1    **Community diversity is associated with intra-species genetic diversity and**
2    **gene loss in the human gut microbiome**
3

4    Naïma Madi[1], Daisy Chen[2,3,^], Richard Wolff[4,^], B. Jesse Shapiro[1,5,6,7,*], and Nandita Garud[4,8,*]

5

6    1.  Département de sciences biologiques, Université de Montréal, Canada;
7    2.  Computational and Systems Biology, University of California, Los Angeles
8    3.  Bioinformatics and Systems Biology Program, University of California, San Diego[SEP]
9    4.  Department of Ecology and Evolutionary Biology, University of California, Los Angeles
10   5.  Department of Microbiology and Immunology, McGill University, Canada;
11   6.  McGill Genome Centre, McGill University, Canada
12   7.  Quebec Centre for Biodiversity Science, Canada
13   8.  Department of Human Genetics, University of California, Los Angeles
14

15   * Correspondence to jesse.shapiro@mcgill.ca and ngarud@ucla.edu. These authors contributed

16   equally.

17   ^ These authors contributed equally.

# Abstract

The human gut microbiome contains a diversity of microbial species that varies in composition over time and across individuals. These species are comprised of diverse strains, which are known to evolve by mutation and recombination within hosts. How the ecological process of community assembly interacts with sub-species diversity and evolutionary change is a longstanding question. Two hypotheses have been proposed based on ecological observations and theory: Diversity Begets Diversity (DBD), where taxa tend to become more diverse in already diverse communities, and Ecological Controls (EC), where higher community diversity impedes diversification within taxa. Recently we showed with 16S rRNA gene amplicon data from the Earth Microbiome Project that DBD is detectable in natural bacterial communities from a range of environments at high taxonomic levels (ranging from phylum to species-level), but that this positive relationship between community diversity and within-taxon diversity plateaus at high levels of community diversity. Whether increasing community diversity is associated with sub-species genetic diversity within microbiomes, however, is not yet known. To test the DBD and EC hypotheses at a finer genetic resolution, we analyzed sub-species strain and nucleotide variation in static and temporally sampled shotgun sequenced fecal metagenomes from a panel of healthy human hosts. We find that both sub-species single nucleotide variation and strain number are positively correlated with community diversity, supporting DBD. We also show that higher community diversity predicts gene loss in a focal species at a future time point and that community metabolic pathway richness is inversely correlated with the pathway richness of a focal species. These observations are consistent with the Black Queen Hypothesis, which posits that genes with functions provided by the community are less likely to be retained in a focal species' genome. Together, our results show that DBD and Black Queen may operate simultaneously in the human gut microbiome, adding to a growing body of evidence that these eco-evolutionary processes are key drivers of biodiversity and ecosystem function.

## Introduction

44

45  Our understanding of the evolution and diversification has been enriched by experimental studies
46  of bacterial isolates in the laboratory, but it remains a challenge to study evolution in the context
47  of more complex communities (Lenski, 2017). Ongoing advances in culture-independent
48  technologies have allowed us to study bacteria in the complex and dense communities in which
49  they naturally occur (Garud and Pollard, 2020). Within a community, individual players engage in
50  many negative and positive ecological interactions. Negative interactions can originate from
51  competition for resources and biomolecular warfare (Hibbing et al., 2010), while positive
52  interactions can stem from secreted metabolites that are used by other members of the community
53  (cross feeding) (Venturelli et al., 2018). These ecological interactions can create new niches and
54  selective pressures, leading to eco-evolutionary feedbacks whose nature are yet to be fully
55  understood.

56

57  Ecological interactions can yield positive or negative effects on the diversification of a focal
58  species. Under the "Diversity Begets Diversity" (DBD) hypothesis, higher levels of community
59  diversity increase the rate of speciation (or diversification, more generally) due to positive
60  feedback mechanisms such as niche construction (Calcagno et al., 2017; Schluter and Pennell,
61  2017). By contrast, the "Ecological Controls" (EC) hypothesis posits that competition for a limited
62  number of niches at high levels of community diversity results in a negative effect on further
63  diversification. Metabolic models predict that DBD may initially spur diversification due to cross
64  feeding, but the diversification rate eventually slows and reaches a plateau as metabolic niches are
65  filled (San Roman and Wagner, 2021). These theoretical predictions are largely supported by our
66  previous study involving 16S rRNA gene amplicon sequencing data from the Earth Microbiome
67  Project, in which we observed a generally positive relationship between community diversity and
68  focal-taxon diversity at most taxonomic levels, reaching a plateau at the highest levels of diversity
69  (Madi et al., 2020). A recent experiment on soil bacteria also found evidence of DBD at the family
70  level, most likely driven by niche construction and metabolic cross-feeding (Estrela et al., 2022).
71  Both of these studies show that DBD shapes microbial communities at higher taxonomic levels –
72  involving community assembly and species sorting – but lacked the genetic resolution to
73  interrogate sub-species strain-level dynamics such as strain colonization dynamics, polymorphism
74  levels, and gene gain and loss events. Moreover, these studies also lacked time series data to enable

75   directly tracking the dynamics of DBD, in which community diversity at one time point influences

76   diversity of a focal species in a future time point.

77

78   Like DBD and EC, the Black Queen Hypothesis (BQH) also makes predictions about the effects

79   of community diversity on the evolutionary dynamics of a focal species. BQH predicts that a focal

80   species will be less likely to encode genes with functions provided by other members of the

81   surrounding community, if such functions are "leaky" and available as diffusible public goods

82   (Morris et al., 2014, 2012). Gene loss may even be adaptive, provided that there is a cost to

83   encoding and expressing the relevant genes (Albalat and Cañestro, 2016; Koskiniemi et al., 2012;

84   Simonsen, 2022). The BQH has been invoked to explain the distribution of genes involved in

85   vitamin B metabolism (Sharma et al., 2019) and iron acquisition (Vatanen et al., 2019) in the gut

86   microbiome, but we still lack a complete picture of how the BQH applies to natural microbial

87   communities.

88

89   Here we examine evidence for DBD and BQH in human gut microbiome data at a sub-species

90   level. We use static and temporal shotgun metagenomic data from a large panel of healthy hosts

91   from the Human Microbiome Project as well as from the same individual sampled almost daily

92   over the course of one year (Poyet et al., 2019). In our previous study, we found strong support for

93   DBD in the animal gut compared to more diverse microbiomes such as soils and sediments which

94   were closer to a plateau of diversity (Madi et al., 2020). As such, the human microbiome represents

95   an ideal model for further studying DBD dynamics. Using metagenomic data to track gene gain

96   and loss events within a focal species allows us to simultaneously test the predictions of DBD and

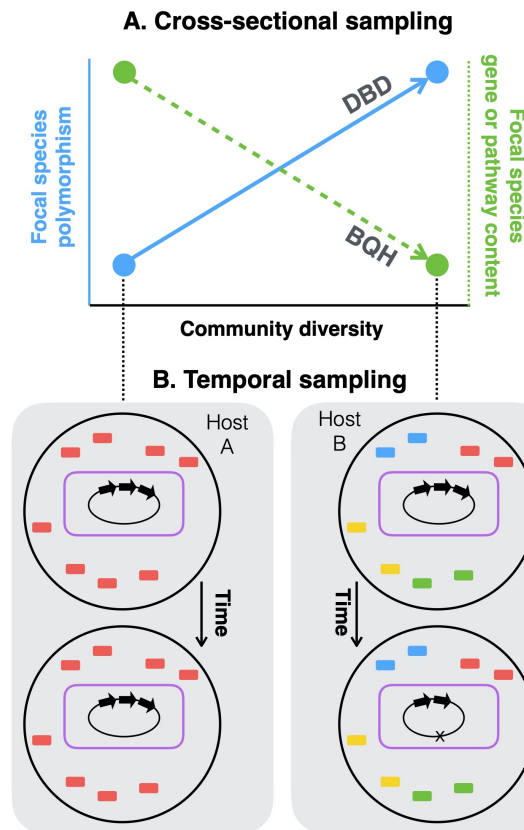97   BQH, which are not mutually exclusive (**Figure 1**).

98

4

**Figure. 1. Diversity begets diversity (DBD) and the Black Queen Hypothesis (BQH) illustrated. (A)** Each point represents a different individual microbiome sampled from a cross-sectional cohort, ranging from low to high community diversity along the x-axis. Under DBD, high community diversity is associated with high focal species polymorphism (blue line) yielding a positive diversity slope. Alternatively (not shown) a niche filling model as described by the EC hypothesis would yield a flat or negative diversity slope. Under the BQH, high community diversity is associated with fewer genes or metabolic pathways encoded by the focal species (dashed green line). **(B)** Two hypothetical gut microbiomes (hosts) sampled over time, illustrating the predictions of DBD and BQH. The bacterial community is shown as small rectangles with different bacterial species in different colors. The focal species is shown as an enlarged rectangle outlined in violet, containing a circular genome encoding genes (bold black arrows). In Host A, the microbiome is not diverse (only one 'red' species) and does not affect the diversity of the focal species. In Host B, a more diverse community at the first time point selects for increased polymorphism (point mutation shown as an 'X') and the loss of genes or pathways within the focal species. For simplicity these are depicted as *de novo* mutation or gene loss events of resident lineages, but could equally come about due to invasions of new strains or changes in the relative abundance of preexisting strain diversity.

5

# Results

We assess evidence for the DBD hypothesis within species using two shotgun metagenomic datasets. First, we analyze data from a panel of 249 healthy hosts (Human Microbiome Project Consortium, 2012; Lloyd-Price et al., 2017), in which stool samples were collected 1-3 times at approximately 6-month intervals. Second, we analyze data from a single individual sampled more densely (206 samples) over the course of ~18 months (Poyet et al., 2019). We analyze both cross-sectional and temporal data to understand the relationship between community diversity and genetic diversity at the sub-species level.

To test the DBD hypothesis, we examined several metrics of community diversity and intra-species diversity and calculated the diversity slope (**Figure 1**). To quantify community diversity, we calculated Shannon diversity and richness at the species level. To control for variation in sequencing depth across samples, richness was computed on rarefied data. We also used unrarefied data and included the number of reads per sample as a covariate in our models, yielding similar results (described below). To quantify intra-species diversity, we used a reference genome-based approach to call single nucleotide variants (SNVs) and gene copy number variants (CNVs) within each focal species and computed polymorphism rates, measured as the fraction of synonymous nucleotide sites in the core genome with intermediate allele frequencies (between 0.2 and 0.8) within a host (Methods). As an additional metric of intra-species diversity, we inferred the number of strains within each species using StrainFinder (Smillie et al., 2018).

**Community diversity is associated with sub-species polymorphism in the human gut microbiome**

We began by plotting the relationship between community diversity and intra-species polymorphism rate in cross-sectional HMP metagenomes (**Figure 2A and B**). The slope of this relationship (which we call the diversity slope; **Figure 1**) provides an indicator of the extent of DBD (positive slope) or EC (flat or negative slope). As a descriptive analysis, we first computed Spearman correlations between Shannon diversity and intra-species polymorphism rate. Out of the 68 bacterial species with sufficient prevalence (present in at least four samples), we found 15 significant correlations ($P < 0.05$, uncorrected for multiple tests), of which 14 were positive (**Fig S1**). Similarly, we found 18 significant correlations between richness and intra-species

6

147    polymorphism rate (**Fig S2**), all of which were positive. These positive associations are broadly

148    consistent with the DBD hypothesis, although we cannot establish the direction of causality in this

149    cross-sectional data.

150

151    The relationship between polymorphism rate and community diversity was found to be non-linear

152    (**Figures 2A, B, S1 and S2**). Polymorphism rates across HMP hosts span several orders of

153    magnitude ($10^{-5}$/bp to $10^{-2}$/bp), largely due to the fact that strain content is variable across hosts.

154    Polymorphism rates of ~$10^{-2}$/bp or more are inconsistent with within-host diversification of a

155    single colonizing lineage, and instead represent mixtures of multiple strains that diverged before

156    colonizing a host. By contrast, rates <$10^{-4}$/bp are more consistent with a single strain colonizing

157    the host (Garud et al., 2019).

158

159    To more formally test the predictions of DBD and to account for non-linear relationships between

160    polymorphism and community diversity, we used generalized additive models (GAMs). Using

161    GAMs, we are able to model non-linear relationships and account for random variation in the

162    strength of the diversity slope across samples, human hosts, and bacterial species (Methods). We

163    find that GAMs support the overall positive association between within-species polymorphism and

164    Shannon diversity (GAM, *P*=0.031, Chi-square test) as well as between within-species

165    polymorphism and community richness after controlling for coverage as a covariate (*P*=0.017) or

166    rarefying samples to an equal number of reads (*P*=1.93e-05) (**Fig. S3**). While the polymorphism-

167    community diversity relationships were generally positive, it appears that polymorphism reaches

168    a plateau at high levels of community richness (**Fig S2**), as supported by GAMs (**Fig S3 B,C**); see

169    **Table S1** and **Supplementary File 1** for additional model details.

170

171    These generally positive correlations between focal species polymorphism and species-level

172    measures of community diversity also hold when community diversity is measured at higher

173    taxonomic levels; specifically, polymorphism rate was significantly positively associated with

174    Shannon diversity calculated at the genus and family levels (GAMs, *P*<0.05, Chi-square test) (**Fig**

175    **S4**). However, polymorphism rate was not significantly associated with Shannon diversity

176    calculated at the highest taxonomic levels (order, class and phylum, GAMs, *P*>0.05, Chi-square

177    test). The positive correlation between polymorphism rate and richness held at all taxonomic levels

178    (GAMs, *P*<0.05, Chi-square test) (**Fig S4**, **Table S2**, **Supplementary File 1**). Overall, these results

179    are consistent with DBD acting within the human gut microbiome at most taxonomic levels, as

180    previously observed in environmental samples (Madi et al., 2020) and experimental soil

181    communities (Estrela et al., 2022).

182

183    **Community diversity is associated with sub-species strain diversity**

184    To more explicitly account for the strain structure within hosts, we next inferred the number of

185    strains per focal species with StrainFinder (Smillie et al., 2018) (Methods) and used strain number

186    as another quantifier of intra-species diversity. Strain-level variation has important functional and

187    ecological consequences; among other things, strains are known to engage in interactions that

188    cannot be predicted from their species identity alone (Goyal et al. 2021). How ecological processes

189    at the strain level affect and are affected by community composition and dynamics, however,

190    remains poorly characterized.

191

192    We found that the number of strains per focal species follows an approximately linear relationship

193    with community Shannon diversity (**Figure 1C** and **S5)**. We therefore calculated Pearson

194    correlations between community diversity and the number of strains per focal species. Out of the

195    134 species for which strains were inferred (Methods), we found a total of 23 significant

196    correlations between Shannon diversity and strain number (*P*<0.05), of which 21 were positive

197    (**Fig S5**). By contrast, out of the 16 significant correlations between richness and strain number,

198    13 were negative (**Figures 1D** and **S6**). We note that the 16 species with a significantly positive

199    Shannon-strain number correlation were completely non-overlapping with the 13 with a significant

200    negative richness-strain number correlation, suggesting species-specific effects.

201

202    We next used generalized linear mixed models (GLMMs) to investigate the relationship between

203    the number of strains per focal species and community diversity, while taking into account

204    coverage per sample as a covariate and variation between species, hosts and samples as random

205    effects. The number of strains per focal species was positively correlated with community Shannon

206    diversity (GLMM, *P*= 1.549e-04, likelihood ratio test (LRT)) (Table S3, supplementary file 1).

207    This is consistent with the positive correlation between polymorphism rates and Shannon diversity

208    and is also generally concordant with DBD.

209

210    While Shannon diversity was positively correlated with strain number, species richness was

211    negatively correlated with strain number (GLMM, $P$=1.5e-05, LRT) (**Table S3**, **Supplementary**

212    **File 1**). The negative relationship with richness was unlikely to be confounded by sequencing

213    depth, since the same result was obtained using rarefied data, albeit with a weaker negative

214    relationship (GLMM, $P$=0.037, LRT) (**Table S3**, **Supplementary File 1**). The negative strain

215    number-richness relationship also held at all other taxonomic ranks (GLMM, $P$<0.05, LRT) (**Table**

216    **S4**, **Supplementary File 1**), while the strain number-Shannon diversity relationship was generally

217    positive (**Fig S7**). Together, these results show that richness and Shannon diversity are both

218    positively correlated with polymorphism rates, consistent with DBD, whereas richness and

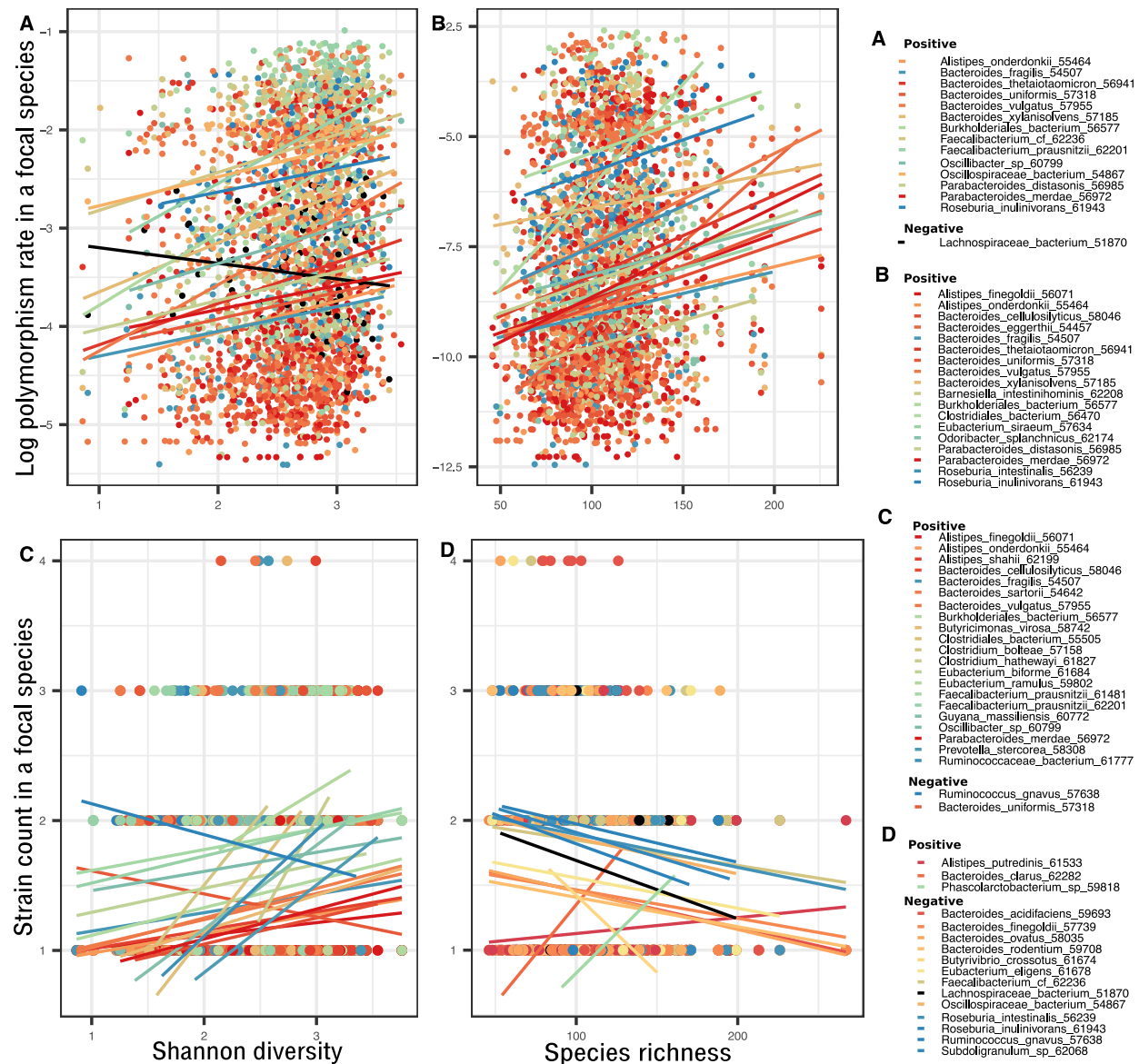219    Shannon diversity have contrasting correlations with strain number.

**Fig 2. Relationship between intra-species diversity and community diversity in human gut microbiomes. (A)** Intra-species polymorphism rate versus Shannon diversity (15 out of 68 species have a significant Spearman correlation with Shannon alpha diversity; 14 are positive and 1 is negative, shown with a black trendline). **(B)** Intra-species polymorphism rate versus rarefied richness (18 species out of 68 are significantly correlated; all are positive). **(C)** Intra-species strain number versus Shannon diversity (23 out of 134 species are significantly correlated with richness; 21 are positive). **(D)** Intra-species strain number versus rarefied richness (16 out of 134 species are significantly correlated with richness; 13 are negative). Each species is represented by a different color. Only significant correlations (points and lines) are plotted ($P<0.05$, Spearman in A and B and Pearson in C and D). Species present in fewer than 4 samples were not considered. The Y-axis in A and B is on a log10 scale.

10

232 **Testing DBD over time in the gut**

233 Our analyses thus far have considered only individual time points, which represent static snapshots

234 of the dynamic processes of community assembly and evolution in the microbiome. To test the

235 effects of DBD over time, we analyzed 160 HMP hosts with multiple time points, in which the

236 same person was sampled 2-3 times ~6 months apart. Under a DBD model, we expect community

237 diversity at an earlier time point to result in higher within-species polymorphism at a future time-

238 point. To test this expectation, we defined 'polymorphism change' as the difference between

239 polymorphism rates at the two time points (Methods). We also investigated the effects of

240 community diversity on gene loss and gain events within a focal species, as such changes in gene

241 content are known to occur frequently within host gut microbiomes (Garud et al., 2019; Groussin

242 et al., 2021; Zhao et al., 2019). Here a gene was considered absent if its coverage ($c$) was <0.05

243 and present if $0.6 \leq c \leq 1.2$ (Methods). As in the cross-sectional analyses above, we also

244 controlled for sequencing depth of the sample and excluded genes with aberrant coverage or that

245 were present in multiple species.

246

247 In HMP samples, polymorphism change showed no significant relationships with community

248 diversity at the earlier time point (**Fig S8**, GAM, P=0.4, P=0.64 and P=0.497, for Shannon, richness

249 and rarefied richness respectively), nor did gene gains show any relationships (**Fig S9**, GLMM,

250 *P*= 0.733, *P*= 0.617 and *P*= 0.508, LRT for Shannon, richness and rarefied richness respectively)

251 (**Supplementary File 1**). These results suggest that DBD is negligible or undetectable over ~6-

252 month time lags in the human gut. By contrast, we found that gene loss in a focal species between

253 two consecutive time points was positively correlated with community diversity at the earlier time

254 point (**Figure 3**, **S10**, GLMM, *P*= 0.028, *P*= 0.036 and *P*= 0.01, LRT for Shannon, richness and

255 rarefied richness respectively) (**Table S5**, supplementary file 1). Elevation of gene loss in more

256 diverse communities is consistent with the BQH, which we investigate in further detail below.

257 Most species in the HMP hosts lost fewer than ten genes over ~6 months, but occasionally

258 hundreds of genes were lost (**Figure 3**), suggesting a mixture of *de novo* deletion of a few genes

259 as well as selection of strains encoding fewer genes in more diverse communities.
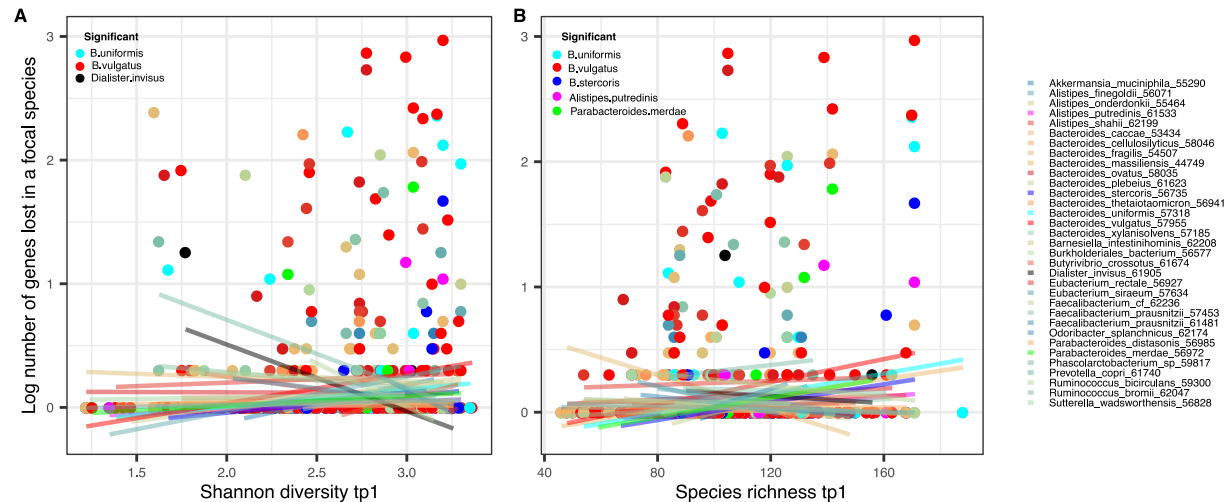
260

261

11

**Figure 3. Community diversity predicts focal species gene loss in the human microbiome over time.** Plots show the number of genes lost between two time points in a focal species as a function of **(A)** Shannon diversity, or **(B)** richness measured at an earlier time point in the HMP. The species listed in the top left of each panel are those with individually significant correlation (Pearson correlation, $P<0.05$) out of 33 listed on the right. Only the 33 species present in more than three samples with at least one gene loss event were tested. All significant correlations were positive except for the correlation between gene loss and Shannon diversity in *Dialister invisus* in **A** (black). The Y-axis in A and B is on a log10 scale.

To assess the evidence for DBD at higher temporal resolution, we analyzed shotgun metagenomic data from the most frequently sampled healthy individual (host *am*) from a previous study (Poyet et al., 2019). This individual donated stool samples that were sequenced over 18 months with a median of one day (mean of 2.6 days) between time points. In this data, we tracked both polymorphism change and gene gains and losses between two successive time points in *Bacteroides vulgatus*, the most abundant species across samples (mean coverage=58.5; median=54.2; Methods). Polymorphism and strain-level diversity within *B. vulgatus* were positively associated with community diversity in the HMP cross-sectional data (**Figures S1**, **S2**, **S5**). We would therefore expect similar associations in time series data.

We asked whether community diversity in the gut microbiome at one time point could predict increases in *B. vulgatus* polymorphism at the next time point, typically a few days later. Consistent with DBD, Shannon diversity at the earlier time point was positively correlated with changes in polymorphism (**Figure 4A**, Pearson correlation $P$=0.002). Notably, this positive correlation was

12

285    not evident in the HMP time series, perhaps due to insufficient density of sampling to capture rapid

286    dynamics. Even when individual correlations were tested in HMP data, *B.vulgatus* did not show a

287    significant relationship (Pearson, *P*>0.05).

288

289    Consistent with observations from HMP time series (**Figure 3**), we found a positive relationship

290    between gene loss and Shannon diversity in *B. vulgatus* in individual *am* (**Figure 4B,** Pearson

291    correlation, *P*=0.06). The positive association with gene loss was mirrored by a negative

292    association with gene gain, although both with borderline statistical significance due to relatively

293    few observed gain or loss events over these short time intervals (**Figure 4C,** Pearson correlation,

294    *P*=0.09). All genes gained and lost in *B. vulgatus* in *am* were annotated as hypothetical proteins.

295    Neither polymorphism change nor gene gains or losses in *B. vulgatus* were correlated with species

296    richness in individual *am* (**Figure S11A, B, and C,** Pearson correlation, *P*>0.2).

297

298    Overall, these results are consistent with community diversity promoting changes within the *B.*

299    *vulgatus* genome over timescales of a few days. We note that this is an example of one abundant

300    species in one well-sampled individual and may not generalize to other species and hosts.

301    However, it does suggest that changes in polymorphism captured over daily time scales could be

302    obscured over the timescales on the order of months, as reflected in the HMP samples. The

303    association between Shannon diversity and gene loss, in both HMP and Poyet time series, is

304    suggestive of adaptive gene loss as posited by the Black Queen Hypothesis (BQH). Under BQH,

305    genes are lost from an individual genome when their functions are provided by other members of
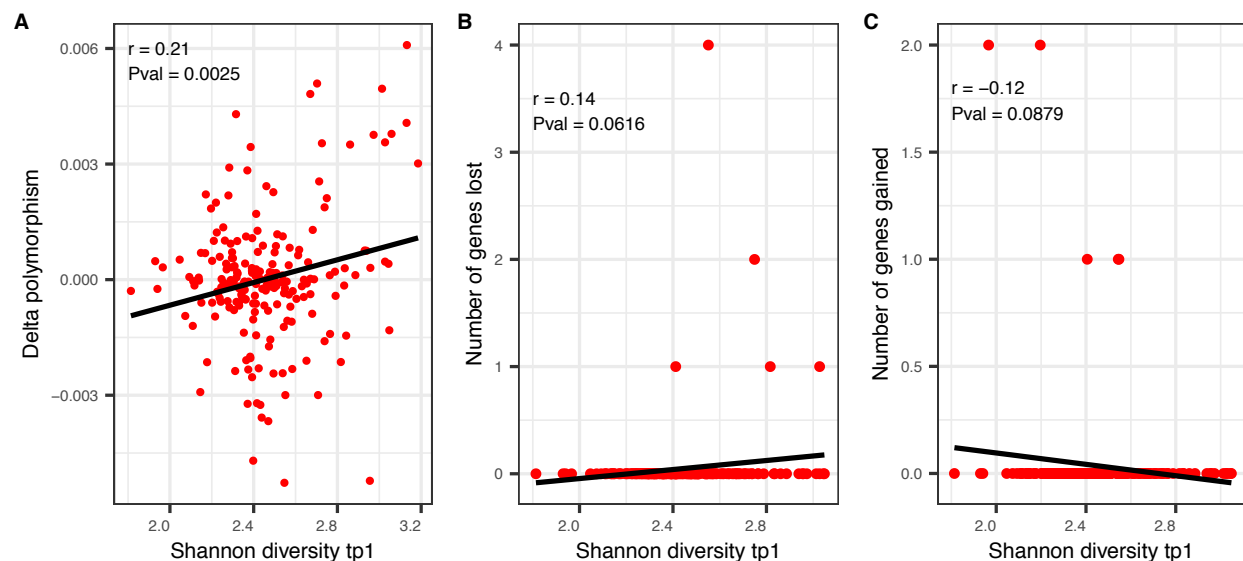
306    a (diverse) community.

307

**Figure 4**. **Correlations between polymorphism and gene content changes in *B. vulgatus* and Shannon diversity in one human gut microbiome over 18 months**. A) Correlation between polymorphism change (time point 2 – time point 1) and Shannon diversity at time point 1 (tp1; the earlier time point). B) Correlation between the number of genes lost between tp1 and tp2 and the Shannon diversity at tp1. C) Correlation between the number of genes gained between tp1 and tp2 and the Shannon diversity at tp1. Each point represents the change measured between a pair of consecutive time points. Pearson correlations and *P*-values are reported in each panel.

## Testing the Black Queen hypothesis in the human gut microbiome

To further assess evidence for the BQH in the HMP data, we tested the hypothesis that a focal species encodes fewer genes in a community that collectively harbors more genes. This would be expected under adaptive gene loss, provided that the genes encoded by the community provide 'leaky' functions to the focal species. Contrary to this simple expectation, we observed a significant positive relationship between community gene richness and focal species gene richness (see Methods for computation of gene richness) (**Figures 5A and S12A**; GAM, *P*=2.92e-06, Chi-square test) (Table S6, supplementary file 1). By estimating Spearman correlation between gene richness per focal species and community gene richness, we found that out of 134 species, 42 had significant correlations, of which 39 were positive (**Fig S13**). This result is inconsistent with a simple version of the BQH acting on individual gene families assuming that all gene functions are equally 'leaky'. It is, however, broadly consistent with DBD, provided that gene content is

14

330    correlated with polymorphism rate, which we already showed to be correlated with community

331    diversity (**Figure 2**). In other words, DBD is supported both in terms of within-species single

332    nucleotide polymorphism and gene content variation.

333

334    Next, we tested the hypothesis that the BQH acts at the level of metabolic pathways rather than

335    individual gene families. Specifically, cellular pathways that are encoded by the community need

336    not be encoded by a focal species provided that the pathway product or function is leaky.

337    Consistent with the BQH acting at the pathway level, we found that community pathway richness,

338    measured as the number of pathways present with non-zero abundance inferred with HUMAnN2

339    (Franzosa et al., 2018) (Methods) was negatively correlated with focal species pathway richness

340    (**Figures 5B, S12B;** GAM, $P$<2e-16, Chi-square test) (Table S6, supplementary file 1). When

341    testing 239 prevalent species, we found 107 significant Spearman correlations ($P < 0.05$), of which

342    95 (89%) were positive (**Fig S14**). Note that three species (*Escherichia coli, Enterobacter cloacae*

343    and *Klebsiella pneumoniae*, shown respectively with green, orange, and red points and trendlines

344    in **Figure 5B**) with particularly high pathway richness had much steeper negative slopes, but they

345    are not solely responsible for the overall negative trend **(Fig 5B)**.
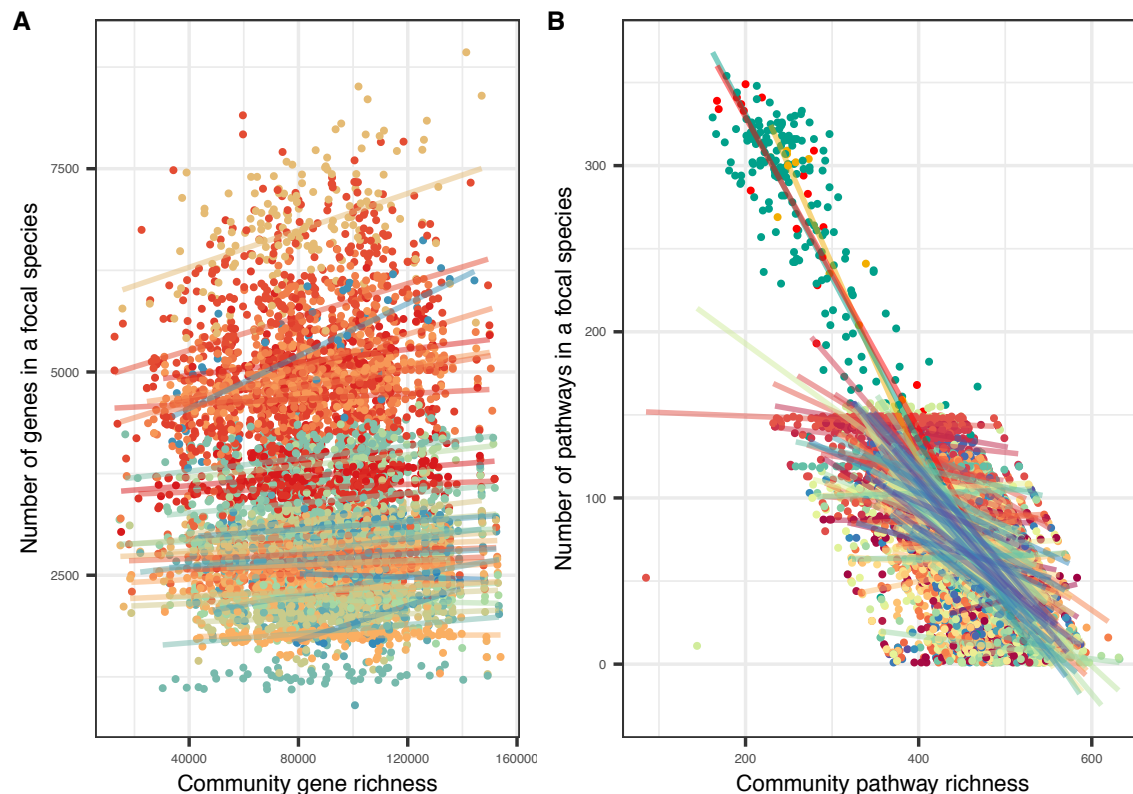
346

347

348

15

**Figure 5. Testing predictions of the Black Queen Hypothesis in the human microbiome.** Plots show correlations between (A) the number of genes present in a focal species and genes present in other members of the community (out of 124 species tested, 42 had significant correlations, of which 39 positive and 3 negative), and B) the number of pathways present in focal species and community pathway richness (out of 239 species tested, 107 had significant correlations, of which 95 negative and 12 positive). Only species present in at least 4 samples with a significant Spearman correlation ($P<0.05$) are plotted.

# Discussion

In this paper we investigated whether community diversity begets genetic diversity within species in gut microbiota using static and temporally resolved fecal shotgun metagenomic data from a panel of healthy hosts. In support of the DBD hypothesis, we found that focal species often have higher polymorphism rates and strain counts in more diverse communities, whether community diversity was estimated with Shannon index or species richness. The same pattern held when community diversity was estimated at higher taxonomic ranks, consistent with our previous analysis of amplicon sequence data across environments (Madi et al., 2020) and a recent experimental study of soil bacteria community assembly (Estrela et al., 2022). Together, these

16

366    results indicate that the DBD hypothesis is relevant at multiple taxonomic levels, and extends past
367    the species level to the sub-species genetic level.

368         While sub-species strain diversity is generally positively correlated with Shannon
369    diversity, it is inversely correlated with species richness, suggesting that the ability of strains to
370    colonize a host may be associated with higher community evenness rather than their total count.
371    Although Shannon diversity is considered to be more robust and informative than richness in
372    estimating bacterial diversity (He et al., 2013; Reese and Dunn, 2018), we observe the same
373    contrasting results between Shannon and richness when community diversity is calculated at
374    higher taxonomic levels, suggesting that this pattern is not due to artifacts such as sequencing
375    effort.

376         Another study also recently found evidence for eco-evolutionary feedbacks in the HMP, in
377    the form of a positive relationship between evolutionary modifications or strain replacements in a
378    focal species and community diversity (Good and Rosenfeld, 2022). Using a model, they further
379    showed that these eco-evolutionary dynamics could be explained by resource competition and did
380    not require the cross-feeding interactions previously invoked to explain DBD at higher taxonomic
381    levels (Estrela et al., 2022; San Roman and Wagner, 2021, 2018). This could be because cross-
382    feeding operates at the family- or genus- level, and is less relevant as a finer-scale evolutionary
383    process.

384         Perhaps compatible with the recent work, we found that community diversity predicts gene
385    loss in a future time point and that community pathway richness is negatively correlated with
386    pathway richness of a focal species. This suggests that both DBD and BQH might be at play in the
387    gut microbiome, in which high community diversity may simultaneously select for diversification
388    (at the SNV and strain level) while also selecting for adaptive gene loss as predicted by BQH (that
389    is, relaxed selective pressure to maintain pathways already provided by the community). While it
390    is possible that gene deletion events could explain the loss of functional metabolic pathways, it is
391    also possible that there is a propensity for strains with fewer pathways to colonize hosts with more
392    complex communities. Higher resolution time series data can help to disentangle these possibilities
393    as well as to more deeply quantify the effect of BQH on microbiome diversity.

394         The tendency for reductive genome evolution in bacteria has already been reported by
395    comparing hundreds of genomes (Albalat and Cañestro, 2016; Koskiniemi et al., 2012; Puigbò et
396    al., 2014). Genome reduction is also a hallmark of endosymbiotic bacteria, which receive many

17

397    metabolites from their hosts (McCutcheon and Moran, 2012; Nikoh et al., 2011). It has been shown

398    that uncultivated bacteria from the gut have undergone considerable genome reduction, which may

399    be an adaptive process that results from use of public goods (Nayfach et al., 2019). Our findings

400    suggest that genome reduction in the gut is higher in more diverse gut communities, and future

401    work could establish whether this effect is indeed due to metabolic cross-feeding as posited by

402    some models (Estrela et al., 2022; San Roman and Wagner, 2021, 2018), but not others (Good and

403    Rosenfeld, 2022).

404         The BQH may help explain why the majority of gut microbial species remain recalcitrant

405    to cultivating under laboratory conditions (Nayfach et al., 2019; Walker et al., 2014). Specifically,

406    gut microbes may lack the necessary pathways to survive in culture in absence of their natural

407    counterparts that may otherwise provide essential goods. For instance, menaquinone and fatty

408    acids have been shown to promote the growth of uncultured bacteria, and both pathways were

409    missing from many uncultured bacteria identified in (Nayfach et al., 2019). Additionally, more

410    than 70% of the recent created Unified Human Gastrointestinal Genome (UHGG) collection lack

411    cultured representatives (Almeida et al., 2019).

412         As noted in our previous study (Madi et al., 2020), we cannot establish causal relationships

413    between community diversity and focal species diversity using cross-sectional survey data; doing

414    so requires controlled experiments. In the case of DBD, the correlations observed in naturally

415    occurring microbiomes are generally concordant with experimental (Estrela et al., 2022; Jousset

416    et al., 2016) and metabolic modeling studies (San Roman and Wagner, 2021), strengthening the

417    plausibility of the hypothesis. Although they also note that causality is difficult to establish, Good

418    and Rosenfeld (2022) suggest the importance of focal species evolution as a driver of changes in

419    community structure, as shown in an experimental study of *Pseudomonas* in compost communities

420    (Padfield et al., 2020). Further work is therefore needed to establish the extent and relative rates

421    of eco-evolutionary feedbacks in both directions. How these feedbacks among bacteria are

422    influenced by abiotic factors and by interactions with fungi, archaea, and phages also deserve

423    further study.

424         In summary, our results show support for both DBD and the BQH within the human gut

425    microbiome. Using metagenomic time series data, we find a positive association between

426    community diversity and sub-species strain-level diversity. Higher community diversity is also

427    associated with losses of genes and metabolic pathways in a focal species. Whether these reductive

428  genome evolution events are adaptive, as predicted by BQH, and if they can be explained by

429  metabolic cross-feeding, remains to be seen.

430
431
432

# Data and materials availability

434  The raw sequencing reads for the metagenomic samples used in this study were downloaded

435  from Human Microbiome Project Consortium 2012 and Lloyd-Price et al. (2017)

436  (URL: https://aws.amazon.com/datasets/human-microbiome-project/); and Poyet et al. 2019

437  (NCBI accession number PRJNA544527). All computer code for this paper is available at

438  https://github.com/Naima16/DBD_in_gut_microbiome.

439

440

# Methods

442

## Metagenomic analyses

444

### Estimation of species, gene, and SNV content of metagenomic samples

446  We used MIDAS (Metagenomic Intra-Species Diversity Analysis System, version 1.2,

447  downloaded on November 21, 2016) (Nayfach et al., 2016) to estimate within-species nucleotide

448  and gene content of raw metagenomic whole genome shotgun sequencing data for HMP1-2 and

449  Poyet et al. 2019 data. MIDAS relies on a reference database comprised of 31,007 bacterial

450  genomes that are clustered into 5,952 species, covering roughly 50% of species found in human

451  stool metagenomes from "urban" individuals. Described below are the parameters used to estimate

452  species abundances, SNVs, and gene copy numbers variants (CNVs) with MIDAS:

453

#### *Estimation of species content:*

455  To assess evidence for community diversity begetting genetic diversity, we estimated

456  species diversity and SNVs and CNVs by mapping reads to reference genomes. Since a component

457  of this work relies on quantifying polymorphism and CNV changes over time, we constructed a

19

458    "personal" reference database to avoid spurious inferences of allele frequency and CNV changes

459    due to errors in mapping of reads to regions of the genome shared by multiple species. This per-

460    host reference database was comprised of the union of all species present at one or more timepoints

461    so as to be as inclusive as possible to prevent reads from being "donated" to reference genome,

462    while also being selective to prevent a reference genome from "stealing" reads from a species truly

463    present.

464    To estimate the species relative abundances for each host x timepoint sample, we mapped

465    reads to 15 universal single-copy marker genes that are a part of the MIDAS pipeline (Nayfach et

466    al., 2016; Wu et al., 2013) and belong to the 5,952 species. A species with an average marker gene

467    coverage $\geq 3$ was considered present for the purposes of inferring SNVs and CNVs below. The

468    per-host database was constructed by including all species present at one or more timepoints with

469    coverage $\geq 3$.

470

471    ***Estimation of copy number variation:***

472    To estimate gene copy number variation (CNV) we mapped reads to the pangenomes of

473    species present in a host's personal database using Bowtie2 (Langmead and Salzberg, 2012) with

474    default MIDAS settings (local alignment, MAPID$\geq$94.0%, READQ$\geq$20, and ALN_COV$\geq$0.75).

475    Each gene's coverage was estimated by dividing the total number of reads mapped to a given gene

476    by the gene length. These genes included the aforementioned 15 universal single-copy marker

477    genes. A given gene's copy number ($c$) was estimated by taking the ratio of its coverage and the

478    median coverage of the single-copy marker genes.

479    With these copy number values, we estimated the prevalence of genes in the broader

480    population, defined as the fraction of samples with copy number $c \leq 3$ and $c \geq 0.3$ (conditional on

481    the mean single gene marker coverage being $\geq 5x$). For each species, we computed "core genes",

482    defined as genes in the MIDAS reference database that are present in at least 90% of samples

483    within a given cohort. Within-host polymorphism rates were computed in core genes.

484    However, orthologous genes present in multiple species can result in read stealing and read

485    donating. Thus, we excluded a set of genes belonging to a 'blacklist' comprised of genes present

486    in multiple species. This blacklist was constructed in Garud et al. 2019 using USEARCH (Edgar,

487    2010) to cluster all genes in human-associated reference genomes with a 95% identity threshold.

488    Since some genes may be absent from the MIDAS database that may also be shared across species,

489    we implemented another filter in Garud et al. 2019 in which genes with $c \geq 3$ in at least one sample

490    in our cohort was excluded from analysis of polymorphism rate or gene changes over time.

491

*Inferring single nucleotide variants (SNVs) within bacterial species*

493    To call SNVs, we mapped reads to a single representative reference genome as per the

494    default MIDAS software. Reads were mapped with Bowtie2, with default MIDAS mapping

495    thresholds: global alignment, MAPID≥94.0%, READQ≥20, ALN_COV≥0.75, and MAPQ≥20.

496    Species were excluded from further analysis if reads mapped to $\leq 40\%$ of their genome. We further

497    excluded samples from further analysis if they had low median read coverage $(\overline{D})$ at protein coding

498    sites. Specifically, samples with $\overline{D} < 5$ of across all protein coding sites with nonzero coverage

499    were excluded. This MIDAS SNV output was then used subsequently for computing within-

500    species polymorphism rates and inferring the number of strains present for each species in each

501    sample (see below).

502    To compute polymorphic rates, additional bioinformatic filters were imposed to avoid read

503    stealing and donating across different species. First, we did not call SNVs in blacklisted genes

504    present in multiple species. Additionally, we excluded sites in a given sample if $D < 0.3\overline{D}$ or $D >$

505    $3\overline{D}$ as these sites harbor coverage anomalously low or high compared to the genome-wide average

506    $\overline{D}$. An additional coverage threshold requirement of 20 reads/site was imposed for inclusion of

507    SNVs in the polymorphism rate computation.

508

509

*Shannon diversity, species richness and polymorphism rate calculations*

511    Shannon diversity and richness were computed within each sample by including any species with

512    abundance greater than zero. Rarefied species richness estimates are based on HMP1-2 samples

513    rarefied to 20 million reads and Poyet samples rarefied to 5 million reads.

514

515    The polymorphism rate of a species in a sample was computed as the proportion of synonymous

516    sites in core genes with intermediate allele frequencies ($0.2 \leq f \leq 0.8$). This is quantitatively similar

517    to the more traditional population genetic measure of heterozygosity, $H=E[2f(1-f)]$, in which

518    intermediate frequency alleles contribute the most weight. By computing polymorphism with the

21

519    criteria $0.2 \leq f \leq 0.8$, we avoid inclusion of low frequency sequencing errors, which can otherwise

520    greatly influence the mean heterozygosity.

521

522    *Temporal changes in polymorphism rates and gene content*

523    Delta polymorphism (or changes in polymorphism) was computed as the difference in

524    polymorphism rates between time points. Gene gains and losses between time points were

525    computed by identifying genes with copy number c $<=0.05$ (indicating gene absence) in one

526    sample and $0.6 <= c <= 1.2$ in another (indicating single copy gene presence). These thresholds

527    were used in Garud et al. 2019 when inferring gene changes in temporal data and reflect a range

528    of copy numbers expected in either the absence of a gene or presence of a single copy of a gene.

529    Higher copy numbers were not considered to avoid confounding our analysis with read stealing or

530    donating among different species. Filters for coverage and blacklisted genes were applied as

531    described above.

532

533    *Strain number inference*

534    We used StrainFinder (Smillie et al., 2018) to infer the number of strains present for each species

535    in each HMP1-2 metagenomic sample. To do so, we used allele frequencies from MIDAS SNV

536    output, generated as described above. For each species in each host, all multi-allelic sites with

537    coverage of 20x or greater were passed as input to StrainFinder. Species in which no sites passed

538    the 20x threshold were assumed to have only a single strain. StrainFinder was then run on each

539    sample separately for strain number 1, 2, 3, and 4, and the optimal strain number was chosen based

540    on BIC. This range of strain number was chosen for biological reasons. Based on multiple analyses

541    of the densely longitudinally sampled metagenomic data from four healthy hosts in Poyet et al, a

542    maximum of three strains were shown to be present at any one time within a host for the ~30 most

543    prevalent species (Poyet et al. 2019, Wolff et al. 2021, Zheng et al. 2020). Thus, four strains were

544    chosen as the maximum to accommodate the range of observed possibilities, as well as possible

545    rare cases outside of this, without overfitting.

546

547    *Gene and pathway richness*

548    To determine gene richness of each sample, we used the default MIDAS threshold of 0.35 copy

549    number to define gene presence and absence. All genes from the species' pangenome with

550    minimum read-depth of 1, including core and accessory genes, were considered for this analysis.

551    Finally, we define "community gene richness" of a sample, with respect to a focal species, as the

552    number of gene clusters present in any of the species in the sample, excluding the focal species.

553    Gene clusters are defined as any set of genes with 95% nucleotide identity.

554

555    In addition to examining gene sets, we utilized previously generated functional profiling output

556    from      HUMAnN      2.0      (Franzosa      et      al.,      2018)      (downloaded      from

557    https://www.hmpdacc.org/hmmrc2/) to estimate pathway richness in each species present in a

558    sample. HUMAnN 2.0 takes in whole genome metagenomes and reports gene family (UniRef) and

559    metabolic pathway (MetaCyc) abundances in reads per kilobase (RPK); here, we count all

560    pathways with nonzero RPK as present in a sample.

561

562    **Statistical analyses**

563

564    *Model construction and evaluation*

565    Using data from the HMP and Poyet et al. 2019, we examined the relationship between intra-

566    species diversity and gut microbiome community diversity. Intra-species diversity was estimated

567    with polymorphism rate and strain count within each species at individual time points. When two

568    or more time points were available from the same person, delta polymorphism and gene content

569    variation (gain and loss) between time points were used to track DBD over time. Community

570    diversity was estimated with the Shannon index, species richness and rarefied richness (to 20

571    million reads per sample). When the relationship between the response variable (intra species

572    genetic diversity) and the predictor (community diversity) was approximately linear by visual

573    inspection, we fit generalized linear mixed models (GLMMs) (glmmTMB function from the

574    glmmTMB R package - RStudio version 1.2.5042) with community diversity as the predictor of

575    within-species genetic diversity, otherwise we fit Generalized additive mixed models (GAMs)

576    (mgcv function from the mgcv R package - RStudio version 1.2.5042) to account for the non-

577    linearity of the relationships.

578

579    To account for variation in sequencing depth, we added read count per sample (coverage) as a

580    covariate to all generalized mixed models except when richness was calculated on the rarefied

581    data. Species name and sample identifier nested within subject identifier were added as random

582    effects to account for variation between different species, subjects, and samples.

583

584    In generalized mixed models, the predictors were standardized to zero mean and unit variance

585    before analyses. We first assessed random effects significance by comparing nested models where

586    each random effect was dropped one at a time using the likelihood-ratio test (LRT, anova function

587    from the R stats package). We then assessed the fixed effects significance with LRTs implemented

588    in drop1 function in the stats package (this function drops individual terms from the full model

589    and report the AIC and the LRT *p*-value). We again used LRTs to compare the full significant

590    models to null models including all random effects but no fixed effects other than the intercept.

591    The difference in Akaike information criterion (ΔAIC) between full and null model and their

592    associated *p*-value are reported in Supplementary Tables S3, S4 and S5. As an additional

593    evaluation of the goodness of fits, we estimated the coefficient of determination ($R^2$) using the r2

594    function from the performance R package. Two values are reported: the marginal $R^2$, a measure of

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

596    explained by the entire model (Supplementary Table S5). We evaluated the GLMM fits by

597    inspecting the residuals using the DHARMa library in R (simulateResiduals and plot functions).

598    In generalized additive mixed models (GAMs), we evaluated the fits by inspecting residual

599    distributions and fitted-observed values plots using the gam.check function from the mgcv R

600    package. Adjusted $R^2$ (from summary function from the mgcv R package) values are reported as

601    a goodness of fits. All model outputs (summary function from mgcv and glmmTMB R packages)

602    are reported in the **Supplementary File 1**.

603

604    *Correlation analyses and scatter plots between community diversity and within-species genetic*

605    *diversity*

606    Only species present in at least four samples were retained to produce the scatter plots (ggplot

607    function in the ggplot2 R package) and to test the relationship between community diversity and

608    within-species genetic diversity with correlation analyses (Pearson when the relationship is linear

609    and Spearman otherwise; cor.test function from the stats R package).

610

611    *Community diversity is correlated with strain-level diversity*

612    To assess evidence for DBD in the gut microbiome, we first tested the relationship between
613    community diversity and within-species polymorphism rate. Because scatter plots (**Figs 2A,B,**
614    **S1,S2**) showed non-linear trends, we fitedt a separate generalized additive mixed model (GAM)
615    with polymorphism rate in a focal species as a function of each of the community diversity metrics
616    (Shannon index, species richness and rarefied species richness).

617

618    We then sought to test this relationship with community diversity calculated at higher taxonomic
619    ranks (from genus to phylum). We used GTDBK and the Genome Taxonomy Database (GTDB)
620    (Chaumeil et al., 2020) to annotate MIDAS reference genomes. Richness at each level was
621    estimated with the total number of distinct units in the sample. Shannon index was calculated based
622    on the relative abundances table from MIDAS (469 samples*5952 species). At each level and for
623    every distinct unit from the sample, we used the sum of the abundances of all species belonging to
624    the focal unit to calculate the Shannon index (using the diversity function from R vegan library).
625    We then fit two GAMs for each taxonomic rank (from genus to phylum) with Shannon diversity
626    and richness as the predictors of polymorphism rate in a focal species (with the coverage per
627    sample as a covariate and species name, sample and subject identifiers as random effects). All the
628    GAMs in this section were fitted with a beta error distribution with logit-link function because
629    polymorphism rate is a continuous value strictly bounded by 1, and all the terms were smooth
630    terms (See **Table S2** and **Supplementary File 1** for additional model details).

631

632    As a second test of DBD in the HMP data, we looked at the relationship between strain count in a
633    focal species and community diversity. Because scatter plots (**Figs 2C,D, S5,S6**) showed a linear
634    trend, we fit separate generalized linear mixed models (GLMMs) with strain count in a focal
635    species as a function of community diversity estimated with Shannon diversity, species richness,
636    or rarefied species richness. As strain number is positive count data, we compared many zero-
637    truncated count models based on the Akaike information criterion (AIC) score (AICtab function
638    from bbmle R library) (Brooks et al., 2017). We fit the model with the truncated negative binomial
639    distribution (truncated_nbinom2 in glmmTMB; the second best fit) in order to resolve the
640    overdispersion detected in the best fit (the truncated Poisson model) using the
641    check_overdispersion function from the performance R package as described here:
642    https://bbolker.github.io/mixedmodels-misc/glmmFAQ.html.

643

644    As in the previous section, we tested the relationship between strain count and community diversity

645    at higher taxonomic levels from genus to phylum, fitting a separate GLMM with strain diversity

646    in a focal species as a function of each metric of diversity (Shannon and richness) at higher

647    taxonomic levels. All GLMMs details are reported in **Table S4** and **Supplementary File 1**.

648

649    ***Genetic diversity as a function of community diversity over time***

650    To test DBD over time, we used HMP samples with multiple time points from the same person to

651    look at the relationship between polymorphism change (delta polymorphism) between two time

652    points and community diversity at the earlier time point. We fit Generalized additive mixed models

653    with delta polymorphism as a function of community diversity at the earlier time point, and added

654    the coverage per sample at the earlier time point as a covariate when diversity was not estimated

655    on rarefied data, as well as species name, sample and subject identifiers as random effects. We

656    used a Gaussian GAM since delta polymorphism is a continuous number that can take on negative

657    values (**Supplementary File 1**).

658

659    In addition, we investigated the effect of community diversity at one time point on gene variation

660    at the subsequent time point. We used separate negative binomial generalized linear mixed models

661    with gene gain as the response and each of the metrics of community diversity as the predictor

662    with the same covariates and random effects used in the previous models (**Supplementary File**

663    **1**). The same method was used to test how gene loss was related to community diversity (**Table**

664    **S5**, **Supplementary File 1**).

665

666    HMP longitudinal data were sampled at a time lag of ~6 months. To analyze time series at higher

667    resolution, we used longitudinal metagenomic data from a highly sampled healthy donor (host *am,*

668    sampled 206 times spanning 539 days between 2014-12-03 and 2016-05-25) (Poyet et al., 2019).

669    We tested the relationship between community diversity and genetic variation (polymorphism

670    change and gene content variation) in *B. vulgatus*. *B. vulgatus* is the most abundant species in all

671    *am* samples (mean coverage=58.46 and median=54.22). Community diversity was estimated with

672    richness and Shannon index calculated on rarefied data to 5 million reads per sample. We used a

673    Spearman correlation test (cor.test function from the stats R package) for the diversity-delta

26

674  polymorphism relationship (a nonlinear relationship) and Pearson correlations for both diversity-

675  gene loss and diversity-gene gain relationships (linear relationships) (**Figures 4** and **S11**).

676

677  ***Testing the Black Queen Hypothesis in HMP***

678  The negative relationship between gene loss in focal species and community diversity observed in

679  HMP and Poyet et al. (2019) data suggested the Black Queen Hypothesis (BQH) in the gut

680  microbiome. We sought to further test the BQH by comparing the content in genes and pathways

681  in a focal species to those present in the surrounding community. We used generalized additive

682  models (GAMs) to account for the non-linearity of the relationships (**Figures 5**, **S13**, **S14**). As in

683  all our models, we added the coverage per sample as a covariate as well as species name, sample,

684  and subject identifiers as random effects. Because both responses were count data, we compared

685  Poisson and negative binomial GAMs in both cases by looking at residual distribution and fitted-

686  observed values plots (gam.check function from the mgcv R package). We used a negative

687  binomial GAM for gene richness and a Poisson GAM for pathway richness, both with log-link

688  function. All the terms were specified as smooth terms, see **Table S6** and **Supplementary File 1**

689  for additional model details.

690

691

692  **Acknowledgements**

699

700  **References**

701  Albalat R, Cañestro C. 2016. Evolution by gene loss. *Nat Rev Genet*. doi:10.1038/nrg.2016.39

702  Almeida A, Mitchell AL, Boland M, Forster SC, Gloor GB, Tarkowska A, Lawley TD, Finn RD.

703      2019. A new genomic blueprint of the human gut microbiota. *Nature* **568**:499–504.

704      doi:10.1038/s41586-019-0965-1

705  Brooks ME, Kristensen K, Benthem KJ van, Magnusson A, Berg CW, Nielsen A, Skaug HJ,

706      Mächler M, Bolker BM. 2017. Modeling zero-inflated count data with glmmTMB. *BioRxiv*.

707  Calcagno V, Jarne P, Loreau M, Mouquet N, David P. 2017. Diversity spurs diversification in

708      ecological communities. *Nat Commun* **8**. doi:10.1038/ncomms15810

709  Chaumeil PA, Mussig AJ, Hugenholtz P, Parks DH. 2020. GTDB-Tk: A toolkit to classify

710      genomes with the genome taxonomy database. *Bioinformatics* **36**.

711      doi:10.1093/bioinformatics/btz848

712  Edgar RC. 2010. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*

713      **26**. doi:10.1093/bioinformatics/btq461

714  Estrela S, Diaz-Colunga J, Vila JCC, Sanchez-Gorostiaga A, Sanchez A. 2022. Diversity begets

715      diversity under microbial niche construction. *BioRxiv*.

716  Franzosa EA, McIver LJ, Rahnavard G, Thompson LR, Schirmer M, Weingart G, Lipson KS,

717      Knight R, Caporaso JG, Segata N, Huttenhower C. 2018. Species-level functional profiling

718      of metagenomes and metatranscriptomes. *Nat Methods* **15**:962–968. doi:10.1038/s41592-

719      018-0176-y

720  Garud NR, Good BH, Hallatschek O, Pollard KS. 2019. Evolutionary dynamics of bacteria in the

721      gut microbiome within and across hosts. *PLoS Biol* **17**:e3000102.

722      doi:10.1371/JOURNAL.PBIO.3000102

723  Garud NR, Pollard KS. 2020. Population Genetics in the Human Microbiome. *Trends Genet*

724      **36**:53–67. doi:10.1016/j.tig.2019.10.010

725  Good BH, Rosenfeld LB. 2022. Eco-evolutionary feedbacks in the human gut microbiome.

726      *bioRxiv*.

727  Groussin M, Poyet M, Sistiaga A, Kearney SM, Moniz K, Noel M, Hooker J, Gibbons SM,

728      Segurel L, Froment A, Mohamed RS, Fezeu A, Juimo VA, Lafosse S, Tabe FE, Girard C,

729      Iqaluk D, Nguyen LTT, Shapiro BJ, Lehtimäki J, Ruokolainen L, Kettunen PP, Vatanen T,

730      Sigwazi S, Mabulla A, Domínguez-Rodrigo M, Nartey YA, Agyei-Nkansah A, Duah A,

731      Awuku YA, Valles KA, Asibey SO, Afihene MY, Roberts LR, Plymoth A, Onyekwere CA,

732      Summons RE, Xavier RJ, Alm EJ. 2021. Elevated rates of horizontal gene transfer in the

733      industrialized human microbiome. *Cell* **184**. doi:10.1016/j.cell.2021.02.052

734  He Y, Zhou BJ, Deng GH, Jiang XT, Zhang H, Zhou HW. 2013. Comparison of microbial

735      diversity determined with the same variable tag sequence extracted from two different PCR

736 amplicons. *BMC Microbiol* **13**. doi:10.1186/1471-2180-13-208

737 Hibbing ME, Fuqua C, Parsek MR, Peterson SB. 2010. Bacterial competition: Surviving and

738  thriving in the microbial jungle. *Nat Rev Microbiol*. doi:10.1038/nrmicro2259

739 Human Microbiome Project Consortium T. 2012. A framework for human microbiome research

740  The Human Microbiome Project Consortium*. *Nature* **486**.

741 Jousset A, Eisenhauer N, Merker M, Mouquet N, Scheu S. 2016. High functional diversity

742  stimulates diversification in experimental microbial communities. *Sci Adv* **2**.

743  doi:10.1126/sciadv.1600124

744 Koskiniemi S, Sun S, Berg OG, Andersson DI. 2012. Selection-driven gene loss in bacteria.

745  *PLoS Genet* **8**. doi:10.1371/journal.pgen.1002787

746 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**.

747  doi:10.1038/nmeth.1923

748 Lenski RE. 2017. Experimental evolution and the dynamics of adaptation and genome evolution

749  in microbial populations. *ISME J*. doi:10.1038/ismej.2017.69

750 Lloyd-Price J, Mahurkar A, Rahnavard G, Crabtree J, Orvis J, Hall AB, Brady A, Creasy HH,

751  McCracken C, Giglio MG, McDonald D, Franzosa EA, Knight R, White O, Huttenhower C.

752  2017. Strains, functions and dynamics in the expanded Human Microbiome Project. *Nature*

753  **550**:61–66. doi:10.1038/nature23889

754 Madi N, Vos M, Murall CL, Legendre P, Shapiro BJ. 2020. Does diversity beget diversity in

755  microbiomes? *Elife* **9**. doi:10.7554/eLife.58999

756 McCutcheon JP, Moran NA. 2012. Extreme genome reduction in symbiotic bacteria. *Nat Rev*

757  *Microbiol*. doi:10.1038/nrmicro2670

758 Morris JJ, Lenski RE, Zinser ER. 2012. The black queen hypothesis: Evolution of dependencies

759  through adaptive gene loss. *MBio* **3**. doi:10.1128/mBio.00036-12

760 Morris JJ, Papoulis SE, Lenski RE. 2014. Coexistence of evolving bacteria stabilized by a shared

761  Black Queen function. *Evolution (N Y)* **68**. doi:10.1111/evo.12485

762 Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. 2016. An integrated metagenomics

763  pipeline for strain profiling reveals novel patterns of bacterial transmission and

764  biogeography. *Genome Res* **26**:1612–1625. doi:10.1101/gr.201863.115

765 Nayfach S, Shi ZJ, Seshadri R, Pollard KS, Kyrpides NC. 2019. New insights from uncultivated

766  genomes of the global human gut microbiome. *Nature* **568**:505–510. doi:10.1038/s41586-

767    019-1058-x

768    Nikoh N, Hosokawa T, Oshima K, Hattori M, Fukatsu T. 2011. Reductive evolution of bacterial

769        genome in insect gut environment. *Genome Biol Evol* **3**. doi:10.1093/gbe/evr064

770    Padfield D, Vujakovic A, Paterson S, Griffiths R, Buckling A, Hesse E. 2020. Evolution of

771        diversity explains the impact of pre-adaptation of a focal species on the structure of a

772        natural microbial community. *ISME J* **14**. doi:10.1038/s41396-020-00755-3

773    Poyet M, Groussin M, Gibbons SM, Avila-Pacheco J, Jiang X, Kearney SM, Perrotta AR, Berdy

774        B, Zhao S, Lieberman TD, Swanson PK, Smith M, Roesemann S, Alexander JE, Rich SA,

775        Livny J, Vlamakis H, Clish C, Bullock K, Deik A, Scott J, Pierce KA, Xavier RJ, Alm EJ.

776        2019. A library of human gut bacterial isolates paired with longitudinal multiomics data

777        enables mechanistic microbiome research. *Nat Med* **25**:1442–1452. doi:10.1038/s41591-

778        019-0559-3

779    Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin E V. 2014. Genomes in turmoil:

780        Quantification of genome dynamics in prokaryote supergenomes. *BMC Med* **12**.

781        doi:10.1186/s12915-014-0066-4

782    Reese AT, Dunn RR. 2018. Drivers of microbiome biodiversity: A review of general rules, feces,

783        and ignorance. *MBio* **9**. doi:10.1128/mBio.01294-18

784    San Roman M, Wagner A. 2021. Diversity begets diversity during community assembly until

785        ecological limits impose a diversity ceiling. *Mol Ecol* **30**. doi:10.1111/mec.16161

786    San Roman M, Wagner A. 2018. An enormous potential for niche construction through bacterial

787        cross-feeding in a homogeneous environment. *PLoS Comput Biol* **14**.

788        doi:10.1371/journal.pcbi.1006340

789    Schluter D, Pennell MW. 2017. Speciation gradients and the distribution of biodiversity. *Nature*.

790        doi:10.1038/nature22897

791    Sharma V, Rodionov DA, Leyn SA, Tran D, Iablokov SN, Ding H, Peterson DA, Osterman AL,

792        Peterson SN. 2019. B-Vitamin sharing promotes stability of gut microbial communities.

793        *Front Microbiol* **10**. doi:10.3389/fmicb.2019.01485

794    Simonsen AK. 2022. Environmental stress leads to genome streamlining in a widely distributed

795        species of soil bacteria. *ISME J* **16**. doi:10.1038/s41396-021-01082-x

796    Smillie CS, Sauk J, Gevers D, Friedman J, Sung J, Youngster I, Hohmann EL, Staley C, Khoruts

797        A, Sadowsky MJ, Allegretti JR, Smith MB, Xavier RJ, Alm EJ. 2018. Strain Tracking

798     Reveals the Determinants of Bacterial Engraftment in the Human Gut Following Fecal

799     Microbiota Transplantation. *Cell Host Microbe*. doi:10.1016/j.chom.2018.01.003

800 Vatanen T, Plichta DR, Somani J, Münch PC, Arthur TD, Hall AB, Rudolf S, Oakeley EJ, Ke X,

801     Young RA, Haiser HJ, Kolde R, Yassour M, Luopajärvi K, Siljander H, Virtanen SM,

802     Ilonen J, Uibo R, Tillmann V, Mokurov S, Dorshakova N, Porter JA, McHardy AC,

803     Lähdesmäki H, Vlamakis H, Huttenhower C, Knip M, Xavier RJ. 2019. Genomic variation

804     and strain-specific functional adaptation in the human gut microbiome during early life. *Nat*

805     *Microbiol* **4**. doi:10.1038/s41564-018-0321-5

806 Venturelli OS, Carr A V, Fisher G, Hsu RH, Lau R, Bowen BP, Hromada S, Northen T, Arkin

807     AP. 2018. Deciphering microbial interactions in synthetic human gut microbiome

808     communities. *Mol Syst Biol* **14**. doi:10.15252/msb.20178157

809 Walker AW, Duncan SH, Louis P, Flint HJ. 2014. Phylogeny, culturing, and metagenomics of

810     the human gut microbiota. *Trends Microbiol*. doi:10.1016/j.tim.2014.03.001

811 Wu D, Jospin G, Eisen JA. 2013. Systematic Identification of Gene Families for Use as

812     "Markers" for Phylogenetic and Phylogeny-Driven Ecological Studies of Bacteria and

813     Archaea and Their Major Subgroups. *PLoS One* **8**:e77033.

814     doi:10.1371/journal.pone.0077033

815 Zhao S, Lieberman TD, Poyet M, Kauffman KM, Gibbons SM, Groussin M, Xavier RJ, Alm EJ.

816     2019. Adaptive Evolution within Gut Microbiomes of Healthy People. *Cell Host Microbe*

817     **25**:656-667.e8. doi:10.1016/j.chom.2019.03.007

818