1    # Diversity begets diversity in microbiomes

2

3

4

5    **Authors**: Naïma Madi[1], Michiel Vos[2], Pierre Legendre[1] and B. Jesse Shapiro[1*]

6

7         1. Departement de sciences biologiques, Universite de Montreal, Canada

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

9            Penryn, UK

10

11        *correspondence: jesse.shapiro@umontreal.ca

12

13    **keywords**: microbiome, diversification, evolution, ecology, Earth Microbiome Project,

14    16S rRNA

15

16      **Abstract**

17

18      Microbes are embedded in complex microbiomes where they engage in a wide array of

19      inter- and intra-specific interactions[1–4]. However, whether these interactions are a

20      significant driver of natural biodiversity is not well understood. Two contrasting

21      hypotheses have been put forward to explain how species interactions could influence

22      diversification. 'Ecological Controls' (EC) predicts a negative diversity-diversification

23      relationship, where the evolution of novel types becomes constrained as available niches

24      become filled[5]. In contrast, 'Diversity Begets Diversity' (DBD) predicts a positive

25      relationship, with diversity promoting diversification via niche construction and other

26      species interactions[6]. Using the Earth Microbiome Project, the largest standardized

27      survey of global biodiversity to date[7], we provide support for DBD as the dominant

28      driver of microbiome diversity. Only in the most diverse microbiomes does DBD reach a

29      plateau, consistent with increasingly saturated niche space. Genera that are strongly

30      associated with a particular biome show a stronger DBD relationship than non-residents,

31      consistent with prolonged evolutionary interactions driving diversification. Genera with

32      larger genomes also experience a stronger DBD response, which could be due to a higher

33      potential for metabolic interactions and niche construction offered by more diverse gene

34      repertoires. Our results demonstrate that the rate at which microbiomes accumulate

35      diversity is crucially dependent on resident diversity. This fits a scenario in which species

36      interactions are important drivers of microbiome diversity. Further (population genomic

37      or metagenomic) data are needed to elucidate the nature of these biotic interactions in

38      order to more fully inform predictive models of biodiversity and ecosystem stability[4,5].

**Main text**

39

40   The majority of the genetic diversity on Earth is encoded by microbes[8–10] and the

41   functioning of all Earth's ecosystems is reliant on diverse microbial communities [11].

42   High-throughput 16S rRNA gene amplicon sequencing studies continue to yield

43   unprecedented insight into the taxonomic richness of microbiomes (e.g. [12,13]), and abiotic

44   drivers of community composition (e.g. pH[14,15]) are increasingly characterised. Although

45   it is known that biotic (microbe-microbe) interactions can also be important in

46   determining community composition[16], comparatively little is known about how such

47   interactions (e.g. cross-feeding[1] or toxin-mediated interference competition[2,3]) shape

48   microbiome diversity.

49        The dearth of studies exploring how microbial interactions could influence

50   diversification and diversity stands in marked contrast to a long research tradition on

51   biotic controls of plant and animal diversity[17,18]. In an early study of 49 animal

52   (vertebrate and invertebrate) community samples, Elton plotted the number of species

53   versus the number of genera and observed a ~1:1 ratio in each individual sample, but a

54   ~4:1 ratio when all samples were pooled[18]. He took this observation as evidence for

55   competitive exclusion preventing related species, more likely to overlap in niche space, to

56   co-exist. This concept, more recently referred to as niche filling or Ecological Controls

57   (EC)[5] predicts speciation (or, more generally, diversification) rates to decrease with

58   increasing standing species diversity because of diminished available niche space[19]. In

59   contrast, the Diversity Begets Diversity (DBD) model predicts that when species

60   interactions create novel niches, standing biodiversity favors further diversification[6,20].

61   For example, niche construction (i.e. the physical, chemical or biological alteration of the

**3**

62    environment) could influence the evolution of the species constructing the niche, and/or

63    that of co-occurring species[21,22].

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

65    communities has been mixed[20,23-26]. Laboratory evolution experiments have sought

66    general principles by tracking the diversification of a focal bacterial lineage in

67    communities of varying complexity – but the results have also been varied[27,28]. For

68    example, diversification of a focal *Pseudomonas* clone was favoured by increasing

69    community diversity in the range of 0-20 species within the same genus[20,29] but

70    diversification was inhibited by very diverse communities (*e.g.* hundreds or thousands of

71    species in natural soil[30]). These experimental results show how interspecific competition

72    can initially drive diversification[31], and eventually inhibit diversification as niches are

73    filled. However, these experiments were restricted to very short evolutionary time scales

74    (*i.e.* a few dozen mutations at most) in a small number of lineages, and it is unclear if

75    they can be generalized to natural communities evolving over longer periods, spanning

76    multiple speciation events and large-scale genomic changes.

77           To test whether natural microbial communities conform to EC or DBD models of

78    diversification, we used 2,000 microbiome samples from the Earth Microbiome Project

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

80    and sequencing protocols[7]. All samples were rarefied to 5,000 observations (counts of

81    16S rRNA gene sequences), as diversity estimates are highly sensitive to sampling

82    effort[32]. Instead of a phylogenetic approach requiring complex assumptions[33,34], we use

83    the equivalent of the Species:Genus (S:G) ratios that Elton used three quarters of a

84    century ago[18] to infer bacterial diversification rates. Rather than species, we considered

85  16S rRNA gene Amplicon Sequence Variants (ASVs) as our finest taxonomic unit. We

86  then used a range of taxonomic ratios (ASV:Genus, Genus:Family, Family:Order,

87  Order:Class, and Class:Phylum) as proxies for diversification of a focal lineage, from

88  shallow to deep evolutionary time, and plot these as a function of the number of non-focal

89  lineages (Genera, Families, Orders, Classes, and Phyla, respectively) with which the focal

90  lineage could interact. A negative relationship is consistent with the EC hypothesis,

91  whereas a positive relationship is consistent with the DBD hypothesis (**Fig. 1**). We used

92  generalized linear mixed models (GLMMs) to determine how the diversification of a

93  focal lineage (*e.g.* its ASV:Genus ratio) is affected by the diversity of other lineages (*e.g.*

94  non-focal genera) in the community. The effects of environment (as defined by the EMP

95  Ontology 'level 3 biomes;' Methods) and the identity of the focal lineage were included

96  by fitting these as random effects on the slope and intercept. We also controlled for the

97  submitting laboratory (identified by the principal investigator) and the EMP unique

98  sample identifier (i.e. if two taxa were part of the same sample). Finally, we repeated

99  these analyses using a taxonomy-free method based on nucleotide sequence identity

100  cutoffs (Methods).

101      The DBD model was supported across taxonomic ratios, which all had

102  significantly positive slopes fitting the diversity-diversification relationship (**Table S1,**

103  **Supplementary Data file 1 Section 1**), and the vast majority of slope estimates across

104  different lineages and environments were positive (**Fig. S1**). For example, the most

105  prevalent phylum across all samples, Proteobacteria, had significantly positive slopes

106  when fitted with linear models in all environments, except hypersaline and non-saline

107  sediments (**Fig. 2a**). For each taxonomic ratio, the three most prevalent taxa followed

108 positive slopes in most environments (**Fig. S2-S6**), with only a few instances of

109 significantly negative slopes (**Fig. 2b**). The predominance of positive slopes is robust and

110 remains after controlling for data structure and taxonomic assignment (**Fig. S7, S8;**

111 Supplementary Text), nor are they explained by widely measured abiotic drivers (*e.g.*

112 pH) that could simultaneously increase both diversity and diversification (**Table S2**;

113 **Supplementary Data file 1 Section 2**; Supplementary Text). Thus, the EMP data are

114 broadly consistent with the predictions of a DBD model.

115 The DBD hypothesis rests on the premise that species interactions drive

116 diversification[5,20]. We therefore expect that lineages that are more tightly associated with

117 a specific biome (i.e. long-term residents) are more likely to have had a long history of

118 interaction with community members and thus are more likely to experience DBD than

119 lineages that are not tightly associated with that biome (i.e. poorly adapted migrants or

120 broadly adapted generalists). To test this prediction, we clustered environmental samples

121 by their genus-level community composition using fuzzy *k*-means clustering (**Fig. 3a**),

122 which identified three clusters: 'animal-associated', 'saline', and 'non-saline'. The

123 clustering included some outliers (*e.g.* plant corpus grouping with animals), but were

124 generally intuitive and consistent with known distinctions between host-associated vs.

125 free-living[7], and saline vs. non-saline communities[35]. Resident genera were defined as

126 those with a strong preference for a particular environment cluster, using indicator

127 species analysis (permutation test, $P<0.05$; **Fig. 3a**; **Fig. S9**; **Supplementary Data file**

128 **2**), and genera without a strong preference were considered generalists. For each

129 environment cluster, we ran a GLMM with resident genus-level diversity (number of

130 non-focal genera) as a predictor of diversification (ASV:Genus ratio) for residents,

**6**

131    generalists, or migrants (residents of one cluster found in a different cluster)

132    (**Supplementary Data file 1 Section 3**). Resident diversity had no significant effect on

133    the diversification of generalists ($z$=0.646, $P$=0.518; $z$=0.279, $P$=0.780; $z$=0.347,

134    $P$=0.729, respectively for animal-associated, saline and non-saline clusters), but did

135    significantly increase resident diversification ($z$=7.1, $P$= 1.25e-12; $z$=3.316, $P$=0.0009;

136    $z$=7.109, $P$=1.17e-12, respectively). Resident diversity significantly decreased migrant

137    diversification in saline ($z$=-3.194, $P$=0.0014) and non-saline environment clusters ($z$=-

138    2.840, $P$=0.0045), but had no significant effect in the animal-associated cluster ($z$=-0.566,

139    $P$=0.571) (**Fig. 3b**). These results suggest that diversity begets diversification among

140    lineages sharing the same environment over a long evolutionary time period, but that this

141    is not the case for lineages that do not consistently occur in the same microbiome and

142    presumably interact less frequently. The diversification of migrants in a new environment

143    might even be impeded, presumably because most niches are already occupied by

144    residents.

145         The positive effect of diversity on diversification should eventually reach a

146    plateau as niches, including those constructed by biotic interactions, become

147    saturated[27,30]. In the animal distal gut, a relatively low-diversity biome, we observed a

148    strong linear DBD relationship at most sequence identity ratios; in contrast, the more

149    diverse soil biome clearly attained a plateau (**Fig. S10**). To further test the hypothesis that

150    increasingly diverse microbiomes experience weaker DBD due to saturated niche space,

151    we used a GLMM including the interaction between diversity and environment type as a

152    fixed effect. We considered this model only for taxonomic ratios with evidence for

153    significant DBD slope variation by environment (**Table S1**): Family:Order, Order:Class

154    and Class:Phylum. Consistent with our hypothesis, DBD slopes were significantly more

155    positive in less diverse (often host-associated) biomes (**Fig. 4a**, **Figure S11,**

156    **Supplementary Data file 1 Section 4**).

157         The Black Queen hypothesis posits that microbes embedded in complex

158    communities can exploit the production of extracellular public goods produced by other

159    species, resulting in selection for loss of genes encoding these goods – as long as the

160    essential trait is not lost from the community as a whole[36]. Lineages that interact more

161    frequently with other lineages through such public good exploitation would be expected

162    to experience greater loss of function and thus greater genome reduction. These reduced

163    genome would also be expected to experience stronger DBD, because their survival and

164    diversification is dependent on other community members. To test this expectation, we

165    assigned genome sizes to 576 genera for which at least one whole-genome sequence was

166    available and added an interaction term between genome size and diversity as a fixed

167    effect to the GLMM (Methods). Contrary to expectation, we observed a slight but

168    significant positive effect of genome size on the slope ($z=2.5$, $P=0.01$; **Fig. 4b,**

169    **Supplementary Data file 1 Section 5**). The positive relationship may even be stronger

170    than estimated, because genus-level genome size estimates are likely quite noisy. This

171    result supports a model in which biotic interactions (and resulting diversification) drive

172    genome expansion (*e.g.* through the accumulation of toxin- and resistance-gene diversity

173    during antagonistic coevolution[2]). Alternatively (or additionally), species with larger

174    biosynthetic gene repertoires and greater opportunity to engage in niche construction[21]

175    could be more prone to interact with other species, driving DBD.

176    Using 10 million individual marker sequences, we demonstrated a pervasive

177    positive relationship between prokaryotic diversity and diversification, which holds

178    across a broad range of environments and taxa. The strength of the DBD relationship

179    dissipates with increasing microbiome diversity which might be due to niche saturation,

180    or potentially due to the fact that highly diverse communities prevent species from

181    reliably interacting with each other. DBD appears to be particularly strong among deeply

182    diverged lineages (*e.g.* phyla), suggesting the importance of DBD in the ancient

183    diversification of bacterial lineages and supporting the view that high taxonomic ranks

184    are ecologically coherent[37,38]. We note that the very early stages of diversification are

185    inaccessible at the resolution of 16S ASVs, but this could be addressed in the future using

186    (meta-)genomic approaches. At the limited resolution of 16S sequences, we do not expect

187    measurable diversification within an individual microbiome sample; however community

188    diversity could still select for (as in DBD) or against (as in EC) standing diversity in a

189    focal lineages, even if this lineage diversified before the sampled community assembled.

190    Due to the correlational nature of our data, it is not possible to test whether the positive

191    relationship between diversification and diversity is primarily due to the creation of novel

192    niches via biotic interactions and niche construction[22], or potentially due to increased

193    competition leading to specialisation on underexploited resources[3,29]. Despite their

194    importance in shaping microbiome diversity and community structure, abiotic factors

195    such as pH and temperature do not appear to be driving the DB relationship; this could be

196    further tested in studies with more extensive abiotic metadata. Regardless of the

197    underlying mechanisms, our results demonstrate the importance of biotic interactions in

198    shaping microbiome diversity, which has important implications for modelling and

**9**

199  predicting their function and stability[4,39]. The answer to the question 'why are

200  microbiomes so diverse?' might in a large part be because microbiomes are so diverse[25].

201

202  **Acknowledgements**

203  We thank Luke Thompson for assistance obtaining EMP data and Zofia Ecaterina

204  Taranu, Vincent Fugère and Guillaume Larocque for advice on Generalized Linear

205  Mixed Models. We are also grateful to Steven Kembel and Tom Battin for critical

206  comments that improved the manuscript. **Funding:** This project was made possible by an

207  NSERC Discovery Grant and Canada Research Chair to BJS.

208

209  **Author contributions**

210  Conceptualization: BJS, MV. Data curation: NM. Formal analysis: NM, MV, BJS.

211  Funding acquisition: BJS. Investigation: NM, MV, PL, BJS. Methodology: NM, MV, PL,

212  BJS. Resources: BJS, PL. Supervision: PL, BJS. Software: NM. Visualization: NM.

213  Writing original draft: NM, MV, BJS. Writing - review & editing: NM, MV, PL, BJS.

214

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

216

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

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

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

220

221

222     **Supplementary Materials**

223

224     Supplementary text

225     Methods

226     Tables S1 – S2

227     Fig S1 – S11

228     File 1. Full GLMM outputs.

229     File 2. Indicator species analysis.

230

**Fig. 1. Contrasting the Diversity Begets Diversity (DBD) and Ecological Controls**

**(EC) models of diversification.** We consider the diversification of a focal lineage as a

function of initial diversity present at the time of diversification.

**(A)** For example, sample 1 contains one non-focal genus, and two ASVs diversify within

the focal genus (point at x=1, y=2 in the plot). Sample 2 contains three non-focal genera,

and four ASVs diversify within the focal genus (point at x=3, y=4). Tracing a line

through these points yields a positive slope, supporting the Diversity Begets

Diversification (DBD) model (red).

**(B)** Alternatively, a negative slope would support the Ecological Controls (EC) model

(blue line).

**Fig. 2. Diversification as a function of diversity across biomes in the phylum Proteobacteria.**

**(A) Linear models for diversification** (the number of classes within Proteobacteria, y-axis) as a function of diversity (the number of non-proteobacterial phyla, x-axis) in each of the 17 environments (EMPO3 biomes). P-values are Bonferroni corrected for 17 tests. Significant ($P < 0.05$) models are shown with red trend lines; non-significant ($P > 0.05$) trends are shown in blue.

**(B) Summary of linear model slopes across taxonomic ratios.** The number of significant positive (+) or negative (−) slope estimates are shown for each taxonomic ratio, summed across biomes. Significant slopes are those with $P < 0.05$ (Bonferroni corrected). Non-significant slope estimated are excluded.

**Fig. 3. Diversity begets diversification in resident versus non resident genera.**

**(A) PCA showing genera clustering into their preferred environment clusters**.

Circles indicate genera and triangles indicate environments (EMPO 3 biomes). The three

258     environment clusters identified by fuzzy *k*-means clustering are: Non-saline (NS, blue),

259     saline (S, green) and animal-associated (purple). Resident genera were identified by

260     indicator species analysis.

261     **(B) DBD in resident versus non resident genera across environment clusters.** Results

262     of GLMMs modeling diversification as a function of diversity in resident, migrant, or

263     generalist groups. The x-axis shows the standardized number of non-focal resident genera

264     (diversity); the y-axis shows the number of ASVs per focal genus (diversification).

265     Resident focal genera are shown in orange, migrant focal genera in red, and generalist

266     focal genera in black.

**Fig. 4. Ecological and evolutionary mechanisms to explain variation in the strength of DBD.**

**(A) DBD slope is higher in low-diversity (often host-associated) microbiomes.** The x-axis shows the mean number of phyla in each biome. On the y-axis, DBD slope was estimated by the GLMM predicting diversification as a function of the interaction between diversity and environment type at the Class:Phylum ratio (**Supplementary Data file 1 Section 4.3**). The line represents a regression line; the shaded area depicts 95% confidence limits of the fitted values.

**(B) Positive correlation between genome size and DBD slope.** Results are shown from a GLMM predicting diversification as a function of the interaction between diversity and genome size at the ASV:Genus ratio (**Supplementary Data file 1 Section 5**). The x-axis is genus-level genome size in Mbp (min=0.97, max=14.78); the y-axis is DBD slope (the effect of diversity on diversification). Vertical bars indicate 95% confidence limits of the fitted values.

**16**

283    **References**

284    1. Seth, E. C. & Taga, M. E. Nutrient cross-feeding in the microbial world. *Front.*

285    *Microbiol.* **5**, 350 (2014).

286    2. Czárán, T. L., Hoekstra, R. F. & Pagie, L. Chemical warfare between microbes

287    promotes biodiversity. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 786–790 (2002).

288    3. Hibbing, M. E., Fuqua, C., Parsek, M. R. & Peterson, S. B. Bacterial competition:

289    surviving and thriving in the microbial jungle. *Nat. Rev. Microbiol.* **8**, 15–25 (2010).

290    4. Coyte, K. Z., Schluter, J. & Foster, K. R. The ecology of the microbiome:

291    Networks, competition, and stability. *Science* **350**, 663–666 (2015).

292    5. Schluter, D. & Pennell, M. W. Speciation gradients and the distribution of

293    biodiversity. *Nature* **546**, 48–55 (2017).

294    6. Whittaker, R. H. Evolution and Measurement of Species Diversity. *Taxon* **21**,

295    213–251 (1972).

296    7. Thompson, L. R. *et al.* A communal catalogue reveals Earth's multiscale

297    microbial diversity. *Nature* **551**, 457-463 (2017).

298    8. Sunagawa, S. *et al.* Ocean plankton. Structure and function of the global ocean

299    microbiome. *Science* **348**, 1261359 (2015).

300    9. Lapierre, P. & Gogarten, J. P. Estimating the size of the bacterial pan-genome.

301    *Trends Genet.* **25**, 107–110 (2009).

302    10. Hug, L. A. *et al.* A new view of the tree and life's diversity. *Nature*

303    *Microbiology* **1**, 16048 (2016).

304    11. Falkowski, P. G., Fenchel, T. & Delong, E. F. The microbial engines that drive

305    Earth's biogeochemical cycles. *Science* **320**, 1034–1039 (2008).

306    12.  Sogin, M. L. *et al.* Microbial diversity in the deep sea and the underexplored

307    'rare biosphere'. *Proc. Natl. Acad. Sci. U. S. A.* **103**, 12115–12120 (2006).

308    13.    Louca, S., Mazel, F., Doebeli, M. & Parfrey, L. W. A census-based

309    estimate of Earth's bacterial and archaeal diversity. *PLoS Biol.* **17**, e3000106 (2019).

310    14.    Lauber, C. L., Hamady, M., Knight, R. & Fierer, N. Soil pH as a predictor

311    of soil bacterial community structure at the continental scale: a pyrosequencing-

312    based assessment. *Appl. Environ. Microbiol.* **75**, 5111-5120 (2009).

313    15.    Power, J. F. *et al.* Microbial biogeography of 925 geothermal springs in

314    New Zealand. *Nat. Commun.* **9**, 2876 (2018).

315    16.    Needham, D. M. & Fuhrman, J. A. Pronounced daily succession of

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

317    *Microbiology* **1**, 16005 (2016).

318    17.    Gause, G. F. *The Struggle for Existence*. (Courier Corporation, 2003).

319    18.    Elton, C. Competition and the Structure of Ecological Communities. *J.*

320    *Anim. Ecol.* **15**, 54–68 (1946).

321    19.    Rabosky, D. L. & Hurlbert, A. H. Species richness at continental scales is

322    dominated by ecological limits. *Am. Nat.* **185**, 572–583 (2015).

323    20.    Calcagno, V., Jarne, P., Loreau, M., Mouquet, N. & David, P. Diversity

324    spurs diversification in ecological communities. *Nat. Commun.* **8**, 15810 (2017).

325    21.    San Roman, M. & Wagner, A. An enormous potential for niche

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

327    *Comput. Biol.* **14**, e1006340 (2018).

328    22.    Laland, K. N., Odling-Smee, F. J. & Feldman, M. W. Evolutionary

329    consequences of niche construction and their implications for ecology. *Proc. Natl.*

330    *Acad. Sci. U. S. A.* **96**, 10242–10247 (1999).

331    23.    Price, T. D. *et al.* Niche filling slows the diversification of Himalayan

332    songbirds. *Nature* **509**, 222-225 (2014).

333    24.    Rabosky, D. L. *et al.* An inverse latitudinal gradient in speciation rate for

334    marine fishes. *Nature* **559**, 392–395 (2018).

335    25.    Emerson, B. C. & Kolm, N. Species diversity can drive speciation. *Nature*

336    **434**, 1015–1017 (2005).

337    26.    Palmer, M. W. & Maurer, T. A. Does Diversity Beget Diversity? A Case

338    Study of Crops and Weeds. *J. Veg. Sci.* **8**, 235–240 (1997).

339    27.    Brockhurst, M. A., Colegrave, N., Hodgson, D. J. & Buckling, A. Niche

340    occupation limits adaptive radiation in experimental microcosms. *PLoS One* **2**, e193

341    (2007).

342    28.    Meyer, J. R. & Kassen, R. The effects of competition and predation on

343    diversification in a model adaptive radiation. *Nature* **446**, 432–435 (2007).

344    29.    Jousset, A., Eisenhauer, N., Merker, M., Mouquet, N. & Scheu, S. High

345    functional diversity stimulates diversification in experimental microbial

346    communities. *Sci Adv* **2**, e1600124 (2016).

347    30.    Gómez, P. & Buckling, A. Real-time microbial adaptive diversification in

348    soil. *Ecol. Lett.* **16**, 650–655 (2013).

349    31.    Bailey, S. F., Dettman, J. R., Rainey, P. B. & Kassen, R. Competition both

350    drives and impedes diversification in a model adaptive radiation. *Proc. Biol. Sci.*

351    **280**, 20131253 (2013).

352    32.    Gotelli, N. J. & Colwell, R. K. Quantifying biodiversity: procedures and

353    pitfalls in the measurement and comparison of species richness. *Ecol. Lett.* **4**, 379–

354    391 (2001).

355    33.    Etienne, R. S., Pigot, A. L. & Phillimore, A. B. How reliably can we infer

356    diversity-dependent diversification from phylogenies? *Methods Ecol. Evol.* **7**, 1092–

357    1099 (2016).

358    34.    Louca, S. *et al.* Bacterial diversification through geological time. *Nat Ecol*

359    *Evol* **2**, 1458–1467 (2018).

360    35.    Lozupone, C. A. & Knight, R. Global patterns in bacterial diversity. *Proc.*

361    *Natl. Acad. Sci. U. S. A.* **104**, 11436–11440 (2007).

362    36.    Morris, J. J. & Lenski, R. E. The Black Queen Hypothesis: evolution of

363    dependencies through adaptive gene loss. *MBio* **3**, e00036-12 (2012).

364    37.    Philippot, L. *et al.* The ecological coherence of high bacterial taxonomic

365    ranks. *Nat. Rev. Microbiol.* **8**, 523–529 (2010).

366    38.    Martiny, J. B. H., Jones, S. E., Lennon, J. T. & Martiny, A. C.

367    Microbiomes in light of traits: A phylogenetic perspective. *Science* **350**, aac9323

368    (2015).

369    39.    Pennekamp, F. *et al.* Biodiversity increases and decreases ecosystem

370    stability. *Nature* **563**, 109–112 (2018).

# Supplementary Materials

# Diversity begets diversity in microbiomes

**Authors**: Naïma Madi[1], Michiel Vos[2], Pierre Legendre[1] and B. Jesse Shapiro[1*]

1. Departement de sciences biologiques, Universite de Montreal, Canada

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

    *correspondence: jesse.shapiro@umontreal.ca

**Supplementary Text**

**Methods**

**Tables S1 – S2**

**Figures S1 – S11**

## Supplementary Text

384

385         To test for any potential confounding effects of data structure or sampling bias,

386 we sought to remove any patterns of co-occurrence between ASVs in the same sample

387 via permutation. We took 2,000 simulated samples by selecting from the overall

388 distribution of 155,002 unique ASVs across all samples, weighted by their abundance

389 (total number of sequence counts). This resulted in a slightly negative diversity-

390 diversification relationship (slope = −0.002; Pearson correlation = −0.61; P<2.2.e−16;

391 **Fig. S7**), indicating that the observed positive relationships (**Table S1; Fig. 2**) are not the

392 effect of data structure.

393         We sought to further validate the results with a taxonomy-independent approach,

394 because not all taxonomic ranks have the same phylogenetic depth [40] and not all named

395 taxa are monophyletic [41]. Therefore, we clustered ASVs at decreasing levels of nucleotide

396 identity, from 100% identical ASVs down to 75% identity (roughly equivalent to phyla

397 [42]). We estimated diversification as the mean number of descendants per cluster (e.g.

398 number of 100% clusters per 97% cluster) and plotted this against the total number of

399 non focal clusters (97% identity in this example). For each of the six nucleotide

400 divergence ratios tested, the relationship between diversity and diversification was

401 positive (**Fig. S8**), consistent with DBD and suggesting that the taxonomic analyses were

402 largely unbiased.

403         To exclude the possibility that our results were driven by abiotic confounders, we

404 repeated the taxonomic analysis on a subset of 192 EMP samples for which

405 measurements of four important abiotic drivers of diversity, temperature, pH, latitude,

406 and elevation [5,14,15,43] were available. We fitted a GLMM with diversification rate as the

**22**

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

408    and their interactions as predictors (fixed effects). As in the full dataset (**Table S1**),

409    diversification was positively associated with diversity at all taxonomic ratios (**Table S2**).

410    As expected, certain abiotic factors, alone or in combination with diversity, had

411    significant effects on diversification. However, the effects of abiotic factors were always

412    weaker than the effect of community diversity (**Table S2**; **Supplementary Data file 1**

413    **Section 2**). Although only a small subset of abiotic factors was considered, this analysis

414    suggests that the DBD trend is unlikely to be mainly driven by variation in the abiotic

415    environment.

# Methods

416

417

418 **16S rRNA marker data acquisition and preprocessing.**

419 16S rRNA-V4 region reads (90 bp, GreenGenes 13.8 taxonomy) along with

420 environmental data and EMPO3 designations

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

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

423 summaries were downloaded from :

424 ftp://ftp.microbio.me/emp/release1/otu_distributions/otu_summary.emp_deblur_90bp.sub

425 set_2k.rare_5000.tsv, environmental data from :

426 ftp://ftp.microbio.me/emp/release1/mapping_files/emp_qiime_mapping_release1.tsv, and

427 EMPO3 designations from :

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

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

430 were downloaded from https://www.nature.com/articles/nature24621#s1.

431 We used the EMP '2000 subset' rarefied to 5000 sequences per sample. This subset

432 contains 155 002 ASVs from 2000 samples with an even distribution across 17 natural

433 environments (EMP Ontology level 3) (Thompson et al,. 2017). Based on the ASVs

434 annotations across samples, we estimated diversification for every taxonomic ratio

435 (ASV:Genus, Genus:Family, Family:Order, Order:Class and Class:Phylum), along with

436 the number of non-focal lineages (Python script, Python Version 2.7).

437

438

**24**

**Generalized Linear Mixed Models (GLMMs)**

All models were fitted in Rstudio (Version 1.1.442, R Version 3.5.2) using the glmer function of the lme4 package [44]. Data standardization (transformation to a mean of zero and a standard deviation of one) was applied to all predictors to get comparable estimates. In models with only one predictor, applying standardization resolved convergence warnings and considerably sped up the optimization. Standardization has previously been reported to improve model performance and solve convergence problems[45].

We used likelihood-ratio tests (anova R function from stats package) as follows: 1) on nested models to assess the significance of random effects (in the nested models, each effect was dropped one at a time); 2) on the full model and the null model comprising only random effects, to assess the significance of fixed effects[46]; 3) on the full model and the model without the interaction term, to assess the significance of interactions. All models reported here were found to be significant ($P<0.05$).

Diagnostic plots (plot and qqnorm R functions in base and stats packages) were checked for each model to ensure that residual homoscedasticity (homogeneity of variance) was fulfilled: no increase of the variance with fitted values and residuals were symmetrically distributed tending to cluster around the 0 of the ordinate, but with an expected pattern due to count data. Normality plots were imperfect, but they generally showed that the residuals were close to being normally distributed. The assumption of normality is often difficult to fulfill with high numbers of observations, as is the case in our models (https://www.statisticshowto.datasciencecentral.com/shapiro-wilk-test/), and

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

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

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

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

465    were not overdispersed, but rather were underdispersed. The ratio of the sum of squared

466    Pearson residuals to residual degrees of freedom was < 1 and non-significant when tested

467    with a chi-squared test. Given that underdispersion leads to more conservative results, we

468    retained the GLMMs with Poisson error distribution, despite the underdispersion.

469    (GLMM FAQ; Ben Bolker and others; 25 September 2018;

470    https://bbolker.github.io/mixedmodels-misc/glmmFAQ.html#underdispersion).

471

472    **Taxonomy-based generalized linear mixed models**

473    The effect of diversity on diversification was tested for different environment types and

474    lineages using generalized linear mixed models (GLMMs) fitted on the EMP dataset, for

475    all taxonomic ratios. As the dependent variable (diversification, defined as taxonomic

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

477    count response, we used a Poisson error distribution with a log link function. Diversity

478    (number of non-focal lineages: non-focal Genera, Families, Orders, Classes, and Phyla),

479    standardized to a mean of zero and a standard deviation of one, was specified as the

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

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

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

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

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

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

486    variation across groups of each categorical variable. We included the EMP unique sample

487    ID as a random effect to control for dependencies between observations (if two taxa were

488    part of the same sample).

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

490    diversification across different taxonomic ratios, we used a separate GLMM, with

491    Poisson error distribution with a log link function, for every ratio. We fitted the GLMM

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

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

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

495    variables (temperature, pH, latitude and elevation) were available. We defined diversity

496    and the abiotic variables as well as the interactions between diversity and every abiotic

497    variable as predictors (fixed effects) of diversification. All predictors were standardized

498    to a mean of zero and a standard deviation of one to obtain comparable estimates. The

499    GLMM had the same random effects as in the previous analysis, but only on the intercept

500    for simplicity.

501

502    **Nucleotide sequence identity-based analysis**

503    We defined a threshold of percent nucleotide identity between ASVs, corresponding to

504    different taxonomic ranks (from 100% identical ASVs down to 75% identity) [42]. Fasta

505    files for all samples were produced by a python script (Python Version 2.7) from the

506    sequences summary file (otu_summary.emp_deblur_90bp.subset_2k.rare_5000 from

507     EMP ftp server). We clustered sequences from each sample using USEARCH V9.2. We

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

509     diversification as the mean number of descendent clusters (*e.g.* number of 100% clusters

510     per 97% cluster). To describe the relationship between diversity and diversification, we

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

512     were based on the adjusted $R^2$.

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

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

515     effects of this non-independence on the results, we conducted permutation tests by

516     randomizing the associations between $d_i$ and $d_{i+1}$. Using 999 permutations, *P*-values were

517     calculated based on how many times we observed a correlation greater than that seen in

518     the real data (cor.test R function with kendall method). In each permutation, we

519     recalculated the significance test (Wald z) for the correlation in the randomized data, and

520     then computed the P-value based on how many times we observed a z value greater than

521     that of the original data (one tailed test because we wanted to demonstrate that the

522     relationship was positive). At all six levels of nucleotide identity, the real data always

523     showed a significantly stronger positive correlation when compared to permuted data (*P*

524     = 0.001), indicating that the DBD patterns was not an artefact of the dependence structure

525     in the data.

526         The effect of diversity on diversification was also tested across different

527     environments analysed separately. We modelled this relationship with linear, quadratic

528     and cubic fits, and compared those models based on the adjusted $R^2$.

529

## DBD among residents of the same environment

We clustered the environmental samples based on their genus-level community composition using fuzzy *k*-means clustering. Fuzzy clustering is a version of non-hierarchical clustering, where each cluster is a fuzzy set of all biomes and greater membership values indicates higher confidence in the allocation pattern to the cluster. The clustering (cmeans function, package e1071 in R) was done on the 'hellinger' transformed data (decostand function, package vegan in R). To identify resident genera to each cluster, we used indicator species analysis [47] as implemented in the indval function (labdsv R package). Indicators are genera found mostly in a certain environment group and present in the majority of environments of that group. The indicator value (indval index) of a genus is (maximum=1) if the genus is observed in only one environmental cluster and in all samples belonging to that cluster. We defined residents as genera with indval indices between 0.4 and 0.9, with permutation test $P < 0.05$. Genera not been associated with any cluster were considered generalists. We used principal component analysis (PCA) to visualize clustering and indicator genera (rda function, vegan R package). We then ran a separate GLMM for each environmental cluster, with resident genus-level diversity (number of non-focal genera) as a predictor of diversification (ASV:Genus ratio) for resident, migrant (residents of one cluster found in a different cluster) and generalist genera. The fixed effect was specified as the interaction between diversity and a factor defining the genus-cluster association (with three levels: resident, migrant and generalist). Random effects on intercept and slope were kept as in the previous GLMMs.

**DBD variation across biomes**

We tested the variation of DBD slope across different environments by defining environment (EMPO 3 biome type) as fixed effect. We fitted a GLMM with the interaction between diversity and environment type as a predictor of diversification. The main effects of diversity and environment individually were not included for model simplicity and we sought to look at the effect of the interaction alone (diversity*environment). All other random effects on intercept and slope were kept as in the previous GLMMs. DBD variation across environments was tested for Family:Order, Order:Class and Class:Phylum taxonomic ratios, as DBD slope variation by environment was statistically significant (likelihood-ratio test) for these ratios (**Table S1**).

**Genome size analysis**

We chose a subset of genera represented by one or more sequenced genomes in the NCBI microbial genomes database (https://www.ncbi.nlm.nih.gov/genome/browse#!/prokaryotes/). For these genera, a representative genome size was assigned by selecting the genome with the lowest number of scaffolds (if no closed genomes were available). If multiple genomes were available, sequenced to the same level of completion, the largest genome size was used. We fitted a GLMM on the subset of data with known genome size (576 genera) with  the interaction between diversity and genome size as a predictor of diversification (ASV:Genus). All the other random effects on intercept and slope were kept as in the previous GLMMs.

576 **Code availability**

577 All computer code used for analysis are archived on the github repository

578 https://github.com/Naima16/dbd.git.

579    **Supplementary references**

580    40.        Vos, M. A species concept for bacteria based on adaptive divergence.

581    *Trends Microbiol.* **19**, 1–7 (2011).

582    41.        Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome

583    phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004

584    (2018).

585    42.        Konstantinidis, K. T. & Tiedje, J. M. Towards a genome-based taxonomy

586    for prokaryotes. *J. Bacteriol.* **187**, 6258–6264 (2005).

587    43.        Delgado-Baquerizo, M. *et al.* A global atlas of the dominant bacteria

588    found in soil. *Science* **359**, 320–325 (2018).

589    44.        Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-

590    Effects Models Using lme4. *Journal of Statistical Software, Articles* **67**, 1–48

591    (2015).

592    45.        Harrison, X. A. *et al.* A brief introduction to mixed effects modelling and

593    multi-model inference in ecology. *PeerJ* **6**, e4794 (2018).

594    46.        Forstmeier, W. & Schielzeth, H. Cryptic multiple hypotheses testing in

595    linear models: overestimated effect sizes and the winner's curse. *Behav. Ecol.*

596    *Sociobiol.* **65**, 47–55 (2011).

597    47.        Dufrene, M. & Legendre, P. Species Assemblages and Indicator Species:

598    The Need for a Flexible Asymmetrical Approach. *Ecol. Monogr.* **67**, 345–366

599    (1997).

600

601 **Supplementary Tables**

602

603 **Table S1. Diversity has a positive effect on diversification across taxonomic ratios.**

604 The GLMMs showed statistically significant positive effect of diversity on

605 diversification. Each row reports the effect of diversity on diversification, as well as its

606 standard deviation, Wald z-statistic for its effect size and the corresponding $P$-value (left

607 section), or standard deviation on the slope for the significant random effects (right

608 section). SE=standard error, Env=environment type, Lin=lineage type, Lab=Principal

609 Investigator ID, Sample=EMP Sample ID. Interactions are denoted as '*'. n.s.=not

610 significant (likelihood-ratio test).

| | Slope (fixed effects) | | | | Standard deviation on the slope (random effects) | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Diversity | 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. |

611
612

613 **Table S2. Diversity has a stronger effect than abiotic factors on diversification.**

614 Results are shown from GLMMs with diversity, four abiotic factors (temperature,

615 elevation, pH, and latitude), and their interactions with diversity, as predictors of

616 diversification. Random effects on the intercept included environment, lineage, lab ID

617 and sample ID. Results are summarized as the coefficient (slope)±standard error (for

618 fixed effects). Temp=temperature, Lat=latitude, Elev=elevation. Interactions denoted as

619 '*'. Significant terms (Wald test) are shown in bold: ***$P$<2.2e-16; **$P$<0.01, *$P$ <0.05.

620 Random effects are not shown.

621
622

|  | Diversity | Temp | Lat | pH | Elev | Div *Temp | Div *Lat | Div *pH | Div *Elev |
|---|---|---|---|---|---|---|---|---|---|
| **ASV: Genus** | **0.129*** ± 0.013** | **0.044** ±0.016** | 0.017 ±0.019 | 0 ±0.018 | 0 ±0.023 | **0.043** ±0.014** | **0.032* ±0.014** | 0.003 ±0.011 | **-0.032* ±0.016** |
| **Genus: Family** | **0.094*** ±0.009** | **0.04*** ±0.011** | -0.009 ±0.01 | **- 0.049** ±0.009** | - 0.003±0.01 | 0.019 ±0.01 | -0.011 ±0.009 | -0.011 ±0.007 | -0.005 ±0.009 |
| **Family: Order** | **0.12*** ±0.013** | 0.012 ±0.014 | 0.002 ±0.021 | 0 ±0.013 | - 0.011±0.026 | 0.024 ±0.013 | 0.01 ±0.013 | 0.003 ±0.009 | -0.015 ±0.014 |
| **Order: Class** | **0.184*** ±0.01** | 0.001 ±0.013 | -0.011 ±0.012 | -0.002 ±0.012 | - 0.008±0.013 | **0.036** ±0.012** | **0.023* ±0.01** | -0.003 ±0.01 | **-0.02 ±0.01*** |
| **Class: Phylum** | **0.233*** ±0.013** | -0.025 ±0.015 | 0.014 ±0.015 | 0.011 ±0.015 | 0.032 ±0.019 | **0.06*** ±0.015** | **0.039** ±0.013** | **0.029* ±0.013** | 0.004 ±0.016 |

623

## Supplementary Figures

**Figure S1. Distributions of DBD slope estimates across different random effects, from the GLMMs predicting diversification as a function of diversity. (A) Class:Phylum, (B) Order:Class, (C) Family:Order, (D) Genus:Family and (E) ASV:Genus ratios.** Estimation of random effect coefficients from the GLMMs (Table S1), shows that the effect of diversity on diversification (slope estimates) are generally positive but could be negative in some lineages or combinations of environment, lineage (Environment*Lineage), and the laboratory that submitted the dataset (Environment*Lab).

637 **Figure S2. Diversification as a function of diversity across biomes in the two most**
638 **prevalent phyla after Proteobacteria (shown in Figure 2A of the main text).** (A)
639 Bacteroidetes, (B) Actinobacteria. Linear models are shown for diversification (classes
640 number per phylum, y-axis) as a function of diversity (non focal phyla number, x-axis) in
641 each of the 17 environments (EMPO3 biomes). P-values are Bonferroni corrected for 17
642 tests. Significant ($P <0.05$) models are shown with red trend lines, non-significant ($P >$
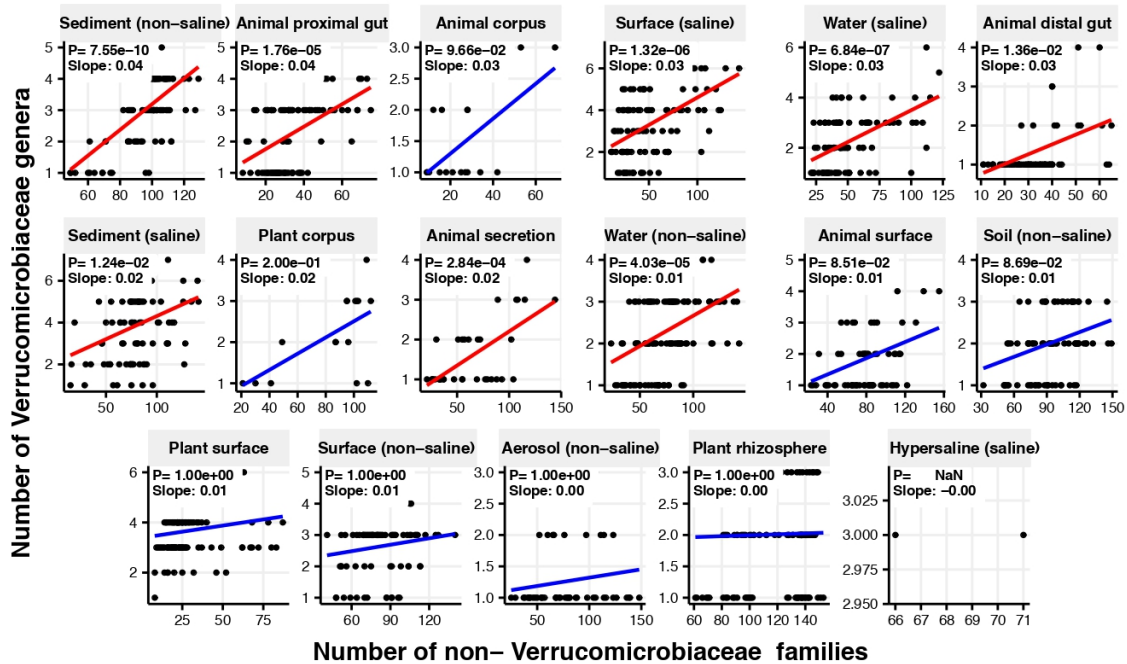643 0.05) trends are shown in blue.



## A. Bacteroidetes

644
645
646

## B. Actinobacteria



647
648
649
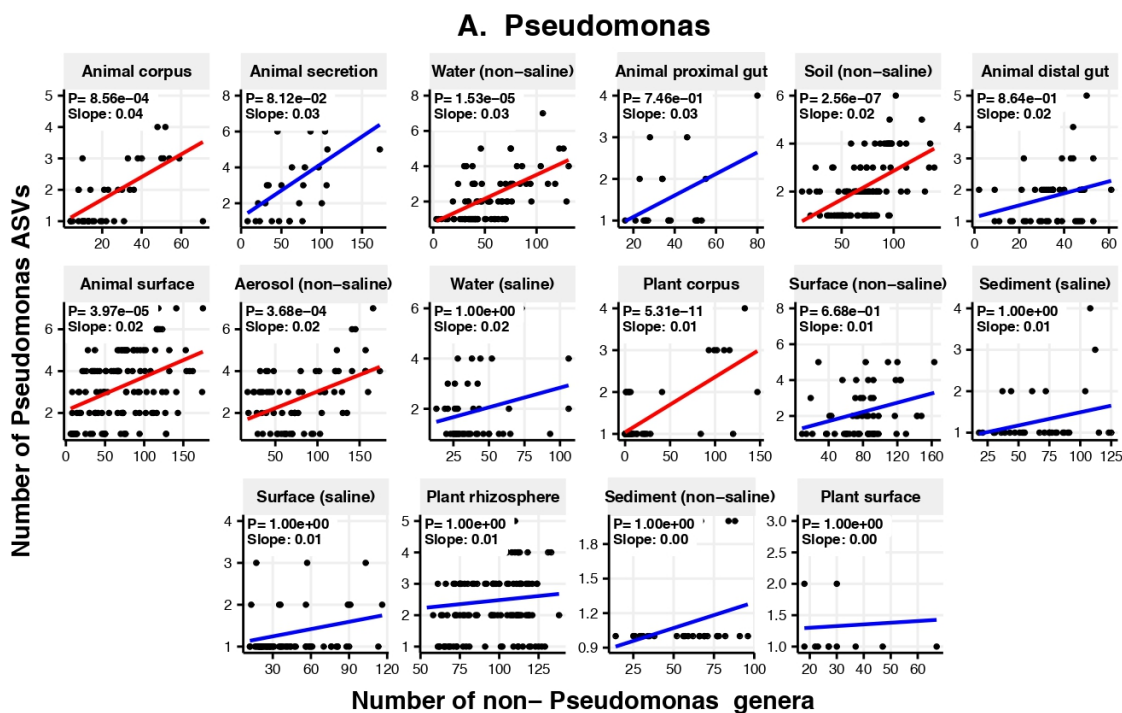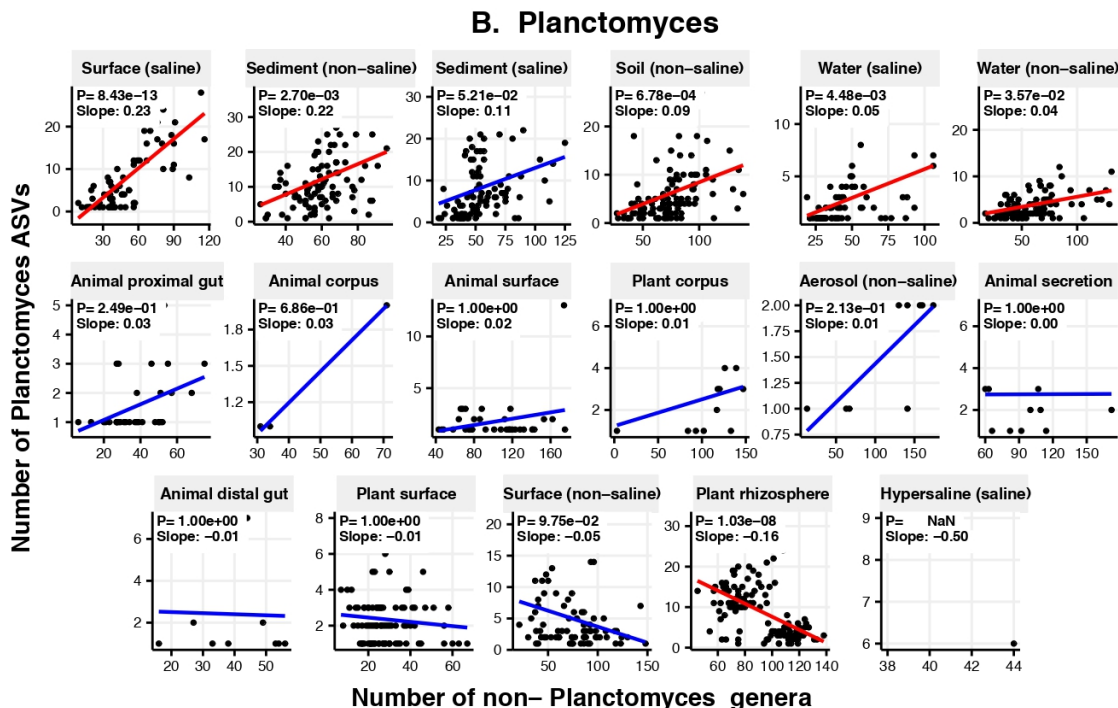650
651

**Figure S3. Diversification as a function of diversity across biomes in the three most prevalent classes.** Linear models are shown for diversification (orders per class, y-axis) as a function of diversity (non-focal classes, x-axis) in each of the 17 environments (EMPO3 biomes). P-values are Bonferroni corrected for 17 tests. Significant ($P <0.05$) models are shown with red trend lines, non-significant ($P > 0.05$) trends are shown in blue.



### A. Gammaproteobacteria

## B. Alphaproteobacteria



659
660
661

## C. Actinobacteria



662
663
664
665

666 **Figure S4. Diversification as a function of diversity across biomes in the three most**
667 **prevalent orders.** Linear models are shown for diversification (families per order, y-
668 axis) as a function of diversity (non-focal orders, x-axis) in each of the 17 environments
669 (EMPO3 biomes). P-values are Bonferroni corrected for 17 tests. Significant ($P < 0.05$)
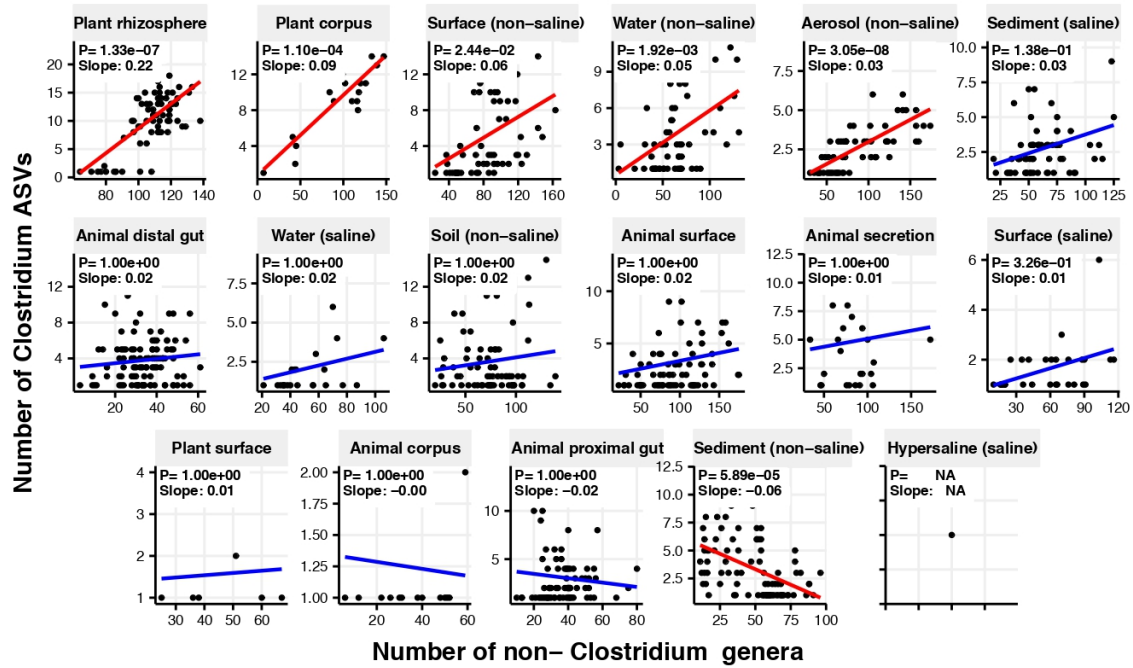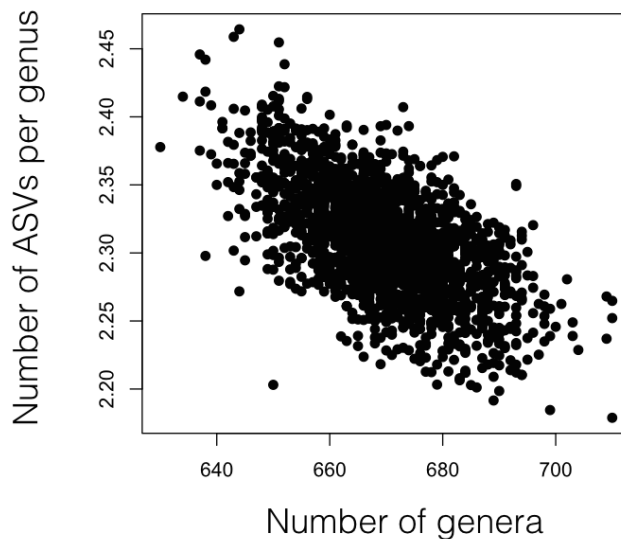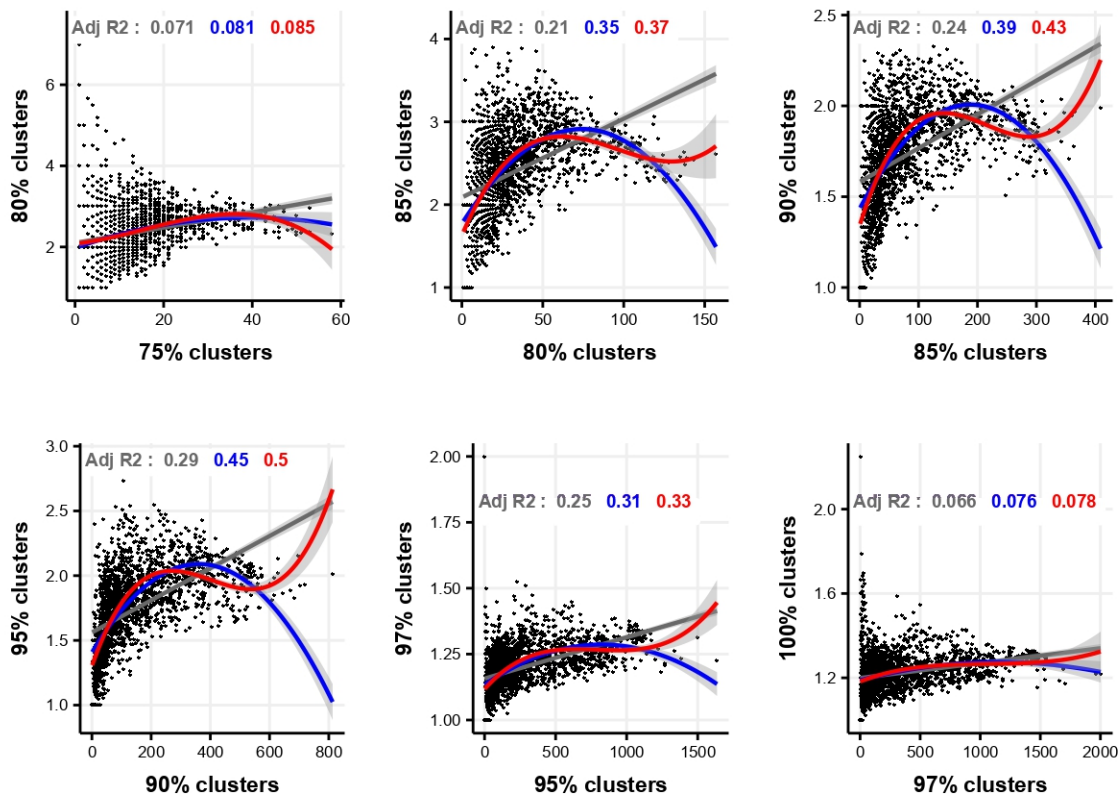670 models are shown with red trend lines, non-significant ($P > 0.05$) trends are shown in
671 blue.



672

## B. Flavobacteriales



673

## C. Rhizobiales



674
675
676
677
678
679
680

681 **Figure S5. Diversification as a function of diversity across biomes in the three most**
682 **prevalent families.** Linear models are shown for diversification (genera per family, y-
683 axis) as a function of diversity (non-focal families, x-axis) in each of the 17 environments
684 (EMPO3 biomes). P-values are Bonferroni corrected. Significant ($P < 0.05$) models are
685 shown with red trend lines, non-significant ($P > 0.05$) trends are shown in blue.



686
687

## B.  Sphingomonadaceae



688

## C.  Verrucomicrobiaceae



689
690
691
692
693

**Figure S6. Diversification as a function of diversity across biomes in the three most prevalent genera.** Linear models are shown for diversification (ASVs per genus, y-axis) as a function of diversity (non-focal genera, x-axis) in each of the 17 environments (EMPO3 biomes). P-values are Bonferroni corrected. Significant ($P < 0.05$) models are shown with red trend lines, non-significant ($P > 0.05$) trends are shown in blue.



**A. Pseudomonas**



**B. Planctomyces**

## C. Clostridium



701
702
703
704
705
706

707 **Figure S7. Permuted EMP data is biased toward a negative diversity-diversification**
708 **relationship.** We permuted the EMP dataset of 2,000 samples each rarefied to 5,000
709 sequences/sample and took 2,000 simulated samples, by picking from the overall
710 distribution of 155,002 unique ASVs across all samples, weighted by their total number
711 of observations. Thus, the 'true' patterns of co-occurrence between ASVs in the same
712 sample (and thus any 'biologically true' pattern of either DBD or EC models) is removed
713 from the data. The permutations yield a negative relationship between diversity (number
714 of genera) and diversification (number of ASVs per genus): slope = -0.002; Pearson
715 correlation = -0.61; $P<2.2.e16$.
716



717
718

719 **Figure S8. Linear, quadratic and cubic models for the relationship between**
720 **diversification and diversity for varying levels of % nucleotide identity.** Diversity
721 was estimated as the number of clusters at a focal level ($d_i$) and diversification as the
722 mean of the clusters at the rank above ($d_{i+1}/d_i$). All *P*-values are < 0.001. Linear fit
723 (grey); quadratic fit (blue), cubic fit (red); same colors for the associated adjusted $R^2$. The
724 x-axis (diversity) shows the number of clusters at the focal percent-identity level ($d_i$), and
725 the y-axis (diversification) is the mean of the clusters at the rank above ($d_{i+1}/d_i$).



726
727
728
729
730
731

47

732 **Figure S9. Resident genera of environment clusters.** Results from indicator species
733 analysis illustrated as a heatmap**.** Only the 25 resident genera with the highest indval
734 indices and *P*<0.05 (permutation test) are shown for every environment cluster (animal-
735 associated, non-saline and saline free). For the full results see **Supplementary Data file**
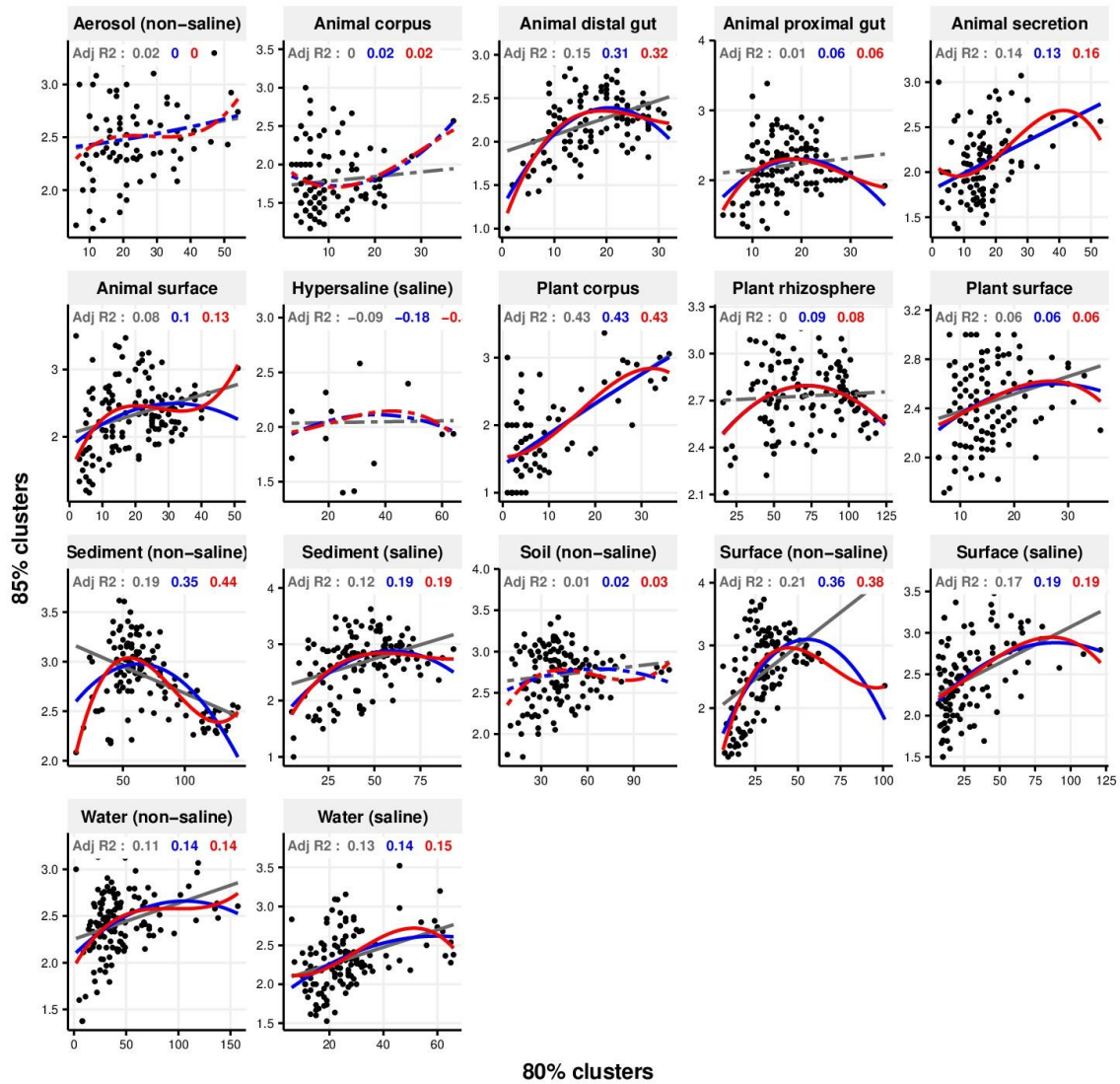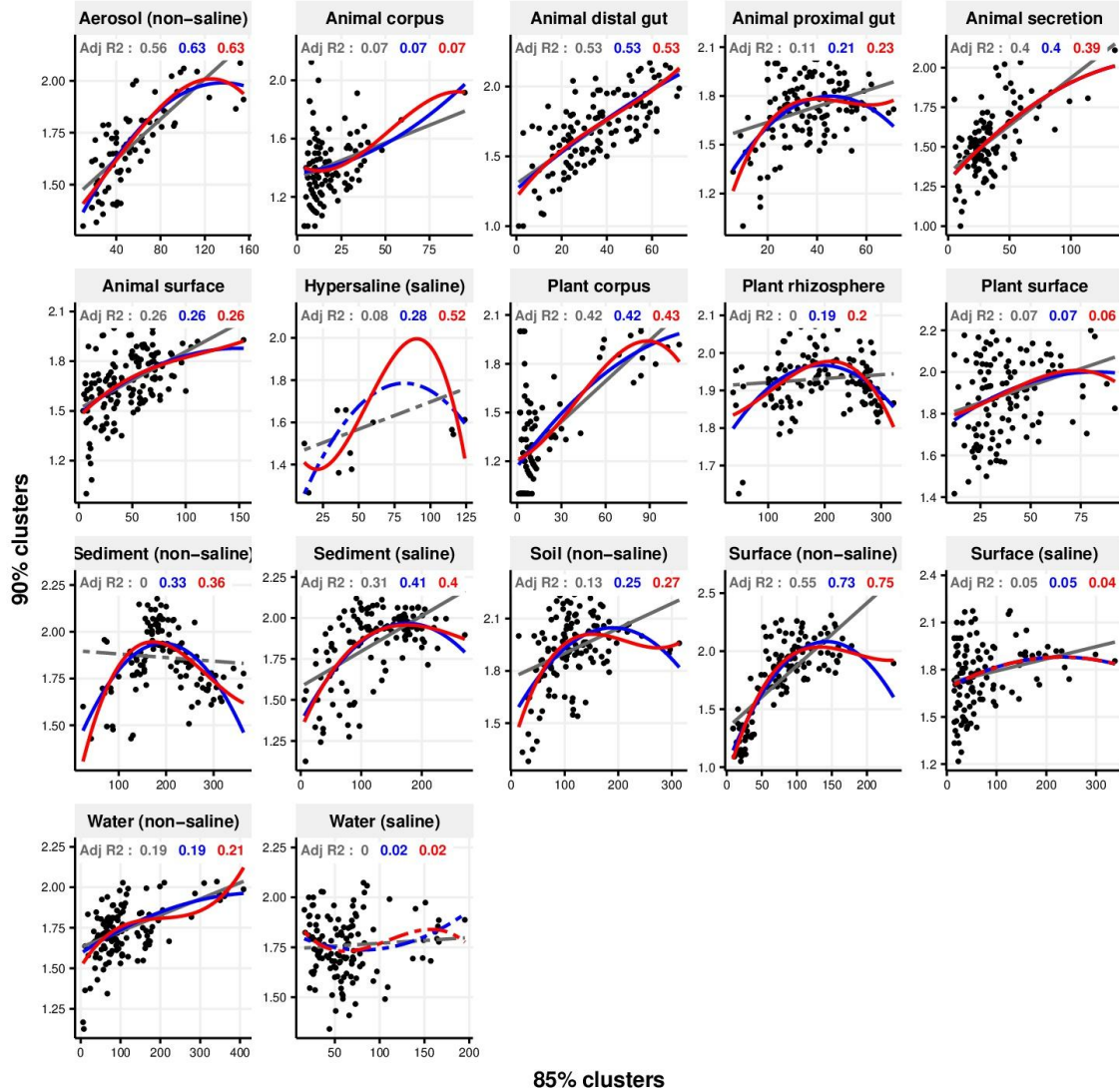736 **2.**
737



738
739
740

741 **Figure S10. Linear, quadratic and cubic models for diversification-diversity**
742 **relationship for each environment type for varying levels of % nucleotide identity.**
743 Diversity was estimated as the number of clusters at a focal level ($d_i$) and diversification
744 as the mean of the clusters at the rank above ($d_{i+1}/d_i$). Linear (grey), quadratic (blue) and
745 cubic (red), with corresponding adjusted R-squared values in the same color. *P*-values are
746 Bonferroni corrected for 17 tests. Significant, $P < 0.05$ (solid lines), non-significant
747 (dashed lines). The x-axis shows the number of clusters at the focal percent-identity level
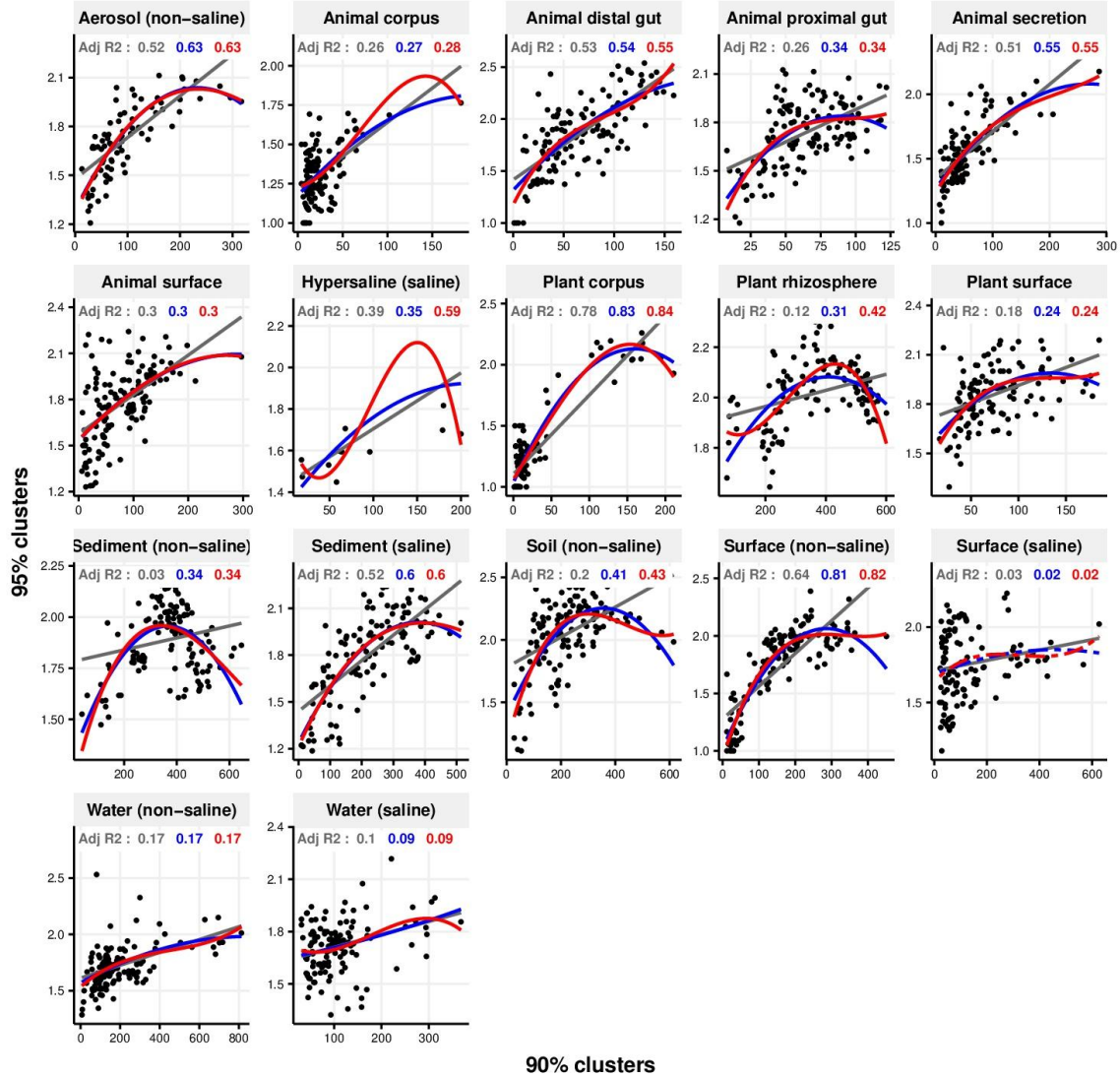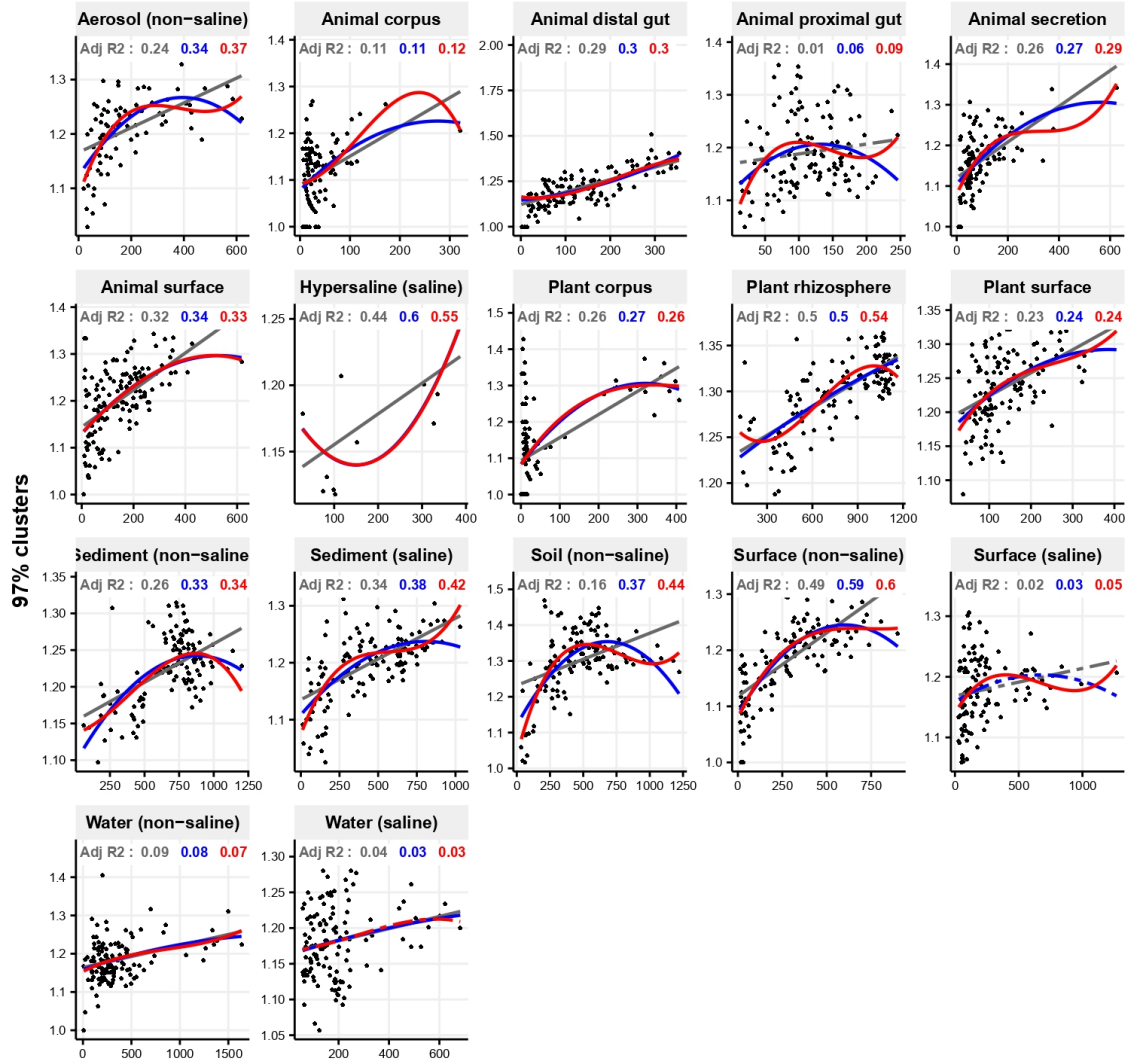748 ($d_i$), and the y-axis is the mean of the clusters at the rank above ($d_{i+1}/d_i$).
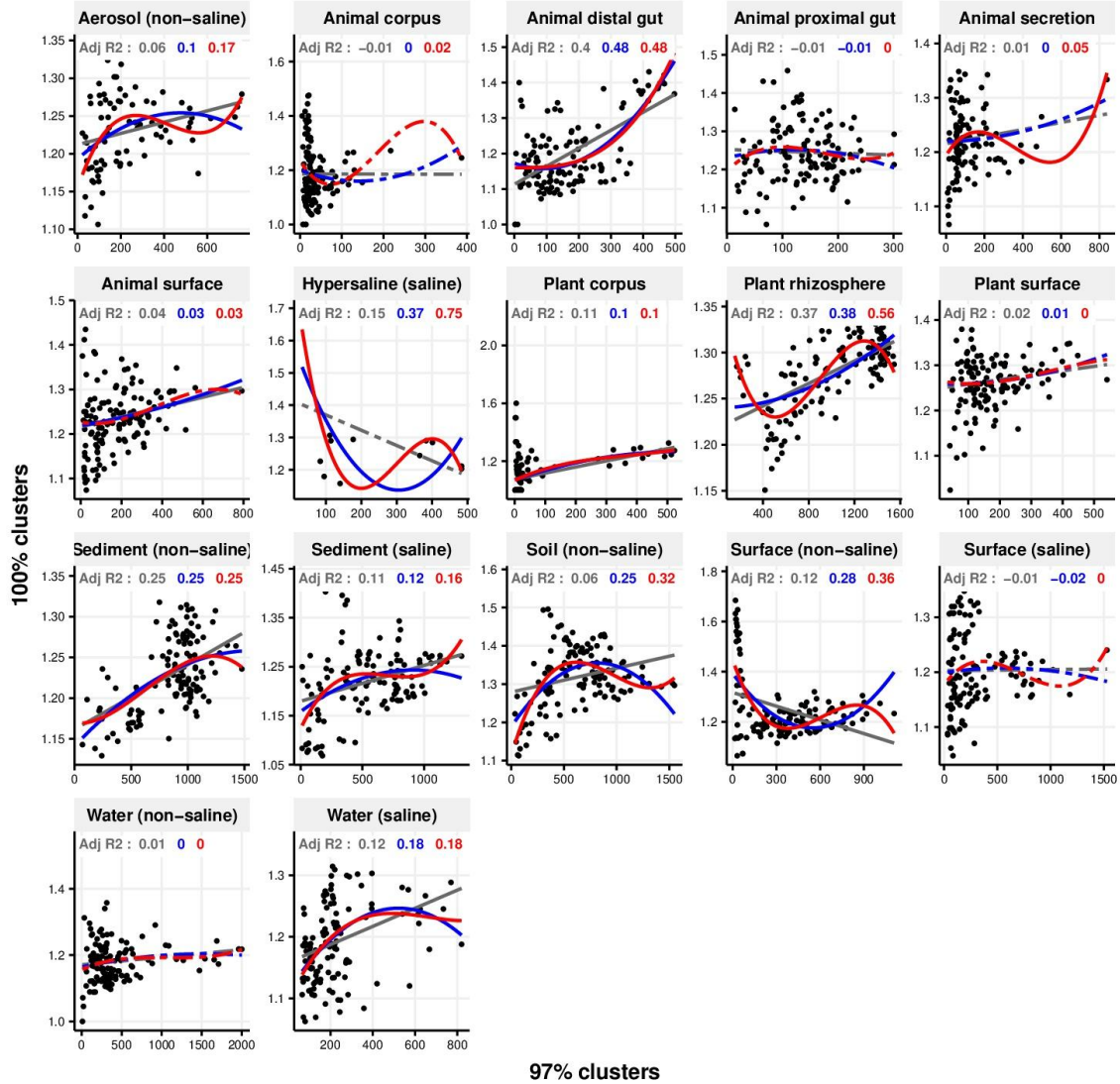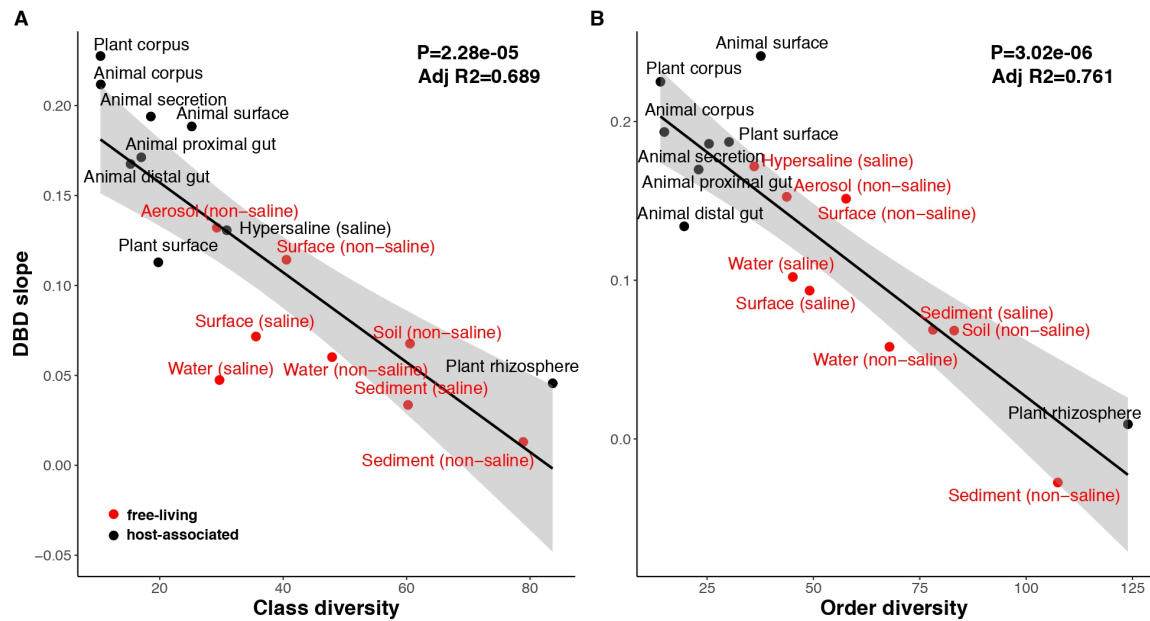


749

750

50

751

**85% clusters**

752

753
754

53

755
756
757

758    **Figure S11. DBD slope is higher in low-diversity (often host-associated)**
759    **environments.** The x-axis shows the mean number of (A) classes and (B) orders in each
760    biome; on the y-axis, DBD slope is the result from the GLMMs predicting diversification
761    as a function of the interaction between diversity and environment type at (A)
762    Order:Class and  (B) Family:Order ratio (**Supplementary Data file 1 Section 4)**. The
763    line represents a regression line, shaded areas depict 95% confidence limits of the fitted
764    values.
765
766



767
768
769
770
771
772
773
774
775
776
777
778
779

780