5**). Known

270  abiotic drivers of soil bacterial diversity such as pH (Lauber et al., 2009) and latitude

271  (Delgado-Baquerizo et al., 2018) had effects of similar or stronger magnitude compared

272  to the effect of community diversity (**Table 5, Supplementary Data file 4**). The

273  relatively weak effect of DBD and strong effect of abiotic drivers on diversity in this soil

274  dataset can be explained by the fact that soils generally are highly diverse and have

275  relatively low diversity slopes (**Figure 3**).

276       We note that it remains possible that unmeasured abiotic effects could explain

277  some of the DBD effects observed in the EMP. Although only a small subset of abiotic

278  factors was considered, the generally positive diversity slopes in the EMP are not likely

279  to be driven by these factors in the abiotic environment (**Table 4**). Specifically, we

280  consider it unlikely that unmeasured abiotic factors would always act similarly, and in the

281  same direction across multiple different environments, to drive DBD. However, as

282  demonstrated in soil (**Table 5**), abiotic factors may become increasingly important in

283  highly diverse biomes with weak DBD.

284

285  **DBD is more pronounced in resident taxa than in migrant- or generalist taxa.**  A

286  recent meta-analysis of 16S sequence data from a variety of biomes suggests there is an

287  important distinction between generalist lineages found in many environments, compared

288  to specialists with a more restricted distribution (Sriswasdi et al., 2017). Generalists were

289  inferred to have higher speciation rates, suggesting that the DBD-EC balance might differ

290    between generalists and specialists (Sriswasdi et al., 2017). To further investigate this

291    difference, we defined 'residents', taxa with a strong preference for a specific biome, in

292    addition to generalists without a strong biome preference in the EMP dataset. We first

293    clustered environmental samples by their genus-level community composition using

294    fuzzy *k*-means clustering (**Fig. 4a**), which identified three major clusters: 'animal-

295    associated', 'saline', and 'non-saline'. The clustering included some outliers (*e.g.* plant

296    corpus grouping with animals), but was generally consistent with known distinctions

297    between host-associated vs. free-living (Thompson et al., 2017), and saline vs. non-saline

298    communities (Auguet et al., 2010; Lozupone & Knight, 2007). Resident genera were

299    defined as those with a strong preference for a particular environment cluster (whether

300    due to dispersal limitation or narrow niche breadth) using indicator species analysis

301    (permutation test, $P<0.05$; **Fig. 4a**; **Figure 4 supplement 1**; **Supplementary Data file 5**),

302    and genera without a strong preference were considered generalists. When residents of

303    one environmental cluster were (relatively infrequently) observed in a different cluster,

304    we defined them as "migrants" in that sample. For each environment cluster, we ran a

305    GLMM with resident genus-level diversity (the number of non-focal genera) as a

306    predictor of focal-lineage diversity (the ASV:Genus ratio) for residents, generalists, or

307    migrants to that sample (**Supplementary Data file 1 Section 5**).

308         Resident community diversity had no significant effect on the diversity of

309    generalists in animal-associated, saline and non-saline clusters (GLMM, Wald test,

310    $P>0.05$), but was positively correlated with lineage-specific resident diversity (GLMM,

311    Wald test, $z=7.1$, $P= 1.25e-12$; $z=3.316$, $P=0.0009$; $z=7.109$, $P=1.17e-12$, respectively).

312    Resident community diversity significantly decreased migrant diversity in saline

313  (GLMM, z=-3.194, *P*=0.0014) and non-saline environment clusters (GLMM, z=-2.840,

314  *P*=0.0045), but had no significant effect in the animal-associated cluster (GLMM,

315  *P*>0.05) (**Fig. 4b**). These results suggest that, although generalist lineages may have

316  higher speciation rates and colonize more habitats than specialists (Sriswasdi et al.,

317  2017), they have lower diversity slopes. Migrants to the "wrong" environment experience

318  even less DBD, and are even subject to EC in two out of three environment types (**Fig.**

319  **4b**). The accumulation of diversity via successful establishment of migrants may thus be

320  limited, presumably because most niches are already occupied by residents.



321

**Fig. 4. The DBD relationship varies between resident and non-resident genera. (A) Ordination showing genera clustering into their preferred environment clusters**. The matrix of 1128 genera (rows) by 17 environments (columns), with the matrix entries indicating the percentage of samples from a given environment in which each genus is present, was subjected to principal components analysis (PCA). Circles indicate genera and triangles indicate environments (EMPO 3 biomes). Colored circles are genera inferred by indicator species analysis to be residents of a certain environmental cluster, and grey circles are generalist genera. The three environment clusters identified by fuzzy *k*-means clustering are: Non-saline (NS, blue), saline (S, green) and animal-associated (purple). Triangles of the same color indicate EMPO 3 biomes clustered into the same environmental cluster. **(B) DBD in resident versus non-resident genera across environment clusters.** Results of GLMMs modeling focal lineage diversity as a function of the interaction between community diversity and resident/migrant/generalist status. The x-axis shows the standardized number of non-focal resident genera (community diversity); the y-axis shows the number of ASVs per focal genus. Resident focal genera are shown in orange, migrant focal genera in red, and generalist focal genera in black. Red stars indicate a significantly positive or negative slope (Wald test, *P*<0.005). See **Supplementary Data file 2** for model goodness of fit.

## Discussion

Using ~10 million individual marker sequences from the EMP, we demonstrate an overall trend for diversity in focal lineages to be positively associated with overall community diversity, albeit with significant variation across lineages and environments. The strength of the DBD relationship dissipates with increasing microbiome diversity, which we hypothesize is caused by niche saturation. In more diverse biomes such as soil, abiotic factors therefore may become relatively more important in driving focal-lineage diversity. The effect of DBD is strongest among habitat specialists (residents), suggesting that long-term niche adaptation tends to select against the establishment of migrant diversity.

350    While most of the DBD literature considers a model of evolutionary

351    diversification (Schluter & Pennell, 2017; Whittaker, 1972), our results pertain mainly to

352    ecological community assembly dynamics. At the limited resolution of 16S rRNA gene

353    sequences, we do not expect measurable diversification within an individual microbiome

354    sample (Kuo & Ochman, 2009b); however, community diversity could still select for (as

355    in DBD) or against (as in EC) increasing diversity in a focal lineage, even if this lineage

356    diversified before the sampled community assembled. Future work with higher resolution

357    genomic or metagenomic data will enable testing if and how DBD arises in microbial

358    communities via evolutionary diversification, and also how prokaryote diversification is

359    affected by other community members including phages (Brockhurst et al., 2005),

360    protists (Meyer & Kassen, 2007), and fungi (Kastman et al., 2016). Predator-prey, cross-

361    feeding, and other biotic interactions with these non-prokaryotic community members

362    could explain some of the unaccounted variation we observed in diversity slopes across

363    environments.

364    Our dataset also provides an opportunity to explore how DBD relates with

365    genome size evolution. Bacteria with larger repertoires of accessory genes, and thus

366    larger genomes, are able to occupy a wider range of niches (Barberán et al., 2014). Taxa

367    with larger genomes might therefore be hypothesized to better survive and thrive when

368    they disperse into a new location, exhibiting stronger DBD. Although a comprehensive

369    test of this hypothesis will require higher resolution genomic or metagenomic data, as a

370    preliminary exploration we assigned genome sizes to 576 focal genera for which at least

371    one whole genome sequence was available (using the largest recorded genome size for

372    each genus) and added an interaction term between genome size and diversity as a fixed

373    effect in the GLMM (**Methods**). Consistent with our expectation, we observed a

374    significant positive effect of genome size on the diversity slope (GLMM, Wald test,

375    $z=2.5$, $P=0.01$; **Fig. 5, Supplementary Data file 1 Section 6**). This effect was not

376    observed in null models, in which the interaction between community diversity and focal

377    genus genome size was never significant (**Supplementary Data file 3 Section 1.3 and**

378    **2.2**) and so this effect of genome size cannot be trivially explained by data structure. The

379    positive relationship between genome size and DBD is likely even stronger than

380    estimated, because assigning genome sizes to entire genera is imprecise (*i.e.* there is

381    variation in genome size within a genus, or even within species), therefore weakening the

382    correlation.

383        The positive correlation between genome size and DBD observed here could be

384    driven by larger metabolic repertoires encoded by larger genomes (40), potentially

385    creating more opportunities to benefit from cross-feeding, niche construction (San Roman

386    & Wagner, 2018), and other interspecies interactions. This tendency appears to be at odds

387    with the Black Queen hypothesis, which predicts that social conflict between interacting

388    species leads to the inactivation and loss of genes involved in shareable metabolites

389    (public goods), eventually resulting in reduced genome size (Morris & Lenski, 2012).

390    Such a process would produce a negative correlation between the degree of species

391    interactions (*i.e.* community diversity) and genome size (Morris & Lenski, 2012). The

392    interaction between genome size, biotic interactions and diversification thus deserves

393    further study.

**Fig. 5. Positive effect of genome size on DBD.** Results are shown from a GLMM predicting focal lineage diversity as a function of the interaction between community diversity and genome size at the ASV:Genus ratio (**Supplementary Data file 1 Section 6**). The x-axis shows the standardized number of non-focal genera (community diversity); the y-axis shows the number of ASVs per focal genus. Variable diversity slopes corresponding to different genome sizes are shown in a blue color gradient; the shaded area depicts 95% confidence limits of the fitted values. See **Supplementary Data file 2** for model goodness of fit.

Alongside theory and experimental data, the EMP survey data provide a window into the biotic drivers of microbial diversity in nature. In particular, our correlational results support previous experimental and theoretical results showing that DBD is strong

407      when community diversity is low (Calcagno et al., 2017; Jousset et al., 2016), driving the

408      accumulation of diversity in a positive feedback loop until niches are filled and EC starts

409      to predominate (Bailey et al., 2013; Brockhurst et al., 2007; Gómez & Buckling, 2013;

410      Meyer & Kassen, 2007). However, due to the correlational nature of the EMP data, it is

411      not possible to test whether DBD is primarily due to the creation of novel niches via

412      biotic interactions and niche construction (Laland et al., 1999), or due to increased

413      competition leading to specialization on underexploited resources (Hibbing et al., 2010;

414      Jousset et al., 2016). We hope future higher resolution genomic studies, and

415      complementary experiments, will be able to elucidate the types of biotic interactions that

416      promote microbiome diversity. Regardless of the underlying mechanisms, our results

417      demonstrate a general scaling between different levels of community diversity, which has

418      important implications for modeling and predicting community function and stability in

419      response to perturbations (Coyte et al., 2015; Pennekamp et al., 2018). The answer to the

420      question 'why are microbiomes so diverse?' might in a large part be because

421      microbiomes are so diverse (Emerson & Kolm, 2005).

422

430 **Competing interests:** none to declare.

431

432 **Data and materials availability:** All data is available from the Earth Microbiome

433 Project (ftp.microbio.me), as detailed in the Methods. All computer code used for

434 analysis are available at https://github.com/Naima16/dbd.git.

435

436     **Tables**

437     **Table 1. Effects of community diversity on focal lineage diversity across taxonomic**
438     **ratios.** The GLMMs showed statistically a significant positive effect of community
439     diversity on focal lineage diversity. Each row reports the effect of community diversity
440     on focal lineage diversity (Div), as well as its standard error, Wald z-statistic for its effect
441     size and the corresponding *P*-value (left section), or standard deviation on the slope for
442     the significant random effects (right section). SE=standard error, Env=environment type,
443     Lin=lineage type, Lab=Principal Investigator ID, Sample=EMP Sample ID. Interactions
444     are denoted as '*'. n.s.=not significant (likelihood-ratio test). All models provide a
445     significantly better fit than null models without fixed effects (ΔAIC > 10 and *P* < 0.05;
446     **Supplementary Data file 2**).

447
448

|  | Slope (fixed effects) | | | | Standard deviation on the slope (random effects) | | | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Div | SE | z | *P* | Env | Lin | Lin*Env | Env*Lab | Sample |
| ASV:Genus | 0.091 | 0.016 | 5.792 | 6.95e-09 | n.s. | 0.074 | 0.142 | 0.114 | 0.067 |
| Genus:Family | 0.047 | 0.008 | 5.911 | 3.41e-09 | n.s. | 0.071 | 0.07 | 0.039 | n.s. |
| Family:Order | 0.119 | 0.017 | 7.001 | 2.54e-12 | 0.023 | 0.094 | 0.092 | 0.106 | n.s. |
| Order:Class | 0.109 | 0.020 | 5.447 | 5.13e-08 | 0.05 | 0.141 | 0.078 | 0.051 | n.s. |
| Class:Phylum | 0.272 | 0.043 | 6.341 | 2.29e-10 | 0.119 | 0.174 | 0.119 | 0.114 | n.s. |

449

450 **Table 2. GLMMs applied to data simulated under null models.** Null models 1 and 2
451 were generated under the ZSM distribution, with a single distribution for the whole
452 dataset (Model 1) or one distribution per environment (Model 2). Model 3 is similar to
453 Model 1, except with a single Poisson distribution for the whole dataset, and +DBD or
454 +EC refer to adding these effects to 100% of ASVs (see **Methods** and **Figure 2**
455 **supplement 7**). Each row reports the effect of community diversity on focal lineage
456 diversity (Div), as well as its standard error, Wald z-statistic for its effect size and the
457 corresponding *P*-value (Wald test) (left section), or standard deviation on the slope for
458 the significant random effects (right section). SE=standard error, Env=environment type,
459 Lin=lineage type, Sample=EMP Sample ID. n.s.=not significant (likelihood-ratio test),
460 n.t.= not tested, because separate environments were not included in Models 1 or 3.
461

| | Slope (fixed effects) | | | | Stand dev on the slope (random effects) | | | |
|---|---|---|---|---|---|---|---|---|
| | Div | SE | z | *P* | Env | Lin | Lin*Env | Sample |
| **Model 1** | -0.005 | 0.000 | -9.807 | <2e -16 | n.t. | 0.639 | n.t. | n.s. |
| **Model 2** | n.s. | | | | | | | |
| **Model 3** | -0.012 | 0.002 | -6.552 | 5.69e-11 | n.t. | 0.021 | n.t. | n.s. |
| **Model 3 + DBD** | 0.016 | 0.001 | 11.48 | <2e-16 | n.t. | 0.008 | n.t. | n.s. |
| **Model 3 + E** | -0.011 | 0.002 | -6.14 | 8.26e-10 | n.t. | ns | n.t. | n.s. |

| C | | | | | | | |
|---|---|---|---|---|---|---|---|

462

463 **Table 3. GLMMs with community diversity measured using Shannon diversity.**
464 Results are shown from GLMMs with Shannon diversity of non-focal taxa (Div) as a
465 predictor of ASVs richness of focal taxa. Each row reports the estimate (Div), as well as
466 its standard error, Wald z-statistic for its effect size and the corresponding *P*-value (Wald
467 test) (left section), or standard deviation on the slope for the significant random effects
468 (right section). SE=standard error, Env=environment type, Lin=lineage type,
469 Lab=Principal Investigator ID, Sample=EMP Sample ID. n.s.=not significant (likelihood-
470 ratio test).
471

| | Fixed effects | | | | Random effects | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Div | SE | z | *P* | Env | Lin | Env*Lin | Env*Lab | Sample |
| **Genus** | 0.055 | 0.013 | 4.33 | 1.49e-05 | n.s. | 0.08 | 0.15 | 0.085 | 0.054 |
| **Family** | 0.148 | 0227 | 6.491 | 8.51e-11 | n.s. | 0.184 | 0.268 | 0.16 | 0.134 |
| **Order** | 0.378 | 0.038 | 9.864 | <2e-16 | n.s. | 0.34 | 0.417 | 0.258 | 0.202 |
| **Class** | 0.398 | 0.05 | 7.973 | 1.54e-15 | n.s. | 0.369 | 0.46 | 0.326 | 0.262 |
| **Phylum** | 0.319 | 0.088 | 3.614 | 0.0003 | 0.169 | 0.316 | 0.5 | 0.495 | 0.378 |

472

473 **Table 4. Community diversity has a stronger effect than abiotic factors on focal lineage**
474 **diversity (EMP dataset).** Results are shown from GLMMs with community diversity, four
475 abiotic factors (temperature, elevation, pH, and latitude), and their interactions with community
476 diversity, as predictors of focal lineage diversity. Random effects on the intercept included
477 environment, lineage, lab ID and sample ID. Each row reports the taxonomic ratio, the predictors
478 used in the GLMM (fixed effects only), their estimate (Est), standard error (SE) and *P*-value (P)
479 (Wald test). Interactions are denoted as '*'. Random effects are not shown.
480

|  | Predictor | Est | SE | P |
|---|---|---|---|---|
| ASV:Genus | Div | 0.128 | 0.013 | < 2e-16 |
|  | Temperature | 0.04 | 0.014 | 0.00479 |
|  | Div*Temperature | 0.043 | 0.014 | 0.00175 |
|  | Div*Latitude | 0.031 | 0.013 | 0.02119 |
|  | Div*Elevation | -0.031 | 0.014 | 0.02829 |
| Genus:Family | Div | 0.094 | 0.009 | < 2e-16 |
|  | Temperature | 0.026 | 0.009 | 0.00268 |
|  | pH | -0.042 | 0.009 | 5.88e-06 |
| Family:Order | Div | 0.131 | 0.01 | < 2e-16 |
| Order:Class | Div | 0.184 | 0.01 | < 2e-16 |
|  | Div*Temperature | 0.032 | 0.009 | 0.000827 |
|  | Div*Latitude | 0.023 | 0.008 | 0.005403 |
| Class:Phylum | Div | 0.236 | 0.011 | < 2e-16 |
|  | Div*Temperature | 0.059 | 0.014 | 2.15e-05 |
|  | Div*Latitude | 0.03 | 0.011 | 0.00884 |

481

**Table 5. GLMMs applied to a soil dataset.** Each row reports the taxonomic ratio, the predictors used in the GLMM (fixed effects only), their estimate (Est), standard error (SE) and *P*-value (P) (Wald test). Left columns: GLMM with community diversity (Div) and all abiotic variables considered separately, as predictors of focal lineage diversity. Right columns: GLMM with community diversity (Div) and the three first principle components (PCs) representing abiotic variables, as predictors of focal lineage diversity. n.s., non-significant (LRT test). All models provide a significantly better fit than null models without fixed effects (ΔAIC > 10 and *P* < 0.05; **Supplementary Data file 2**), except for the GLMM with abiotic factors at the Family:Order level, where latitude has a significant effect on focal lineage diversity but its effect is nearly null, with a ΔAIC between full and null model of 4 and a null marginal $R^2$.

| | GLMMs with abiotic variables | | | | GLMMs with the 3 first PCs | | | |
|---|---|---|---|---|---|---|---|---|
| | Predictor | Est | SE | P | Predictor | Est | SE | P |
| **ASV:Genus** | Div | n.s. | | | Div | 0.064 | 0.016 | 9.47e-05 |
| | Latitude | 0.294 | 0.025 | < 2e-16 | PC1 | -0.065 | 0.007 | < 2e-16 |
| | UV_light | -0.177 | 0.016 | < 2e-16 | PC2 | -0.03 | 0.006 | 1.98e-05 |
| | MDR | 0.028 | 0.006 | 7.12e-06 | | | | |
| | NPP2003_2015 | -0.066 | 0.005 | < 2e-16 | | | | |
| | Latitude^2 | -0.3 | 0.029 | < 2e-16 | | | | |
| | Clay_silt^2 | -0.012 | 0.004 | 0.003 | | | | |
| | Soil_N^2 | -0.007 | 0.001 | 1.66e-06 | | | | |
| | Soil_C_N_ratio | 0.003 | 0.001 | 0.004 | | | | |
| | PSEA^2 | 0.01 | 0.002 | 4.84e-06 | | | | |
| | MDR^2 | 0.017 | 0.003 | 2.40e-08 | | | | |
| | NPP2003_2015 | -0.016 | 0.004 | 0.0001 | | | | |
| **Genus:Family** | Div | 0.032 | 0.01 | 0.0011 | Div | 0.033 | 0.01 | 0.001 |
| | Latitude | -0.035 | 0.006 | 2.04e-09 | PC1 | -0.016 | 0.006 | 0.02 |
| | | | | | PC2 | 0.02 | 0.006 | 0.00089 |
| **Family:Order** | Div | n.s. | | | Div | n.s. | | |
| | Latitude | -0.0005 | 0.0002 | 0.0105 | PC1 | -0.026 | 0.007 | 0.00032 |
| | | | | | Div*PC1 | 0.04 | 0.006 | 2.14e-12 |
| | | | | | Div*PC3 | 0.023 | 0.005 | 1.68e-06 |
| **Order:Class** | Null model with no predictor was significant | | | | | | | |
| **Class:Phylum** | Div | 0.032 | 0.01 | 0.00174 | Div | 0.032 | 0.01 | 0.003 |
| | pH | 0.074 | 0.01 | 4.37e-13 | PC1 | -0.051 | 0.01 | 3.54e-07 |
| | | | | | PC2 | -0.028 | 0.01 | 0.006 |

498     **Supplementary Figure Legends**

499

500     **Figure 2 supplement 1. Distributions of diversity slope estimates across different**
501     **random effects, from the GLMMs predicting focal lineage diversity as a function of**
502     **community diversity. (A)** Class:Phylum, **(B)** Order:Class, **(C)** Family:Order, **(D)**
503     Genus:Family, and **(E)** ASV:Genus. Estimation of random effect coefficients from the
504     GLMMs (Table S1), shows that the effect of diversity on focal lineage diversity (slope
505     estimates) are generally positive but could be negative in some lineages or combinations
506     of environment, lineage (Environment*Lineage), and the laboratory that submitted the
507     dataset (Environment*Lab).

508

509     **Figure 2 supplement 2. Focal lineage diversity as a function of community diversity**
510     **across biomes in the three most prevalent phyla.** (A) Proteobacteria, (B) Bacteroidetes,
511     (C) Actinobacteria. Linear models are shown for the number of classes per phylum (y-
512     axis) as a function of community diversity (number of non-focal phyla, x-axis) in each of
513     the 17 environments (EMPO3 biomes). Only environments containing the focal lineage
514     are shown. $P$-values are Bonferroni corrected for 17 tests. Significant ($P$ <0.05) models
515     are shown with red trend lines, non-significant ($P$ > 0.05) trends are shown in blue.

516

517     **Figure 2 supplement 3. Focal lineage diversity as a function of community diversity**
518     **across biomes in the three most prevalent classes.** Linear models are shown for the
519     number of orders per class (y-axis) as a function of community diversity (non-focal
520     classes, x-axis) in each of the 17 environments (EMPO3 biomes). Only environments
521     containing the focal lineage are shown. Significant positive diversity slopes are shown in
522     red, negative in blue (linear models, $P$ <0.05, Bonferroni corrected for 17 tests), and non-
523     significant in grey.

524

525     **Figure 2 supplement 4. Focal lineage diversity as a function of community diversity**
526     **across biomes in the three most prevalent orders.** Linear models are shown for the
527     number of families per order (y-axis) as a function of community diversity (non-focal
528     orders, x-axis) in each of the 17 environments (EMPO3 biomes). Only environments
529     containing the focal lineage are shown. Significant positive diversity slopes are shown in
530     red, negative in blue (linear models, $P$ <0.05, Bonferroni corrected for 17 tests), and non-
531     significant in grey.

532

533     **Figure 2 supplement 5. Focal lineage diversity as a function of community diversity**
534     **across biomes in the three most prevalent families.** Linear models are shown for
535     genera per family (y-axis) as a function of community diversity (non-focal families, x-
536     axis) in each of the 17 environments (EMPO3 biomes). Only environments containing the
537     focal lineage are shown. Significant positive diversity slopes are shown in red, negative
538     in blue (linear models, $P$ <0.05, Bonferroni corrected for 17 tests), and non-significant in
539     grey.

540

541     **Figure 2 supplement 6. Focal lineage diversity as a function of community diversity**
542     **across biomes in the three most prevalent genera.** Linear models are shown for ASVs
543     per genus (y-axis) as a function of community diversity (non-focal genera, x-axis) in each

544    of the 17 environments (EMPO3 biomes). Only environments containing the focal
545    lineage are shown. Significant positive diversity slopes are shown in red, negative in blue
546    (linear models, $P$ <0.05, Bonferroni corrected for 17 tests), and non-significant in grey.
547
548    **Figure 2 supplement 7. Null models based on Neutral Theory.** Results are shown from
549    data simulated under (A) neutral Model 1, (B) neutral Model 2, or (C) neutral Model 3.
550    Model 1 is sampled from the zero-sum multinomial distribution with a single distribution
551    for the whole dataset, while Model 2 includes a separate distribution for each of the 17
552    different environments (EMPO 3 biomes). In Model 3 (C), the effect of DBD (top rows)
553    or EC (bottom rows) are "spiked in" at different levels, ranging from 0 to 100% of ASVs
554    in a sample. Blue lines show a linear fit, with slopes ($m$) estimated by GLMM in selected
555    panels. See Methods for model details, and Table 2 and Supplementary Data file 3,
556    Section 1.2 for full GLMM results.
557
558    **Figure 2 supplement 8. Lineage diversity (mean ASV:Genus ratio among all**
559    **lineages) as a function of community diversity (number of genera) in the EMP data**.
560    Samples from different environments (EMPO level 3) are shown in different colors, each
561    with their corresponding linear model fit.
562
563    **Figure 2 supplement 9. Taxonomic ratios estimated from simulated rarefied**
564    **sequence data.** Each panel simulates a set of microbiome samples that differ in their
565    diversity (number of genera in left panels **A** and **B**, number of phyla in right panels **C** and
566    **D**) while maintaining a set true taxonomic ratio (horizontal black line). **(A)** True ratio set
567    to 2 ASVs/genus, close to the per-sample mean and median in the real EMP data, in a
568    range of samples between 1 and 1128 named genera, as observed in the real EMP data.
569    **(B)** True ratio set to 20 ASVs/genus, equal to the overall mean of 22,014 named ASVs in
570    1128 named genera, and close to the maximum ratios observed in individual samples
571    (Fig. 2 supplement 6). Insets show the ranges of 1-50 and 51-150 genera, approximating
572    observations from lower- or higher-diversity samples such as gut and soil, respectively
573    (Fig. 2 supplement 6). The insets only show the rarefaction to 5,000 sequences, as used in
574    the real EMP dataset. **(C)** True ratio set to 3 classes/phylum, close to the per-sample
575    mean and median in the real EMP data, in a range of samples between 1 and 84 named
576    phyla, as observed in the real EMP data. **(D)** True ratio set to 10 classes/phylum, close to
577    the maximum ratios observed in individual samples (Fig. S2). Different rarefaction levels
578    are shown as different colored lines.
579
580    **Figure 2 supplement 10. Linear, quadratic and cubic models for the relationship**
581    **between focal lineage diversity and community diversity for varying levels of %**
582    **nucleotide identity.** Community diversity was estimated as the number of clusters at a
583    focal level ($d_i$) and focal lineage diversity as the mean of the clusters at the rank above
584    ($d_{i+1}/d_i$). All $P$-values are < 0.001. Linear fit (grey); quadratic fit (blue), cubic fit (red);
585    same colors for the associated adjusted $R^2$. The x-axis (diversity) shows the number of
586    clusters at the focal percent-identity level ($d_i$), and the y-axis (diversification) is the mean
587    of the clusters at the rank above ($d_{i+1}/d_i$).
588

**Figure 2 supplement 11. Linear, quadratic and cubic models for each environment type for varying levels of % nucleotide identity.** Community diversity was estimated as the number of clusters at a focal level ($d_i$) and focal lineage diversity as the mean of the clusters at the rank above ($d_{i+1}/d_i$). Linear (grey), quadratic (blue) and cubic (red), with corresponding adjusted R-squared values in the same colour. $P$-values are Bonferroni corrected for 17 tests. Significant, $P < 0.05$ (solid lines), non-significant (dashed lines). The x-axis shows the number of clusters at the focal percent-identity level ($d_i$), and the y-axis is the mean of the clusters at the rank above ($d_{i+1}/d_i$).

**Figure 4 supplement 1. Resident genera of environment clusters.** Results from indicator species analysis illustrated as a heatmap**.** Only the 25 resident genera with the highest indval indices and $P<0.05$ (permutation test) are shown for every environment cluster (animal-associated, non-saline and saline free). For the full results see **Supplementary Data file 5.**

**Supplementary File legends**

**File 1. Full GLMM outputs for the EMP data.**

**File 2. Goodness of fit for the GLMMs.**

**File 3. Full GLMM output for simulated data under Neutral Theory models**

**File 4. Full GLMM output for soil data (Delgado et al.)**

**File 5. Indicator species analysis.** The table shows the assignment of each genus to one of three environment types.

**File 6. Genome size assignment.** The table shows genome sizes assigned to each genus.

621 **Materials and Methods**

622 **Earth Microbiome Project dataset.** We used the EMP '2000 subset' of 16S rRNA gene

623 sequences, rarefied to 5000 sequences per sample. This subset contains 155,002 ASVs

624 from 2,000 samples with an even distribution across 17 natural environments (EMP

625 Ontology level 3). Data were downloaded from the EMP FTP server ([ftp.microbio.me](ftp.microbio.me)),

626 on February 9, 2018.

627

628 Specifically, 16S rRNA-V4 region reads (90 bp, GreenGenes 13.8 taxonomy) along with

629 environmental data and EMPO3 designations

630 ([http://press.igsb.anl.gov/earthmicrobiome/protocols-and-standards/empo/](http://press.igsb.anl.gov/earthmicrobiome/protocols-and-standards/empo/)) were

631 downloaded from the EMP FTP server ([ftp.microbio.me](ftp.microbio.me)), on February 9, 2018. Sequence

632 summaries were downloaded from :

633 ftp://[ftp.microbio.me/emp/release1/otu_distributions/otu_summary.emp_deblur_90bp.sub](ftp.microbio.me/emp/release1/otu_distributions/otu_summary.emp_deblur_90bp.sub)

634 [set_2k.rare_5000.tsv](set_2k.rare_5000.tsv), environmental data from:

635 ftp://[ftp.microbio.me/emp/release1/mapping_files/emp_qiime_mapping_release1.tsv](ftp.microbio.me/emp/release1/mapping_files/emp_qiime_mapping_release1.tsv), and

636 EMPO3 designations from :

637 ftp://ftp.microbio.me/emp/release1/mapping_files/emp_qiime_mapping_subset_2k.tsv.

638 The list of the associated 97 studies and 61 corresponding principal investigator identities

639 were downloaded from [https://www.nature.com/articles/nature24621#s1](https://www.nature.com/articles/nature24621#s1).

640 Based on the ASV annotations across samples, we estimated the taxonomic ratio for each

641 focal lineage (ASV:Genus, Genus:Family, Family:Order, Order:Class and Class:Phylum),

642 along with the number of non-focal lineages (dbd_analys_input.py,

643     glmm_analys_input.py, Python Version 2.7). Unclassified ASVs were removed from the

644     analyses.

645

646     **Generalized linear mixed model (GLMM) analyses.** We used GLMMs to determine

647     how focal lineage diversity (*e.g.* its ASV:Genus ratio) is affected by community diversity

648     (*e.g.* non-focal genera). The effects of environment (as defined by the EMP Ontology

649     'level 3 biomes') and the focal lineage identity were included as random effects on the

650     slope and intercept. We also controlled for the submitting laboratory (identified by the

651     principal investigator) and the EMP unique sample identifier (i.e. if two taxa were part of

652     the same sample).

653         All models were fitted in Rstudio (Version 1.1.442, R Version 3.5.2) using the

654     glmer function of the lme4 package (Bates et al., 2015). Data standardization

655     (transformation to a mean of zero and a standard deviation of one) was applied to all

656     predictors to get comparable estimates. In models with only one predictor, applying

657     standardization resolved convergence warnings and considerably sped up the

658     optimization. We first tested the significance of random effects, by using likelihood-ratio

659     tests (LRTs, implemented in the anova function in the R stats package) on nested models

660     where each random effect was dropped one at a time. We then assessed the significance

661     of fixed effects using drop1 function from stats package with the likelihood-ratio test

662     option (this function drops individual terms from the full model and compares models

663     based on the AIC). We calculated the Akaike information criterion (AIC) of each

664     significant model and a null model including all random effects but no fixed effects other

665     than the intercept. We then report the difference in AIC between the full and null models

666    (ΔAIC), along with a likelihood ratio test *p*-value to assess the significance of the full

667    model relative to the null. Only significant models ($P<0.05$) are reported.

668    As an additional test of the goodness of fit for the significant models, we

669    estimated the coefficient of determination ($R^2$) using the r.squaredGLMM function from

670    the MuMIn R package. This function implements a method developed by Nakagawa and

671    Schielzeth and its extension for random slopes (Johnson, 2014; Nakagawa & Schielzeth,

672    2013). Two values were estimated: the marginal $R^2$, as a measure of the variance

673    explained only by fixed effects, and the conditional $R^2$ as a measure of the variance

674    explained by the entire model (both fixed effects and random effects). Only results from

675    $R^2$ estimation based on lognormal and trigamma methods were reported because they are

676    specific to the logarithmic link function used in all GLMMs.

677    Diagnostic plots (plot and qqnorm R functions in base and stats packages) were

678    checked for each model to ensure that residual homoscedasticity (homogeneity of

679    variance) was fulfilled: no increase of the variance with fitted values and residuals were

680    symmetrically distributed tending to cluster around the 0 of the ordinate, but with an

681    expected pattern due to count data. Normality plots were imperfect, but they generally

682    showed that the residuals were close to being normally distributed. The assumption of

683    normality is often difficult to fulfill with high numbers of observations, as is the case in

684    our models (https://www.statisticshowto.datasciencecentral.com/shapiro-wilk-test/), and

685    non-normality is less of concern than heteroscedastic for the validity of GLMMs

686    (https://bbolker.github.io/mixedmodels-misc/ecostats_chap.html#diagnostics).

687    We tested for overdispersion using the overdisp_fun R function available at

688    https://bbolker.github.io/mixedmodels-misc/glmmFAQ.html, and found that all the

689      models were not overdispersed, but rather were underdispersed : the ratio of the sum of

690      squared Pearson residuals to residual degrees of freedom was < 1 and non-significant

691      when tested with a chi-squared test. The only exception was Shannon diversity-based

692      GLMMs. In case of underdispersion and given that underdispersion leads to more

693      conservative results, we retained the GLMMs with Poisson error distribution, despite the

694      underdispersion. (GLMM FAQ; Ben Bolker and others; 25 September 2018;

695      [https://bbolker.github.io/mixedmodels-misc/glmmFAQ.html#underdispersion](https://bbolker.github.io/mixedmodels-misc/glmmFAQ.html#underdispersion)). For

696      Shannon diversity-based GLMMs, we accounted for overdispersion by adding an

697      observation-level random effect to the GLMMs (Elston et al., 2001).

698

699      **Taxonomy-based GLMMs**

700      To test how focal lineage diversity (*e.g.* its ASV:Genus ratio) is affected by community

701      diversity (*e.g.* non-focal genera richness), for different environment types and lineages

702      across all taxonomic ratios, we used generalized linear mixed models (GLMMs) fitted on

703      the EMP dataset. As the dependent variable (focal lineage diversity, defined as taxonomic

704      ratios, ASV:Genus, Genus:Family, Family:Order, Order:Class, and Class:Phylum) was a

705      count response, we used a Poisson error distribution with a log link function. Community

706      diversity (number of non-focal lineages: non-focal Genera, Families, Orders, Classes, and

707      Phyla), standardized to a mean of zero and a standard deviation of one, was specified as

708      the predictor (fixed effect). We included the following random effects on the slope and

709      intercept: lineage (Lin), environment (Env), environment nested within lineage (a lineage

710      may be present in different environments) and lab (the principal investigator who

711      conducted the EMP study) nested within environment (different labs sampled and

712    sequenced a given environment) (as suggested in http://bbolker.github.io/mixedmodels-

713    misc/glmmFAQ.html). Defining random effects on the slope enabled us to test slope

714    variation across groups of each categorical variable (*e.g.* slope variation between different

715    environments or different lineages). We included the EMP unique sample ID as a random

716    effect to control for dependencies between observations (if two taxa were part of the

717    same sample) (**Table 1**, **Supplementary file 1 section 1**).

718

719    **Shannon diversity-based GLMMs**

720    We also tested whether ASV diversity in a focal taxon is dependent on the diversity of all

721    other ASVs in that sample (rather than the diversity at only the focal taxonomic level, as

722    in the taxonomy-based GLMMs above). We used the Shannon diversity index, which is

723    robust to differences in sampling effort, and generally reaches a plateau at 5,000

724    sequences or fewer (48, 49). To do so, we fitted a GLMM with the number of ASVs per

725    focal taxon as the response variable, and the Shannon diversity based on ASVs across all

726    non-focal taxa (z-standardized) as the predictor (fixed effect), the random effects were

727    kept as in the taxonomy-based GLMMs, but we added an observation-level random effect

728    to account for overdispersion (**Table 3**, **Supplementary file 1 section 2**). To avoid

729    dependence between the response and predictor variables, we used the rarefied ASV

730    dataset (5,000 ASVs/sample as above) as the response variable, and the Shannon

731    diversity calculated on unrarefied data from the same samples as the predictor.

732

733    **Null models.** We considered three null models, all of which randomize the associations

734    between ASVs within a sample, thus breaking any true biotic interactions. These null

735    models were randomly generated matrices of the same size as the real EMP dataset, but

736    based on a distribution that arises from the Neutral Theory of Biodiversity. Neutral

737    Theory postulates that the biodiversity of a metacommunity is governed by independent

738    random population dynamics across species. The aggregate behaviour is quantified by the

739    fundamental biodiversity number $\theta$, such that $\theta = 2 J_M \upsilon$, where $J_M$ is the size of the

740    metacommunity and $\upsilon$ is the speciation rate. Parametrized by $\theta$, the metacommunity zero-

741    sum multinomial distribution (mZSM) was developed to obtain random samples of size $J$

742    *(Alonso & McKane, 2004)*. We used this mZSM distribution (implemented with the *sads*

743    package in R; http://search.r-project.org/library/sads/html/dmzsm.html) to generate the

744    counts of the ASVs for each dataset in models 1 and 2. Model 1 assumes that the whole

745    dataset follows the same species abundance distribution (SAD), characterized by a

746    mZSM with $\theta = 50$. Model 2 assumes that each environment has its own SAD and thus

747    all the samples of a single environment are assigned the same $\theta$ but are distinct across

748    environments ($\theta$ was chosen uniformly between 1 and 100). The number of samples per

749    environment were the same as the EMP dataset. To obtain similar mean counts as the real

750    dataset, we set $J = 1000$ for both models 1 and 2, in order to vary $\theta$ from 1 to 100. These

751    values are reasonable based on previous studies that estimated these parameters from

752    microbiome data (Li & Ma, 2016). We included a down-sampling step to replicate the

753    zero-inflated nature of the real dataset (on average there were only 96 ASVs per sample

754    while there was a total of 22,014 ASVs in the entire EMP dataset). To replicate the

755    sampling effect due to rarefaction, we first created a vector of all individuals from a

756    single sample. We then selected 5000 individuals at random whose identities determined

757    which ASVs were found in that sample. These neutrally-derived random matrices, null

758    models 1 and 2, were plotted using the same plots (ASV:Genus vs number of genera) as

759    the real EMP dataset and were then analyzed using GLMMs with community diversity as

760    a predictor of focal lineage diversity (fixed effect), with lineage identity and EMP sample

761    ID as random effects. For Model 1, the slope was significantly negative (GLMM, Wald

762    test, z=-9.807, $P<$2e-16). For Model 2, the null GLMM (including the intercept only) was

763    significant, meaning that the community diversity has no significant effect on focal

764    lineages diversity (Likelihood-ratio test between the model with the predictor and the

765    intercept-only model, $P$=0.9399).

766         To generate a null model for a metacommunity assembled by niche processes,

767    null model 3 was made by sampling from a single Poisson distribution ($\lambda = 0.01$) for each

768    element of the data matrix. We used the Poisson distribution as a sensitivity analysis

769    compared to the ZSM, and found the two behave quite similarly (*i.e.* Model 1 and 3

770    produce qualitatively similar results). The probability of size zero was sufficiently large

771    that the down-sampling step was not needed for this model. Next, DBD and EC effects

772    were added to null model 3 according to the following procedure. An element was chosen

773    at random in a sample and tested if it is empty or full (*i.e.* checks the presence/absence of

774    a particular ASV). If the element is full then the DBD algorithm fills an empty element

775    chosen at random in the same sample, while the EC algorithm empties a filled element in

776    the same sample. This is to mimic the effect of DBD creating a niche for a new ASV, or

777    EC removing a niche based on the existing diversity. The strength of DBD or EC effects

778    were determined by the percent of elements tested. These data were analyzed with

779    GLMMs to test the power of our models to detect DBD or EC (**Table 2**, Supplementary

780    Data file 3 Section 2.1).

781 **Rarefaction simulation**

782 We constructed a simple simulation in which each microbiome sample may differ in total

783 diversity (*e.g.* in the observed range of genera) while maintaining a constant taxonomic

784 ratio (*e.g.* ASV:genus ratio = 2). To mimic rarefaction, we then sampled a set number of

785 sequencing reads from each synthetic community, assuming ASVs are sampled with

786 equal probability and plotted the observed taxonomic ratio (**Fig. 2 supplement 9**). This

787 simple simulation is implemented in permute_ASVs_synthetic.pl.

788

789 **Nucleotide sequence-based analysis.** We clustered ASVs at decreasing levels of

790 nucleotide identity, from 100% identical ASVs down to 75% identity (roughly equivalent

791 to phyla (Konstantinidis & Tiedje, 2005)). We estimated focal cluster diversity as the

792 mean number of descendants per cluster (e.g. number of 100% clusters per 97% cluster)

793 and plotted this against the total number of clusters (97% identity in this example). This

794 approach has the advantage of including sequences even if they come from unnamed

795 taxa. For each of the six nucleotide divergence ratios tested, the relationship between total

796 number of clusters and focal cluster diversity was positive (**Fig. 2 supplement 10**),

797 consistent with DBD and suggesting that the taxonomic analyses were qualitatively

798 unbiased.

799      Fasta files with all ASVs per sample were produced by a python script

800 (Construct_fasta_per_sample.py, Python Version 2.7) from the sequences summary file

801 (otu_summary.emp_deblur_90bp.subset_2k.rare_5000 from EMP ftp server). We

802 clustered sequences from each sample using USEARCH V9.2 and estimated sample

803 diversity as the total number of clusters at a given level (*e.g.* 97% identity) and focal

804     cluster diversity as the mean number of descendent clusters (*e.g.* number of 100%

805     clusters per 97% cluster). To describe the putative DBD or EC relationships, we tested

806     three models: linear, quadratic and cubic (lm function in R). Model comparisons were

807     based on the adjusted $R^2$ (**Figure 2 supplement 10**).

808          We note that diversity at level $i$ ($d_i$) and at level $i$+1 ($d_{i+1}/d_i$) are not independent

809     in this analysis because $d_{i+1}$ must be greater than or equal to $d_i$. To assess the effects of

810     this non-independence on the results, we conducted permutation tests by randomizing the

811     associations between $d_i$ and $d_{i+1}$. Using 999 permutations, *P*-values were calculated based

812     on how many times we observed a correlation greater than that seen in the real data

813     (cor.test R function with kendall method). In each permutation, we recalculated the

814     significance test (Wald z) for the correlation in the randomized data, and then computed

815     the *P*-value based on how many times we observed a z value greater than that of the

816     original data. At all six levels of nucleotide identity, the real data always showed a

817     significantly stronger positive correlation when compared to permuted data ($P = 0.001$),

818     indicating that the DBD patterns was not an artefact of the dependence structure in the

819     data.

820          The effect of community diversity on focal cluster diversity was also tested across

821     different environments analyzed separately. We modelled this relationship with linear,

822     quadratic and cubic fits, and compared those models based on the adjusted $R^2$ (**Figure 2**

823     **supplement 11**).

824

825     **DBD variation across environments**

826    We tested the variation of focal lineage diversity slopes across different environments by

827    including EMPO 3 biome type as a fixed effect. We fitted a GLMM with the interaction

828    between community diversity and environment type as a predictor of focal lineage

829    diversity. All other random effects on intercept and slope were kept as in the previous

830    GLMMs (**Figure 3**, **Supplementary Data file 1 Section 3**). DBD variation across

831    environments was tested for Family:Order, Order:Class and Class:Phylum taxonomic

832    ratios, as diversity slope variation by environment was statistically significant

833    (likelihood-ratio test, $P<0.05$) for these ratios in the taxonomy based models (**Table 1**).

834

835    **Abiotic effects**

836    To test for the relative effect of biotic and abiotic environmental variables on focal

837    lineage diversity across different taxonomic ratios, we used a separate GLMM, with

838    Poisson error distribution and a log link function, for every ratio. We fitted the GLMM on

839    a subset (~10%) of the whole dataset, 192 samples (from water: saline (19) and non-

840    saline (44), surface: saline (42) and non-saline (19), sediment: saline (22) and non-saline

841    (31), soil (8) and plant rhizosphere (7)), for which measurements of four key abiotic

842    variables (temperature, pH, latitude and elevation) were available. As predictors of focal

843    lineage diversity (fixed effects), we included non-focal community diversity and abiotic

844    variables, as well as their interactions. All predictors were standardized to a mean of zero

845    and a standard deviation of one to obtain comparable estimates. The GLMM had the

846    same random effects as in the previous analysis, but only on the intercept for simplicity

847    (**Table 4**, **Supplementary file 1 section 4**).

848

849 **Soil dataset analysis**

850 We used the Delgado-Baquerizo et al. 2018 soil microbiome survey (237 samples from

851 18 countries) to further test the relative impacts of biotic versus abiotic drivers of

852 diversity. Raw data and abiotic measurements were downloaded from Figshare

853 (https://figshare.com/s/82a2d3f5d38ace925492; DOI: 10.6084/m9.figshare.5611321).

854 16S bioinformatic processing was performed using QIIME2 and Deblur with the same

855 protocol as in Thompson et al. 2017. Raw data 16S rRNA gene (V3-V4 region), were

856 processed by trimming the primers (341F/805R primer set) with qiime cutadapt trim-

857 paired, then merged using qiime vsearch join-pairs. Sequences were quality filtered and

858 denoised using Deblur with a trimming length of 400bp. The resulting 400-bp Deblur

859 BIOM table was filtered to keep only ASVs with at least 25 reads total over all samples

860 and rarefied to a depth of 5000. Taxonomy was assigned with a Naive Bayes classifier

861 trained on the V4-V3 region of 99% OTU Greengenes 13.8 sequences with qiime feature-

862 classifier. We obtained a final dataset of 186 samples and 24,252 ASVs which was used

863 as input for all statistical analysis as in the EMP dataset analysis. This data set included

864 14 environmental factors: aridity index (Aridity_Index), minimum and maximum

865 temperature (MINT and MAXT), precipitation seasonality (PSEA), mean diurnal

866 temperature range (MDR), ultra-violet (UV) radiation (UV_Light), net primary

867 productivity (NPP2003_2015), soil texture (Clay_silt), pH; total C (Soil_C), N (Soil_N)

868 and P (Soil_P ) concentrations, C:N ratio (Soil_C_N_ratio) and Latitude.

869 We used a separate GLMM with Poisson error distribution and a log link function to test

870 for the effect of biotic (non-focal community diversity) and abiotic environmental

871 variables on focal lineage diversity (*e.g.* the ASV:Genus ratio for a focal genus), across

872    different taxonomic ratios. We defined non-focal taxa diversity and abiotic variables as

873    predictors (fixed effects) and the lineage identity as a random effect.

874    We also fitted the same model but with the first three principal components (PCs) from

875    the principal component analysis (PCA, rda function, vegan R package) of the abiotic

876    variables (a matrix of 237 samples (rows) by 14 abiotic variables (columns)), as well as

877    the interactions between diversity and each PC, and the interaction between PCs as

878    predictors (fixed effects).

879    Because of possible non-linear relationships between abiotic variables and diversity,

880    GLMMs were fitted with a linear and a quadratic term for every abiotic variable. The

881    quadratic terms were not significant, except for the ASV:genus ratio (**Table 5**; likelihood-

882    ratio test, $P < 2.2e\text{-}16$). The interaction terms were not significant except the interaction

883    between diversity and PCs at Family:Order ratio (likelihood-ratio test, $P= 2.182e\text{-}05$;

884    **Table 5**, **Supplementary file 4**).

885

886    **Defining residents, generalists, and migrants.** We defined a genus-level community

887    composition matrix as a matrix of 1128 genera (rows) by 17 environments (columns),

888    with the matrix entries indicating the percentage of samples from a given environment in

889    which each genus is present. We clustered the environmental samples based on their

890    genus-level community composition using fuzzy $k$-means clustering. The clustering

891    (cmeans function, package e1071 in R) was done on the 'hellinger' transformed data

892    (decostand function, vegan R package). To identify resident genera to each cluster, we

893    used indicator species analysis (Dufrene & Legendre, 1997) as implemented in the indval

894  function (labdsv R package). We defined residents as genera with indval indices between

895  0.4 and 0.9, with permutation test $P < 0.05$. Genera not associated with any cluster were

896  considered generalists. We used principal component analysis (PCA) on the community

897  composition matrix to visualize the clustering and the indicator genera (rda function,

898  vegan R package) (**Figure 4**). We then ran a separate GLMM for each environmental

899  cluster, with resident genus-level diversity (number of non-focal genera) as a predictor of

900  focal genus diversity (ASV:Genus ratio) for resident, migrant (residents of one cluster

901  found in a different cluster) and generalist genera. The fixed effect was specified as the

902  interaction between diversity and a factor defining the genus-cluster association (with

903  three levels: resident, migrant and generalist). Random effects on intercept and slope

904  were kept as in the GLMMs described above.

905

906  **Genome size analysis.** We chose a subset of genera represented by one or more

907  sequenced genomes in the NCBI microbial genomes database

908  (https://www.ncbi.nlm.nih.gov/genome/browse#!/prokaryotes/). For these genera, a

909  representative genome size was assigned by selecting the genome with the lowest number

910  of scaffolds (if no closed genomes were available) (**Supplementary file 6**). If multiple

911  genomes were available with the same level of completion, the largest genome size was

912  used, as smaller genomes could be artefacts of incomplete assembly which would bias the

913  mean and median downward. Moreover, given the deletional bias in bacterial genomes

914  (Kuo & Ochman, 2009a), the largest genome is likely more reflective of the ancestral

915  genome size of the genus. Only genera with two or more ASVs in at least one sample

916  were included in the analysis. Intracellular symbionts were excluded. We fitted a GLMM

917    on the subset of data with known genome size (576 genera, ranging from ~1 to 15 Mbp)

918    with the interaction between community diversity and genome size as a predictor of focal

919    lineage diversity at the ASV:Genus level. All the other random effects on intercept and

920    slope were kept as in the previous GLMMs (**Supplementary file 1 section 6)**.

921

922    **References**

923    Alonso, D., & McKane, A. J. (2004). Sampling Hubbell's neutral theory of biodiversity:

924        Sampling neutral theory. *Ecology Letters*, *7*(10), 901–910.

925        https://doi.org/10.1111/j.1461-0248.2004.00640.x

926    Auguet, J.-C., Barberan, A., & Casamayor, E. O. (2010). Global ecological patterns in

927        uncultured Archaea. *The ISME Journal*, *4*(2), 182–190.

928        https://doi.org/10.1038/ismej.2009.109

929    Azaele, S., Suweis, S., Grilli, J., Volkov, I., Banavar, J. R., & Maritan, A. (2016).

930        Statistical mechanics of ecological systems: Neutral theory and beyond. *Reviews of*

931        *Modern Physics*, *88*(3), 035003. https://doi.org/10.1103/RevModPhys.88.035003

932    Bailey, S. F., Dettman, J. R., Rainey, P. B., & Kassen, R. (2013). Competition both drives

933        and impedes diversification in a model adaptive radiation. *Proceedings. Biological*

934        *Sciences / The Royal Society*, *280*(1766), 20131253.

935        https://doi.org/10.1098/rspb.2013.1253

936    Barberán, A., Ramirez, K. S., Leff, J. W., Bradford, M. a., Wall, D. H., & Fierer, N.

937        (2014). Why are some microbes more ubiquitous than others? Predicting the habitat

938        breadth of soil bacteria. *Ecology Letters*, *17*(7), 794–802.

939        https://doi.org/10.1111/ele.12282

940    Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects

941        Models Using lme4. *Journal of Statistical Software*, *67*(1), 1–48.

942        https://doi.org/10.18637/jss.v067.i01

943    Brockhurst, M. A., Buckling, A., & Rainey, P. B. (2005). The effect of a bacteriophage

944        on diversification of the opportunistic bacterial pathogen, *Pseudomonas aeruginosa*.

945      *Proceedings. Biological Sciences / The Royal Society*, *272*(1570), 1385–1391.

946      https://doi.org/10.1098/rspb.2005.3086

947 Brockhurst, M. A., Colegrave, N., Hodgson, D. J., & Buckling, A. (2007). Niche

948      occupation limits adaptive radiation in experimental microcosms. *PloS One*, *2*(2),

949      e193. https://doi.org/10.1371/journal.pone.0000193

950 Calcagno, V., Jarne, P., Loreau, M., Mouquet, N., & David, P. (2017). Diversity spurs

951      diversification in ecological communities. *Nature Communications*, *8*, 15810.

952      https://doi.org/10.1038/ncomms15810

953 Coyte, K. Z., Schluter, J., & Foster, K. R. (2015). The ecology of the microbiome:

954      Networks, competition, and stability. *Science*, *350*(6261), 663–666.

955      https://doi.org/10.1126/science.aad2602

956 Czárán, T. L., Hoekstra, R. F., & Pagie, L. (2002). Chemical warfare between microbes

957      promotes biodiversity. *Proceedings of the National Academy of Sciences of the*

958      *United States of America*, *99*(2), 786–790. https://doi.org/10.1073/pnas.012399899

959 Delgado-Baquerizo, M., Oliverio, A. M., Brewer, T. E., Benavent-González, A.,

960      Eldridge, D. J., Bardgett, R. D., Maestre, F. T., Singh, B. K., & Fierer, N. (2018). A

961      global atlas of the dominant bacteria found in soil. *Science*, *359*(6373), 320–325.

962      https://doi.org/10.1126/science.aap9516

963 de Wit, R., & Bouvier, T. (2006). "Everything is everywhere, but, the environment

964      selects"; what did Baas Becking and Beijerinck really say? *Environmental*

965      *Microbiology*, *8*(4), 755–758. https://doi.org/10.1111/j.1462-2920.2006.01017.x

966 Dufrene, M., & Legendre, P. (1997). Species Assemblages and Indicator Species: The

967      Need for a Flexible Asymmetrical Approach. *Ecological Monographs*, *67*(3), 345–

968    366. https://doi.org/10.2307/2963459

969    Elston, D. A., Moss, R., Boulinier, T., Arrowsmith, C., & Lambin, X. (2001). Analysis of

970        aggregation, a worked example: numbers of ticks on red grouse chicks.

971        *Parasitology*, *122*(Pt 5), 563–569. https://doi.org/10.1017/s0031182001007740

972    Elton, C. (1946). Competition and the Structure of Ecological Communities. *The Journal*

973        *of Animal Ecology*, *15*(1), 54–68. https://doi.org/10.2307/1625

974    Emerson, B. C., & Kolm, N. (2005). Species diversity can drive speciation. *Nature*,

975        *434*(7036), 1015–1017. https://doi.org/10.1038/nature03450

976    Falkowski, P. G., Fenchel, T., & Delong, E. F. (2008). The microbial engines that drive

977        Earth's biogeochemical cycles. *Science*, *320*(5879), 1034–1039.

978        https://doi.org/10.1126/science.1153213

979    Gause, G. F. (2003). *The Struggle for Existence* (Williams & Wilkins, Baltimore, 1934).

980    Gómez, P., & Buckling, A. (2013). Real-time microbial adaptive diversification in soil.

981        *Ecology Letters*, *16*(5), 650–655. https://doi.org/10.1111/ele.12093

982    Gotelli, N. J., & Colwell, R. K. (2001). Quantifying biodiversity: procedures and pitfalls

983        in the measurement and comparison of species richness. *Ecology Letters*, *4*(4), 379–

984        391. https://doi.org/10.1046/j.1461-0248.2001.00230.x

985    Gotelli, N. J., & McGill, B. J. (2006). Null Versus Neutral Models: What's The

986        Difference? *Ecography*, *29*(5), 793–800. https://doi.org/10.1111/j.2006.0906-

987        7590.04714.x

988    Harris, K., Parsons, T. L., Ijaz, U. Z., Lahti, L., Holmes, I., & Quince, C. (2017). Linking

989        Statistical and Ecological Theory: Hubbell's Unified Neutral Theory of Biodiversity

990        as a Hierarchical Dirichlet Process. *Proceedings of the IEEE*, *105*(3), 516–529.

991     https://doi.org/10.1109/JPROC.2015.2428213

992   Hibbing, M. E., Fuqua, C., Parsek, M. R., & Peterson, S. B. (2010). Bacterial

993     competition: surviving and thriving in the microbial jungle. *Nature Reviews*

994     *Microbiology*, *8*(1), 15–25. https://doi.org/10.1038/nrmicro2259

995   Hubbell, S. P. (2001). *The Unified Neutral Theory of Biodiversity and Biogeography*.

996     Princeton University Press.

997   Hug, L. A., Baker, B. J., Anantharaman, K., Brown, C. T., Probst, A. J., Castelle, C. J.,

998     Butterfield, C. N., Hernsdorf, A. W., Amano, Y., Kotaro, I., Suzuki, Y., Dudek, N.,

999     Relman, D. A., Finstad, K. M., Amundson, R., Thomas, B. C., & Banfield, J. F.

1000     (2016). A new view of the tree and life's diversity. *Nature Microbiology*, 1, 16048

1001     https://doi.org/10.1038/nmicrobiol.2016.48

1002   Jarvinen, O. (1982). Species-To-Genus Ratios in Biogeography: A Historical Note.

1003     *Journal of Biogeography*, *9*(4), 363–370. https://doi.org/10.2307/2844723

1004   Johnson, P. C. (2014). Extension of Nakagawa & Schielzeth's R2GLMM to random

1005     slopes models. *Methods in Ecology and Evolution / British Ecological Society*, *5*(9),

1006     944–946. https://doi.org/10.1111/2041-210X.12225

1007   Jousset, A., Eisenhauer, N., Merker, M., Mouquet, N., & Scheu, S. (2016). High

1008     functional diversity stimulates diversification in experimental microbial

1009     communities. *Science Advances*, *2*(6), e1600124.

1010     https://doi.org/10.1126/sciadv.1600124

1011   Kastman, E. K., Kamelamela, N., Norville, J. W., Cosetta, C. M., Dutton, R. J., & Wolfe,

1012     B. E. (2016). Biotic Interactions Shape the Ecological Distributions of

1013     Staphylococcus Species. *mBio*, *7*(5). https://doi.org/10.1128/mBio.01157-16

1014    Kennedy, A. C., & de Luna, L. Z. (2005). Rhizhosphere. In D. Hillel (Ed.), *Encyclopedia*

1015    *of Soils in the Environment* (pp. 399–406). Elsevier. https://doi.org/10.1016/B0-12-

1016    348530-4/00163-6

1017    Konstantinidis, K. T., & Tiedje, J. M. (2005). Towards a genome-based taxonomy for

1018    prokaryotes. *Journal of Bacteriology*, *187*(18), 6258–6264.

1019    Kuo, C.-H., & Ochman, H. (2009a). Deletional bias across the three domains of life.

1020    *Genome Biology and Evolution*, *1*, 145–152. https://doi.org/10.1093/gbe/evp016

1021    Kuo, C.-H., & Ochman, H. (2009b). Inferring clocks when lacking rocks: the variable

1022    rates of molecular evolution in bacteria. *Biology Direct*, *4*, 35.

1023    https://doi.org/10.1186/1745-6150-4-35

1024    Laland, K. N., Odling-Smee, F. J., & Feldman, M. W. (1999). Evolutionary consequences

1025    of niche construction and their implications for ecology. *Proceedings of the National*

1026    *Academy of Sciences of the United States of America*, *96*(18), 10242–10247.

1027    https://doi.org/10.1073/pnas.96.18.10242

1028    Lapierre, P., & Gogarten, J. P. (2009). Estimating the size of the bacterial pan-genome.

1029    *Trends in Genetics: TIG*, *25*(3), 107–110. https://doi.org/10.1016/j.tig.2008.12.004

1030    Lauber, C. L., Hamady, M., Knight, R., & Fierer, N. (2009). Soil pH as a predictor of soil

1031    bacterial community structure at the continental scale: a pyrosequencing-based

1032    assessment. *Applied and Environmental Microbiology*. 75, 5111-5120.

1033    http://aem.asm.org/content/early/2009/06/05/AEM.00335-09.short

1034    Li, L., & Ma, Z. S. (2016). Testing the Neutral Theory of Biodiversity with Human

1035    Microbiome Datasets. *Scientific Reports*, *6*, 31448.

1036    https://doi.org/10.1038/srep31448

1037 Lindow, S. E., & Brandl, M. T. (2003). Microbiology of the phyllosphere. *Applied and*

1038 *Environmental Microbiology*, *69*(4), 1875–1883.

1039 https://doi.org/10.1128/aem.69.4.1875-1883.2003

1040 Louca, S., Mazel, F., Doebeli, M., & Parfrey, L. W. (2019). A census-based estimate of

1041 Earth's bacterial and archaeal diversity. *PLoS Biology*, *17*(2), e3000106.

1042 https://doi.org/10.1371/journal.pbio.3000106

1043 Louca, S., & Pennell, M. W. (2020). Extant timetrees are consistent with a myriad of

1044 diversification histories. *Nature*, *580*(7804), 502–505.

1045 https://doi.org/10.1038/s41586-020-2176-1

1046 Louca, S., Shih, P. M., Pennell, M. W., Fischer, W. W., Parfrey, L. W., & Doebeli, M.

1047 (2018). Bacterial diversification through geological time. *Nature Ecology &*

1048 *Evolution*, *2*(9), 1458–1467. https://doi.org/10.1038/s41559-018-0625-0

1049 Lozupone, C. A., & Knight, R. (2007). Global patterns in bacterial diversity. *Proceedings*

1050 *of the National Academy of Sciences of the United States of America*, *104*(27),

1051 11436–11440. https://doi.org/10.1073/pnas.0611525104

1052 Marshall, C. R. (2017). Five palaeobiological laws needed to understand the evolution of

1053 the living biota. *Nature Ecology & Evolution*, *1*(6), 165.

1054 https://doi.org/10.1038/s41559-017-0165

1055 Meyer, J. R., & Kassen, R. (2007). The effects of competition and predation on

1056 diversification in a model adaptive radiation. *Nature*, *446*(7134), 432–435.

1057 http://www.nature.com/doifinder/10.1038/nature05599

1058 Morris, J. J., & Lenski, R. E. (2012). The Black Queen Hypothesis: evolution of

1059 dependencies through adaptive gene loss. *mBio*. 3, e00036-12

1060    http://mbio.asm.org/content/3/2/e00036-12.short

1061    Nakagawa, S., & Schielzeth, H. (2013). A general and simple method for obtaining R 2

1062    from generalized linear mixed-effects models. *Methods in Ecology and Evolution /*

1063    *British Ecological Society*, *4*(2), 133–142. https://doi.org/10.1111/j.2041-

1064    210x.2012.00261.x

1065    Needham, D. M., & Fuhrman, J. A. (2016). Pronounced daily succession of

1066    phytoplankton , archaea and bacteria following a spring bloom. *Nature*

1067    *Microbiology*, 1, 16005. https://doi.org/10.1038/NMICROBIOL.2016.5

1068    Palmer, M. W., & Maurer, T. A. (1997). Does Diversity Beget Diversity? A Case Study

1069    of Crops and Weeds. *Journal of Vegetation Science*, *8*(2), 235–240.

1070    https://doi.org/10.2307/3237352

1071    Parks, D. H., Chuvochina, M., Waite, D. W., Rinke, C., Skarshewski, A., Chaumeil, P.-

1072    A., & Hugenholtz, P. (2018). A standardized bacterial taxonomy based on genome

1073    phylogeny substantially revises the tree of life. *Nature Biotechnology*, *36*(10), 996–

1074    1004. https://doi.org/10.1038/nbt.4229

1075    Pennekamp, F., Pontarp, M., Tabi, A., Altermatt, F., Alther, R., Choffat, Y., Fronhofer, E.

1076    A., Ganesanandamoorthy, P., Garnier, A., Griffiths, J. I., Greene, S., Horgan, K.,

1077    Massie, T. M., Mächler, E., Palamara, G. M., Seymour, M., & Petchey, O. L.

1078    (2018). Biodiversity increases and decreases ecosystem stability. *Nature*, *563*(7729),

1079    109–112. https://doi.org/10.1038/s41586-018-0627-8

1080    Power, J. F., Carere, C. R., Lee, C. K., Wakerley, G. L. J., Evans, D. W., Button, M.,

1081    White, D., Climo, M. D., Hinze, A. M., Morgan, X. C., McDonald, I. R., Cary, S.

1082    C., & Stott, M. B. (2018). Microbial biogeography of 925 geothermal springs in

1083        New Zealand. *Nature Communications*, *9*(1), 2876. https://doi.org/10.1038/s41467-

1084        018-05020-y

1085  Price, T. D., Hooper, D. M., Buchanan, C. D., Johansson, U. S., Tietze, D. T., Alström,

1086        P., Olsson, U., Ghosh-Harihar, M., Ishtiaq, F., Gupta, S. K., Martens, J., Harr, B.,

1087        Singh, P., & Mohan, D. (2014). Niche filling slows the diversification of Himalayan

1088        songbirds. *Nature*. 509, 222-225  https://doi.org/10.1038/nature13272

1089  Rabosky, D. L., Chang, J., Title, P. O., Cowman, P. F., Sallan, L., Friedman, M.,

1090        Kaschner, K., Garilao, C., Near, T. J., Coll, M., & Alfaro, M. E. (2018). An inverse

1091        latitudinal gradient in speciation rate for marine fishes. *Nature*, *559*(7714), 392–395.

1092        https://doi.org/10.1038/s41586-018-0273-1

1093  Rabosky, D. L., & Hurlbert, A. H. (2015). Species richness at continental scales is

1094        dominated by ecological limits. *The American Naturalist*, *185*(5), 572–583.

1095        https://doi.org/10.1086/680850

1096  San Roman, M., & Wagner, A. (2018). An enormous potential for niche construction

1097        through bacterial cross-feeding in a homogeneous environment. *PLoS*

1098        *Computational Biology*, *14*(7), e1006340.

1099        https://doi.org/10.1371/journal.pcbi.1006340

1100  Schluter, D., & Pennell, M. W. (2017). Speciation gradients and the distribution of

1101        biodiversity. *Nature*, *546*(7656), 48–55. https://doi.org/10.1038/nature22897

1102  Sender, R., Fuchs, S., & Milo, R. (2016). Revised Estimates for the Number of Human

1103        and Bacteria Cells in the Body. *PLoS Biology*, *14*(8), e1002533.

1104        https://doi.org/10.1371/journal.pbio.1002533

1105  Seth, E. C., & Taga, M. E. (2014). Nutrient cross-feeding in the microbial world.

*Frontiers in Microbiology*, *5*, 350. https://doi.org/10.3389/fmicb.2014.00350

Sogin, M. L., Morrison, H. G., Huber, J. A., Welch, D. M., Huse, S. M., Neal, P. R., Arrieta, J. M., & Herndl, G. J. (2006). Microbial diversity in the deep sea and the underexplored "rare biosphere." *Proceedings of the National Academy of Sciences of the United States of America*, *103*(32), 12115–12120.

Sriswasdi, S., Yang, C.-C., & Iwasaki, W. (2017). Generalist species drive microbial dispersion and evolution. *Nature Communications*, *8*(1), 1162. https://doi.org/10.1038/s41467-017-01265-1

Sunagawa, S., Coelho, L. P., Chaffron, S., Kultima, J. R., Labadie, K., Salazar, G., Djahanschiri, B., Zeller, G., Mende, D. R., Alberti, A., Cornejo-Castillo, F. M., Costea, P. I., Cruaud, C., d'Ovidio, F., Engelen, S., Ferrera, I., Gasol, J. M., Guidi, L., Hildebrand, F., … Bork, P. (2015). Ocean plankton. Structure and function of the global ocean microbiome. *Science*, *348*(6237), 1261359. https://doi.org/10.1126/science.1261359

Thompson, L. R., Sanders, J. G., McDonald, D., Amir, A., Ladau, J., Locey, K. J., Prill, R. J., Tripathi, A., Gibbons, S. M., Ackermann, G., Navas-Molina, J. A., Janssen, S., Kopylova, E., Vázquez-Baeza, Y., González, A., Morton, J. T., Mirarab, S., Xu, Z. Z., Jiang, L., … Zhao, H. (2017). A communal catalogue reveals Earth's multiscale microbial diversity. *Nature*, 551, 457-463. https://doi.org/10.1038/nature24621

Vos, M. (2011). A species concept for bacteria based on adaptive divergence. *Trends in Microbiology*, *19*(1), 1–7. https://doi.org/10.1016/j.tim.2010.10.003

Whitman, W. B., Coleman, D. C., & Wiebe, W. J. (1998). Prokaryotes: the unseen majority. *Proceedings of the National Academy of Sciences of the United States of*

1129    *America*, *95*(12), 6578–6583. https://doi.org/10.1073/pnas.95.12.6578

1130    Whittaker, R. H. (1972). Evolution and Measurement of Species Diversity. *Taxon*,

1131    *21*(2/3), 213–251. https://doi.org/10.2307/1218190

Slope estimate

**A. Proteobacteria**

Number of Proteobacteria classes (y-axis) vs Number of non-Proteobacteria phyla (x-axis)

| Panel | P-value | Slope |
|---|---|---|
| Aerosol (non-saline) | P= 1.04e−14 | 0.15 |
| Animal corpus | P= 1.27e−21 | 0.18 |
| Animal distal gut | P= 1.44e−11 | 0.23 |
| Animal proximal gut | P= 1.03e−08 | 0.30 |
| Animal secretion | P= 8.91e−17 | 0.16 |
| Animal surface | P= 8.98e−20 | 0.15 |
| Hypersaline (saline) | P= 1.00e+00 | 0.09 |
| Plant corpus | P= 4.59e−28 | 0.21 |
| Plant rhizosphere | P= 1.71e−04 | 0.03 |
| Plant surface | P= 1.51e−02 | 0.06 |
| Sediment (non-saline) | P= 6.98e−01 | −0.03 |
| Sediment (saline) | P= 1.48e−12 | 0.06 |
| Soil (non-saline) | P= 8.10e−06 | 0.05 |
| Surface (non-saline) | P= 1.26e−22 | 0.18 |
| Surface (saline) | P= 3.64e−18 | 0.09 |
| Water (non-saline) | P= 2.60e−10 | 0.07 |
| Water (saline) | P= 2.52e−06 | 0.05 |

# B. Bacteroidetes

# C. Actinobacteria



Figure axes: Number of Actinobacteria classes (y-axis) vs Number of non-Actinobacteria phyla (x-axis)

Panels:

**Aerosol (non-saline)** — P= 6.15e-13, Slope: 0.21
**Animal corpus** — P= 1.56e-08, Slope: 0.11
**Animal distal gut** — P= 2.49e-03, Slope: 0.10
**Animal proximal gut** — P= 2.91e-02, Slope: 0.15
**Animal secretion** — P= 1.51e-21, Slope: 0.21
**Animal surface** — P= 7.10e-29, Slope: 0.28
**Hypersaline (saline)** — P= 1.00e+00, Slope: -0.02
**Plant corpus** — P= 1.18e-24, Slope: 0.26
**Plant rhizosphere** — P= 8.72e-24, Slope: 0.09
**Plant surface** — P= 3.55e-13, Slope: 0.16
**Sediment (non-saline)** — P= 1.00e+00, Slope: -0.00
**Sediment (saline)** — P= 6.14e-05, Slope: 0.06
**Soil (non-saline)** — P= 1.55e-01, Slope: 0.04
**Surface (non-saline)** — P= 2.96e-04, Slope: 0.12
**Surface (saline)** — P= 4.02e-14, Slope: 0.11
**Water (non-saline)** — P= 3.51e-13, Slope: 0.09
**Water (saline)** — P= 8.70e-17, Slope: 0.08

**A. Gammaproteobacteria**

Number of Gammaproteobacteria orders (y-axis)

Number of non-Gammaproteobacteria classes (x-axis)

Aerosol (non-saline): P= 1.56e−16, Slope: 0.11
Animal corpus: P= 2.22e−11, Slope: 0.11
Animal distal gut: P= 1.28e−08, Slope: 0.17
Animal proximal gut: P= 6.35e−09, Slope: 0.22
Animal secretion: P= 1.67e−23, Slope: 0.16
Animal surface: P= 1.22e−24, Slope: 0.11
Hypersaline (saline): P= 1.00e+00, Slope: 0.05
Plant corpus: P= 9.29e−30, Slope: 0.11
Plant rhizosphere: P= 3.35e−04, Slope: 0.02
Plant surface: P= 4.28e−09, Slope: 0.09
Sediment (non-saline): P= 3.44e−04, Slope: 0.05
Sediment (saline): P= 2.11e−17, Slope: 0.06
Soil (non-saline): P= 1.08e−09, Slope: 0.06
Surface (non-saline): P= 8.67e−09, Slope: 0.08
Surface (saline): P= 4.70e−18, Slope: 0.06
Water (non-saline): P= 2.38e−19, Slope: 0.07
Water (saline): P= 3.28e−05, Slope: 0.06

**B. Alphaproteobacteria**

Number of Alphaproteobacteria orders (y-axis)

Number of non-Alphaproteobacteria classes (x-axis)

Panels:
- Aerosol (non-saline): P= 1.70e−12, Slope: 0.08
- Animal corpus: P= 1.01e−18, Slope: 0.13
- Animal distal gut: P= 3.45e−05, Slope: 0.13
- Animal proximal gut: P= 3.45e−05, Slope: 0.15
- Animal secretion: P= 2.85e−13, Slope: 0.12
- Animal surface: P= 1.71e−30, Slope: 0.11
- Hypersaline (saline): P= 5.46e−01, Slope: 0.08
- Plant corpus: P= 2.72e−34, Slope: 0.13
- Plant rhizosphere: P= 1.40e−04, Slope: −0.01
- Plant surface: P= 1.04e−15, Slope: 0.11
- Sediment (non-saline): P= 5.04e−02, Slope: 0.04
- Sediment (saline): P= 2.61e−21, Slope: 0.07
- Soil (non-saline): P= 9.70e−12, Slope: 0.05
- Surface (non-saline): P= 2.17e−27, Slope: 0.10
- Surface (saline): P= 3.37e−10, Slope: 0.03
- Water (non-saline): P= 2.25e−08, Slope: 0.03
- Water (saline): P= 8.63e−09, Slope: 0.05

**C. Actinobacteria**

Figure panels with title "C. Actinobacteria". X-axis: Number of non-Actinobacteria classes. Y-axis: Number of Actinobacteria orders.

**Aerosol (non-saline)** — P= 1.05e-09, Slope: 0.03

**Animal corpus** — P= 1.00e+00, Slope: 0.00

**Animal distal gut** — P= 1.00e+00, Slope: 0.01

**Animal proximal gut** — P= 1.00e+00, Slope: -0.01

**Animal secretion** — P= 1.18e-03, Slope: 0.01

**Animal surface** — P= 2.09e-10, Slope: 0.02

**Hypersaline (saline)** — P= 3.68e-01, Slope: 0.02

**Plant corpus** — P= 4.66e-07, Slope: 0.01

**Plant rhizosphere** — P= 5.81e-03, Slope: 0.01

**Plant surface** — P= 5.75e-03, Slope: 0.01

**Sediment (non-saline)** — P= 5.08e-03, Slope: 0.01

**Sediment (saline)** — P= 1.21e-03, Slope: 0.01

**Soil (non-saline)** — P= 3.62e-02, Slope: 0.01

**Surface (non-saline)** — P= 1.25e-02, Slope: 0.01

**Surface (saline)** — P= 3.42e-12, Slope: 0.01

**Water (non-saline)** — P= 1.53e-09, Slope: 0.02

**Water (saline)** — P= 8.92e-06, Slope: 0.01

**A. Actinomycetales**

Figure showing scatter plots of Number of Actinomycetales families (y-axis) versus Number of non-Actinomycetales orders (x-axis) across different environments.

- **Aerosol (non-saline):** P= 6.76e−12, Slope: 0.18
- **Animal corpus:** P= 7.77e−16, Slope: 0.14
- **Animal distal gut:** P= 3.43e−07, Slope: 0.16
- **Animal proximal gut:** P= 1.51e−12, Slope: 0.24
- **Animal secretion:** P= 9.13e−37, Slope: 0.21
- **Animal surface:** P= 2.33e−21, Slope: 0.22
- **Hypersaline (saline):** P= 2.53e−02, Slope: 0.09
- **Plant corpus:** P= 1.03e−26, Slope: 0.23
- **Plant rhizosphere:** P= 1.00e+00, Slope: 0.00
- **Plant surface:** P= 3.45e−11, Slope: 0.16
- **Sediment (non-saline):** P= 8.86e−01, Slope: 0.02
- **Sediment (saline):** P= 1.00e+00, Slope: 0.01
- **Soil (non-saline):** P= 1.99e−06, Slope: 0.06
- **Surface (non-saline):** P= 1.00e+00, Slope: 0.02
- **Surface (saline):** P= 1.20e−04, Slope: 0.02
- **Water (non-saline):** P= 4.14e−07, Slope: 0.05
- **Water (saline):** P= 5.29e−08, Slope: 0.05

**B. Flavobacteriales**

*(Y-axis)* Number of Flavobacteriales families

*(X-axis)* Number of non-Flavobacteriales orders

Aerosol (non-saline): P= 1.50e−03, Slope: 0.02
Animal corpus: P= 8.63e−11, Slope: 0.05
Animal distal gut: P= 7.67e−01, Slope: 0.03
Animal proximal gut: P= 8.87e−03, Slope: 0.02
Animal secretion: P= 2.90e−05, Slope: 0.02
Animal surface: P= 8.34e−07, Slope: 0.02
Hypersaline (saline): P= 1.00e+00, Slope: 0.01
Plant corpus: P= 1.51e−02, Slope: 0.01
Plant rhizosphere: P= 1.63e−06, Slope: −0.01
Plant surface: P= 9.24e−02, Slope: 0.01
Sediment (non-saline): P= 3.78e−01, Slope: 0.01
Sediment (saline): P= 7.27e−06, Slope: 0.01
Soil (non-saline): P= 1.00e+00, Slope: 0.00
Surface (non-saline): P= 5.91e−02, Slope: 0.01
Surface (saline): P= 4.39e−08, Slope: 0.01
Water (non-saline): P= 1.24e−02, Slope: 0.01
Water (saline): P= 1.00e+00, Slope: −0.00

**C. Rhizobiales**

Figure: Scatter plots showing Number of Rhizobiales families (y-axis) versus Number of non-Rhizobiales orders (x-axis) across different environments.

Aerosol (non-saline): P= 2.49e-15, Slope: 0.09
Animal corpus: P= 5.75e-11, Slope: 0.06
Animal distal gut: P= 3.01e-01, Slope: 0.08
Animal proximal gut: P= 3.54e-06, Slope: 0.13
Animal secretion: P= 2.07e-08, Slope: 0.09
Animal surface: P= 4.06e-19, Slope: 0.10
Hypersaline (saline): P= 1.10e-01, Slope: 0.14
Plant corpus: P= 1.10e-25, Slope: 0.13
Plant rhizosphere: P= 3.50e-02, Slope: -0.01
Plant surface: P= 2.21e-11, Slope: 0.04
Sediment (non-saline): P= 1.00e+00, Slope: -0.00
Sediment (saline): P= 1.16e-05, Slope: 0.02
Soil (non-saline): P= 4.85e-08, Slope: 0.04
Surface (non-saline): P= 1.41e-01, Slope: 0.02
Surface (saline): P= 3.70e-12, Slope: 0.03
Water (non-saline): P= 2.02e-09, Slope: 0.03
Water (saline): P= 5.18e-05, Slope: 0.03

**A. Flavobacteriaceae**

X-axis (all panels): Number of non-Flavobacteriaceae families

Y-axis (all panels): Number of Flavobacteriaceae genera

Aerosol (non-saline): P= 6.02e-02, Slope: 0.01

Animal corpus: P= 1.00e+00, Slope: 0.01

Animal distal gut: P= 1.00e+00, Slope: 0.02

Animal proximal gut: P= 9.90e-02, Slope: 0.04

Animal secretion: P= 1.71e-13, Slope: 0.03

Animal surface: P= 1.43e-05, Slope: 0.02

Hypersaline (saline): P= 1.00e+00, Slope: 0.02

Plant corpus: P= 2.29e-02, Slope: 0.02

Plant rhizosphere: P= 9.70e-08, Slope: -0.01

Plant surface: P= 4.01e-01, Slope: 0.02

Sediment (non-saline): P= 2.91e-01, Slope: -0.01

Sediment (saline): P= 4.22e-09, Slope: 0.05

Soil (non-saline): P= 1.00e+00, Slope: -0.01

Surface (non-saline): P= 1.00e-03, Slope: 0.02

Surface (saline): P= 1.00e+00, Slope: 0.02

Water (non-saline): P= 1.00e+00, Slope: 0.00

Water (saline): P= 1.00e+00, Slope: -0.01

# B. Sphingomonadaceae



Figure B. Sphingomonadaceae. Scatter plots of the number of Sphingomonadaceae genera (y-axis) versus the number of non-Sphingomonadaceae families (x-axis) across environment types.

**Aerosol (non-saline):** P= 8.35e−10, Slope: 0.03

**Animal corpus:** P= 9.02e−12, Slope: 0.05

**Animal distal gut:** P= 2.25e−01, Slope: 0.03

**Animal proximal gut:** P= 8.21e−02, Slope: 0.02

**Animal secretion:** P= 4.26e−03, Slope: 0.03

**Animal surface:** P= 8.04e−11, Slope: 0.03

**Plant corpus:** P= 1.39e−39, Slope: 0.04

**Plant rhizosphere:** P= 1.00e+00, Slope: 0.00

**Plant surface:** P= 1.00e+00, Slope: 0.01

**Sediment (non-saline):** P= 4.68e−03, Slope: 0.02

**Sediment (saline):** P= 2.21e−03, Slope: 0.01

**Soil (non-saline):** P= 1.97e−04, Slope: 0.02

**Surface (non-saline):** P= 3.66e−08, Slope: 0.03

**Surface (saline):** P= 6.57e−01, Slope: 0.01

**Water (non-saline):** P= 8.58e−04, Slope: 0.02

**Water (saline):** P= 5.65e−04, Slope: 0.02

Number of non-Sphingomonadaceae families

# C. Verrucomicrobiaceae

**A. Pseudomonas**

Figure showing scatter plots of Number of Pseudomonas ASVs (y-axis) versus Number of non-Pseudomonas genera (x-axis) across 16 environmental panels.

- **Aerosol (non-saline)**: P= 3.68e-04, Slope: 0.02
- **Animal corpus**: P= 8.56e-04, Slope: 0.04
- **Animal distal gut**: P= 8.64e-01, Slope: 0.02
- **Animal proximal gut**: P= 7.46e-01, Slope: 0.03
- **Animal secretion**: P= 8.12e-02, Slope: 0.03
- **Animal surface**: P= 3.97e-05, Slope: 0.02
- **Plant corpus**: P= 5.31e-11, Slope: 0.01
- **Plant rhizosphere**: P= 1.00e+00, Slope: 0.01
- **Plant surface**: P= 1.00e+00, Slope: 0.00
- **Sediment (non-saline)**: P= 1.00e+00, Slope: 0.00
- **Sediment (saline)**: P= 1.00e+00, Slope: 0.01
- **Soil (non-saline)**: P= 2.56e-07, Slope: 0.02
- **Surface (non-saline)**: P= 6.68e-01, Slope: 0.01
- **Surface (saline)**: P= 1.00e+00, Slope: 0.01
- **Water (non-saline)**: P= 1.53e-05, Slope: 0.03
- **Water (saline)**: P= 1.00e+00, Slope: 0.02

**B. Planctomyces**

# C. Clostridium

**A. Model 1**

ASV:Genus

*m* = -0.005

number of genera

**B. Model 2**

**Not significant**

number of genera

**C. Model 3**

25 %    50 %    75 %    100 %

+DBD

*m* = +0.016

0 %

ASV:Genus

*m* = -0.012

+EC

*m* = -0.011

number of genera

**A. True ratio = 2 ASVs/genus**

Number of ASVs per genus vs Number of genera

**B. True ratio = 20 ASVs/genus**

Number of ASVs per genus vs Number of genera

**C. True ratio = 3 classes/phylum**

Number of classes per phylum vs Number of phyla

**D. True ratio = 10 classes/phylum**

Number of classes per phylum vs Number of phyla

**Rarefaction level (# of sequences sampled)**

- 1,000
- 2,000
- 3,000
- 4,000
- 5,000

**80% clusters** (x-axis)

**85% clusters** (y-axis)

| Panel | Adj R2 (grey) | Adj R2 (blue) | Adj R2 (red) |
|---|---|---|---|
| Aerosol (non-saline) | 0.02 | 0 | 0 |
| Animal corpus | 0 | 0.02 | 0.02 |
| Animal distal gut | 0.15 | 0.31 | 0.32 |
| Animal proximal gut | 0.01 | 0.06 | 0.06 |
| Animal secretion | 0.14 | 0.13 | 0.16 |
| Animal surface | 0.08 | 0.1 | 0.13 |
| Hypersaline (saline) | −0.09 | −0.18 | −0. |
| Plant corpus | 0.43 | 0.43 | 0.43 |
| Plant rhizosphere | 0 | 0.09 | 0.08 |
| Plant surface | 0.06 | 0.06 | 0.06 |
| Sediment (non-saline) | 0.19 | 0.35 | 0.44 |
| Sediment (saline) | 0.12 | 0.19 | 0.19 |
| Soil (non-saline) | 0.01 | 0.02 | 0.03 |
| Surface (non-saline) | 0.21 | 0.36 | 0.38 |
| Surface (saline) | 0.17 | 0.19 | 0.19 |
| Water (non-saline) | 0.11 | 0.14 | 0.14 |
| Water (saline) | 0.13 | 0.14 | 0.15 |

Aerosol (non-saline) — Adj R2 : 0.56 0.63 0.63

Animal corpus — Adj R2 : 0.07 0.07 0.07

Animal distal gut — Adj R2 : 0.53 0.53 0.53

Animal proximal gut — Adj R2 : 0.11 0.21 0.23

Animal secretion — Adj R2 : 0.4 0.4 0.39

Animal surface — Adj R2 : 0.26 0.26 0.26

Hypersaline (saline) — Adj R2 : 0.08 0.28 0.52

Plant corpus — Adj R2 : 0.42 0.42 0.43

Plant rhizosphere — Adj R2 : 0 0.19 0.2

Plant surface — Adj R2 : 0.07 0.07 0.06

Sediment (non-saline) — Adj R2 : 0 0.33 0.36

Sediment (saline) — Adj R2 : 0.31 0.41 0.4

Soil (non-saline) — Adj R2 : 0.13 0.25 0.27

Surface (non-saline) — Adj R2 : 0.55 0.73 0.75

Surface (saline) — Adj R2 : 0.05 0.05 0.04

Water (non-saline) — Adj R2 : 0.19 0.19 0.21

Water (saline) — Adj R2 : 0 0.02 0.02

90% clusters

85% clusters

Aerosol (non-saline) — Adj R2 : 0.52 0.63 0.63
Animal corpus — Adj R2 : 0.26 0.27 0.28
Animal distal gut — Adj R2 : 0.53 0.54 0.55
Animal proximal gut — Adj R2 : 0.26 0.34 0.34
Animal secretion — Adj R2 : 0.51 0.55 0.55
Animal surface — Adj R2 : 0.3 0.3 0.3
Hypersaline (saline) — Adj R2 : 0.39 0.35 0.59
Plant corpus — Adj R2 : 0.78 0.83 0.84
Plant rhizosphere — Adj R2 : 0.12 0.31 0.42
Plant surface — Adj R2 : 0.18 0.24 0.24
Sediment (non-saline) — Adj R2 : 0.03 0.34 0.34
Sediment (saline) — Adj R2 : 0.52 0.6 0.6
Soil (non-saline) — Adj R2 : 0.2 0.41 0.43
Surface (non-saline) — Adj R2 : 0.64 0.81 0.82
Surface (saline) — Adj R2 : 0.03 0.02 0.02
Water (non-saline) — Adj R2 : 0.17 0.17 0.17
Water (saline) — Adj R2 : 0.1 0.09 0.09

95% clusters

90% clusters

Aerosol (non-saline) — Adj R2 : 0.24 0.34 0.37
Animal corpus — Adj R2 : 0.11 0.11 0.12
Animal distal gut — Adj R2 : 0.29 0.3 0.3
Animal proximal gut — Adj R2 : 0.01 0.06 0.09
Animal secretion — Adj R2 : 0.26 0.27 0.29
Animal surface — Adj R2 : 0.32 0.34 0.33
Hypersaline (saline) — Adj R2 : 0.44 0.6 0.55
Plant corpus — Adj R2 : 0.26 0.27 0.26
Plant rhizosphere — Adj R2 : 0.5 0.5 0.54
Plant surface — Adj R2 : 0.23 0.24 0.24
Sediment (non-saline) — Adj R2 : 0.26 0.33 0.34
Sediment (saline) — Adj R2 : 0.34 0.38 0.42
Soil (non-saline) — Adj R2 : 0.16 0.37 0.44
Surface (non-saline) — Adj R2 : 0.49 0.59 0.6
Surface (saline) — Adj R2 : 0.02 0.03 0.05
Water (non-saline) — Adj R2 : 0.09 0.08 0.07
Water (saline) — Adj R2 : 0.04 0.03 0.03

97% clusters

95% clusters

**Aerosol (non−saline)** — Adj R2 : 0.06 0.1 0.17

**Animal corpus** — Adj R2 : −0.01 0 0.02

**Animal distal gut** — Adj R2 : 0.4 0.48 0.48

**Animal proximal gut** — Adj R2 : −0.01 −0.01 0

**Animal secretion** — Adj R2 : 0.01 0 0.05

**Animal surface** — Adj R2 : 0.04 0.03 0.03

**Hypersaline (saline)** — Adj R2 : 0.15 0.37 0.75

**Plant corpus** — Adj R2 : 0.11 0.1 0.1

**Plant rhizosphere** — Adj R2 : 0.37 0.38 0.56

**Plant surface** — Adj R2 : 0.02 0.01 0

**Sediment (non−saline)** — Adj R2 : 0.25 0.25 0.25

**Sediment (saline)** — Adj R2 : 0.11 0.12 0.16

**Soil (non−saline)** — Adj R2 : 0.06 0.25 0.32

**Surface (non−saline)** — Adj R2 : 0.12 0.28 0.36

**Surface (saline)** — Adj R2 : −0.01 −0.02 0

**Water (non−saline)** — Adj R2 : 0.01 0 0

**Water (saline)** — Adj R2 : 0.12 0.18 0.18

100% clusters

97% clusters