

RESEARCH ARTICLE

# Phylogeny-corrected identification of microbial gene families relevant to human gut colonization

Patrick H. Bradley<sup>1</sup>, Stephen Nayfach<sup>1,2</sup>, Katherine S. Pollard<sup>1,3,4\*</sup>

**1** Gladstone Institute of Data Science and Biotechnology, San Francisco, CA, USA, **2** Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA, **3** Department of Epidemiology & Biostatistics, Institute for Human Genetics, Quantitative Biology Institute and Institute for Computational Health Sciences, University of California, San Francisco, San Francisco, CA, USA, **4** Chan-Zuckerberg Biohub, San Francisco, CA, USA

\* [katherine.pollard@gladstone.ucsf.edu](mailto:katherine.pollard@gladstone.ucsf.edu)



**OPEN ACCESS**

**Citation:** Bradley PH, Nayfach S, Pollard KS (2018) Phylogeny-corrected identification of microbial gene families relevant to human gut colonization. *PLoS Comput Biol* 14(8): e1006242. <https://doi.org/10.1371/journal.pcbi.1006242>

**Editor:** Morgan Langille, DAL, CANADA

**Received:** November 23, 2017

**Accepted:** May 29, 2018

**Published:** August 9, 2018

**Copyright:** © 2018 Bradley et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** For this study, we used published data obtained from the NCBI SRA database; accession IDs are provided as [S3 Table](#). We also used the MIDAS v1.0 database ([http://lighthouse.ucsf.edu/MIDAS/midas\\_db\\_v1.0.tar.gz](http://lighthouse.ucsf.edu/MIDAS/midas_db_v1.0.tar.gz)).

**Funding:** Funding for this research was provided by NSF grants DMS-1069303 and DMS-1563159 ([www.nsf.gov](http://www.nsf.gov)), Gordon & Betty Moore Foundation grant #3300 ([moore.org](http://moore.org)), and institutional funds from the Gladstone Institutes ([gladstone.org](http://gladstone.org)). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## Abstract

The mechanisms by which different microbes colonize the healthy human gut versus other body sites, the gut in disease states, or other environments remain largely unknown. Identifying microbial genes influencing fitness in the gut could lead to new ways to engineer probiotics or disrupt pathogenesis. We approach this problem by measuring the statistical association between a species having a gene and the probability that the species is present in the gut microbiome. The challenge is that closely related species tend to be jointly present or absent in the microbiome and also share many genes, only a subset of which are involved in gut adaptation. We show that this phylogenetic correlation indeed leads to many false discoveries and propose phylogenetic linear regression as a powerful solution. To apply this method across the bacterial tree of life, where most species have not been experimentally phenotyped, we use metagenomes from hundreds of people to quantify each species' prevalence in and specificity for the gut microbiome. This analysis reveals thousands of genes potentially involved in adaptation to the gut across species, including many novel candidates as well as processes known to contribute to fitness of gut bacteria, such as acid tolerance in Bacteroidetes and sporulation in Firmicutes. We also find microbial genes associated with a preference for the gut over other body sites, which are significantly enriched for genes linked to fitness in an *in vivo* competition experiment. Finally, we identify gene families associated with higher prevalence in patients with Crohn's disease, including Proteobacterial genes involved in conjugation and fimbria regulation, processes previously linked to inflammation. These gene targets may represent new avenues for modulating host colonization and disease. Our strategy of combining metagenomics with phylogenetic modeling is general and can be used to identify genes associated with adaptation to any environment.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

Why do certain microbes and not others colonize our gut, and why do they differ between healthy and sick people? One explanation is the genes in their genomes. If we can find microbial genes involved in gut adaptation, we may be able to keep out pathogens and encourage the growth of beneficial microbes. One could look for genes that were present more often in prevalent microbes, and less often in rare ones. However, this ignores that related species are more likely to share an environment and also share many unrelated phenotypes simply because of common ancestry. To solve this problem, we used a method from ecology that accounts for phylogenetic relatedness. We first calculated gut prevalence for thousands of species using a compendium of shotgun sequencing data, then tested for genes associated with prevalence, adjusting for phylogenetic relationships. We found genes that are associated with overall gut prevalence, with a preference for the gut over other body sites, and with the gut in Crohn's disease versus health. Many of these findings have biological plausibility based on existing literature. We also showed agreement with the results of a previously published high-throughput screen of bacterial gene knockouts in mice. These results, and this type of analysis, may eventually lead to new strategies for maintaining gut health.

## Introduction

Microbes that colonize the human gastrointestinal (GI) tract have a wide variety of effects on their hosts, ranging from beneficial to harmful. Increasing evidence shows that commensal gut microbes are responsible for training and modulating the immune system [1, 2], protecting against inflammation [3] and pathogen invasion (reviewed in Sassone-Corsi and Raffatellu [4]), affecting GI motility [5], maintaining the intestinal barrier [6], and potentially even affecting mood [7]. In contrast, pathogens (and conditionally-pathogenic microbes, or “pathobionts”) can induce and worsen inflammation [8, 9], increase the risk of cancer in mouse models [10], and cause potentially life-threatening infections [11]. Additionally, the transplantation of microbes from a healthy host (fecal microbiota transplant, or FMT) is also a highly effective therapy for some gut infections [12], although it is still an active area of investigation why certain microbes from the donor persist long-term and others do not [13], and how pre-existing inflammatory disease affects FMT efficacy [14]. Which microbes are able to persist in the GI tract, and why some persist instead of others, is therefore a question with consequences that directly impact human health.

Because of this, we are interested in the specific mechanisms by which microbes colonize the gut, avoiding other potential fates such as being killed in the harsh stomach environment, simply passing through the GI tract transiently, or being outcompeted by other gut microbes. Understanding these mechanisms could yield opportunities to design better probiotics and to prevent invasion of the gut community by pathogens. In particular, creating new therapies, whether those are drugs, engineered bacterial strains, or rationally designed communities, will likely require an understanding of gut colonization at the level of individual microbial genes. We also anticipate that these mechanisms may vary in health vs. disease, since, for example, different selective pressures are known to be present in inflamed versus healthy guts [15, 16].

One approach that has been used to link genetic features to a phenotype is to correlate the two using observational data. Most typically, this approach is applied in the form of genome-wide association mapping, in which phenotypes are correlated with genetic markers across individuals in a population. While we are interested in comparing phenotypes and genetic

features *across*, rather than within species, the approach we take in this paper is conceptually similar. In order to perform association mapping, it is necessary to account for population structure, that is, dependencies resulting from common ancestry; otherwise, spurious discoveries can be made in genome-wide association studies [17]. Analogously, we expected it to be important to choose a method that can account for the confounding effect of phylogeny when testing for associations across species.

There is increasing interest in using phylogenetic information to make better inferences about associations between microbes and quantities of interest. For example, co-conservation patterns of genes (“correlogs”) have been used to assign functions to microbial genes [18], and genome-wide association studies have been applied within a genus of soil bacteria [19] as well as across strains of *Neisseria meningitidis* [20]. Recent publications have also described techniques that use information from the taxonomic tree to more accurately link clades in compositional taxonomic data to covariates [21, 22, 23]. However, so far, only one study has attempted to associate genes with a preference for the gut [24]. That study introduced a valuable method based on UniFrac and gene-count distances, which compares how well gut- vs. non-gut-associated microbes cluster on the species tree compared to a composite gene tree. This study also provides an important insight in the form of evidence of convergence of glycoside hydrolase and glycosyltransferase repertoires among gut bacteria, suggesting horizontal gene transfer within the gut community to deal with a common evolutionary pressure. The method described in that study, though, requires a binary phenotype of gut presence vs. absence. Deciding which microbes are “gut” vs. “non-gut” requires manual curation and can be somewhat subjective, as microbes have a continuous range of prevalences and can appear in multiple environments; this binarization could also potentially decrease power by excluding microbes with intermediate phenotypes. The method also requires multiple sequence alignments and trees to be built for every gene family under analysis, which are computationally intensive to generate over a large set of genomes.

We take a complementary approach and use a flexible technique, known as phylogenetic linear modeling, to detect associations between microbial genotype and phenotype while accounting for the fact that microbes are related to one another by vertical descent. Phylogenetic linear models have an extensive history in the ecology literature dating back to seminal works by Felsenstein [25] and Grafen [26]. However, despite their power, genome-scale applications of these models are still few in number [27] and, with the exception of one recent study that applied phylogenetic linear modeling to newly-sequenced isolate genomes from plant-associated microbial communities [28], have typically been used to relate traits of macroorganisms (e.g., anole lizards [29]) to their genotypes. While there is a growing appreciation for the need to explicitly account for phylogeny in microbial community analyses [27, 30, 31], we believe ours is the first study to directly apply this class of methods to metagenomic data.

This approach to accounting for phylogenetic relationships is general and could be applied to measure association of any quantitative phenotype with genotypes or other binary or quantitative characteristics. In this study, we focus on phenotypes related to the ability of bacteria to colonize the human gut: 1. overall prevalence in the guts of hosts from a specific population (e.g., post-industrialized countries), which we expect to capture ease of transmission, how cosmopolitan microbes are, and how efficiently they colonize the gut; 2. a preference for the gut over other human body sites in the same hosts, which we expect to capture gut colonization more specifically; and 3. a preference for the gut in disease (e.g., Crohn’s disease) versus health. We present a novel analytic pipeline in which we estimate these quantitative phenotypes for thousands of bacterial species directly from existing shotgun metagenomics data, both obviating the need for us to draw a cutoff between “gut” and “non-gut” microbes, and also giving us the necessary power to detect associations (S2 Fig). Coupling these phenotype estimates with

phylogenetic linear models, we generate a compendium of thousands of bacterial genes whose functions may be involved in colonizing the human gut.

## Results

We present a phylogeny-aware method for modeling associations between the presence of specific genes in bacterial genomes and quantitative phenotypes that measure how common these species are in the human microbiome. To apply phylogenetic linear modeling to the microbiome, we needed to solve three problems. First, we had to show that these models controlled false positives and had reasonable power on large bacterial phylogenies. Second, we needed to develop estimators that captured meaningful phenotypes related to bacterial colonization of humans for thousands of diverse bacterial species, most of which have never been studied in isolation, much less experimentally assayed for their abilities to colonize a mammalian body site. The third problem was to estimate genotypes (e.g., gene presence-absence) for each species. The analysis framework we describe is quite general and could be easily extended to link other phenotypes to genotypes across the tree of life.

### Phylogenetic linear models solve the problem of high false positive rates when testing for associations on bacterial phylogenies

To test for associations between quantitative phenotypes and binary genotypes across species, we use models with the following form:

$$\vec{\phi}_{x,E,D}(A) = \beta_{0,g} + \beta_{1,g} \vec{I}_g + \vec{\epsilon}_g \quad (1)$$

$\vec{\phi}_{x,E,D}(A)$  is a vector of quantitative phenotypes of interest, assessed in one environment  $e_x$  out of a set of possible environments  $E$ , normalizing out a set of study effects  $D$ , estimated from the dataset  $A$ . The elements of the vector  $\vec{\phi}_{x,E,D}(A)$  are  $\phi_{m,x,E,D}(A)$ , the phenotype value for microbe  $m$ .  $\beta_{0,g}$  is a baseline phenotype value,  $\beta_{1,g}$  is the effect of gene  $g$  on  $\vec{\phi}$ ,  $\vec{I}_g$  is a vector whose elements  $I_{m,g}$  are 0 if gene  $g$  is absent in species  $m$  and 1 if present, and  $\vec{\epsilon}_g$  is the remaining unmodeled variation in  $\vec{\phi}$ . We fit one model per gene  $g$ . The distribution of the residuals  $\vec{\epsilon}_g$  is the key difference between standard and phylogenetic linear models. In the standard model, the residuals are assumed to be independent and normally distributed. In the phylogenetic model, however, the residuals covary, with more closely-related species having greater covariance (see [Methods](#), “Modeling gene-phenotype associations” and “Fitting linear vs. phylogenetic models”; for a glossary of notation, see [S1 Appendix](#)).

To explore the potential pitfalls of failing to correct for phylogenetic structure in cross-species association tests, we generated a species tree for thousands of bacteria with genome sequences (see [Methods](#), “Phylogenetic-tree-construction”). In order to have a consistent operational definition of a microbial species, we used a set of previously defined bacterial taxonomic units with approximately 95% pairwise average nucleotide identity across the entire genome [32]. The methods we describe can be applied to other taxonomic levels or with other species definitions. Using this species tree, we performed simulations (see [Methods](#), “Worked example.”) for each of the four major bacterial phyla in the human gut (Bacteroidetes, Firmicutes, Proteobacteria, and Actinobacteria [33]). Specifically, we generated simulated phenotypes along the species tree, and then, for each phenotype, simulated a binary genotype for each species that covaried with the phenotype to varying degrees, including no association. We used levels of covariation spanning those we observed empirically between prevalence of species in gut metagenomes and presence-absence of genes (see below). (An effect size of 0.5

corresponds approximately to a 50% increase in prevalence, while an effect size of 1.0 corresponds approximately to a 72% increase in prevalence: see S8 Fig). These binary genotypes also had varying levels of overall phylogenetic signal (as measured using Ives-Garland  $\alpha$  [34]).

We then fit phylogenetic and standard linear models to the simulated data and tested for a relationship between each binary genotype and its corresponding continuous phenotype. For both standard and phylogenetic linear models, separate models were fit for each of the four phyla. The results were used to estimate false positive rate (Type I error) and power (1—Type II error) for the two methods across different effect sizes.

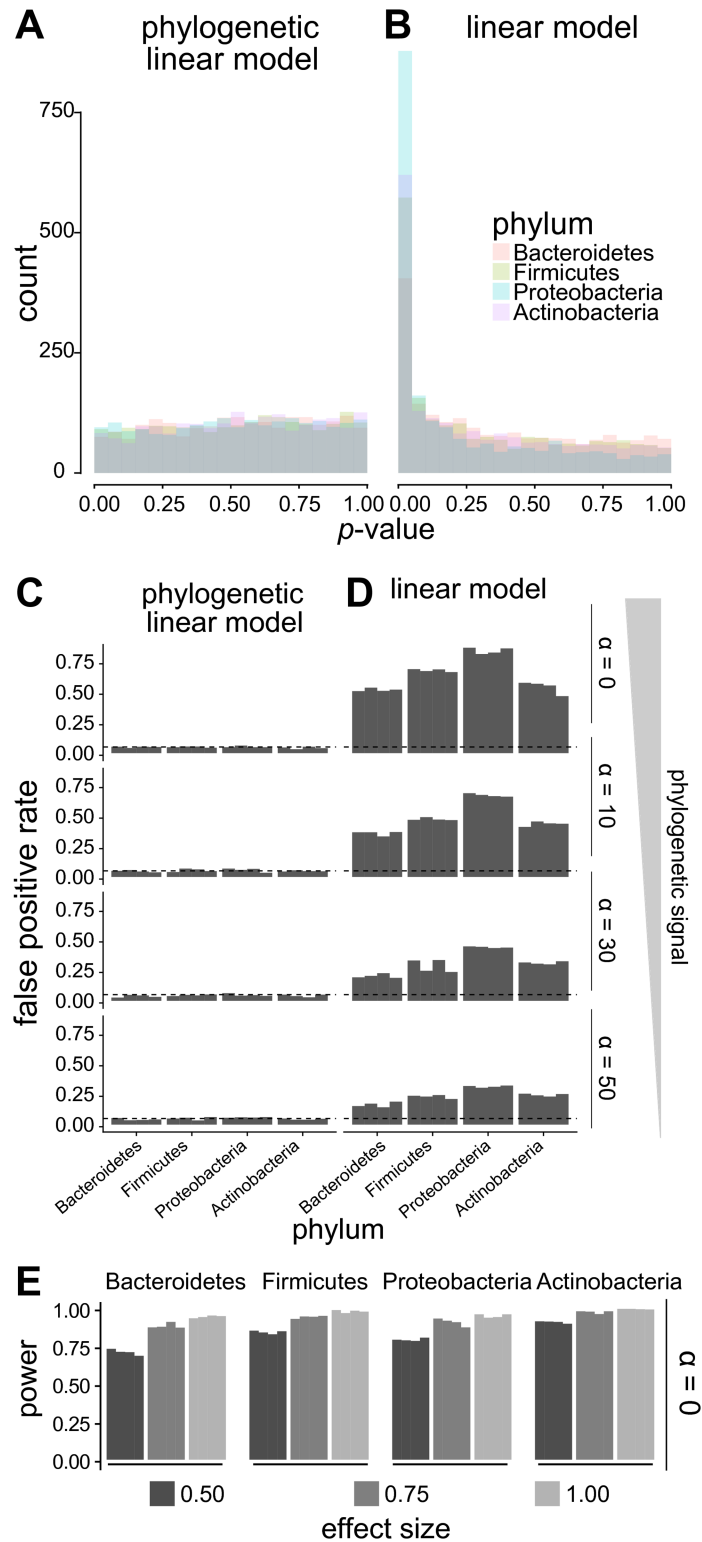
These analyses showed that standard linear models result in many false positive associations. When the binary genotype was specified to be wholly uncorrelated (i.e., under the null),  $p$ -values from the linear model showed a strong anticonservative bias (Fig 1B and 1D) with many more significant  $p$ -values than expected under no correlation. While lower levels of phylogenetic signal (larger Ives-Garland  $\alpha$ ) did result in less bias in the standard linear model, the average false positive rate ranged between 20% and 68% at  $p = 0.05$ . In contrast, the phylogenetic linear model  $p$ -value distribution was flat and Type I error was controlled appropriately (Fig 1A and 1C). This means that at the same  $p$ -value threshold, linear models will identify many spurious relationships compared to phylogenetic linear models. Further, our simulations with non-zero associations showed that the phylogenetic model has high power when applied to gut bacterial phyla, even for small effect sizes (Fig 1E; see Methods, “Power-analysis”). These results emphasize the importance of using models that account for phylogenetic relationships in cross-species association testing and demonstrate the feasibility of applying phylogenetic linear models to the human microbiome.

### Estimating quantitative phenotypes from shotgun data

To apply phylogenetic linear modeling to the microbiome we sought to define meaningful phenotypes for thousands of bacterial species, all of which have genome sequences but most of which have never been experimentally tested for, e.g., their abilities to grow on particular substrates or to colonize a model mammalian gut. We hypothesized that the prevalence and specificity of bacterial species in an environment, such as the human gut, should relate to their ability to colonize that environment and to how well adapted they are to persist there. These quantities can be thought of as phenotypes that can be estimated directly from shotgun metagenomics data. The precise taxonomic composition of a healthy gut microbiome can vary significantly from person to person [35], indicating that the ability of a microbe to colonize the gut is quantitative (and likely context-dependent, and stochastic). This phenotype can be conceptualized differently depending on which aspects of colonization one wishes to capture. We present metagenome-based estimators for two different types of colonization phenotypes. These are described in the context of our goal of studying the gut microbiome, but the approach is general and could be used to quantify how well a given genotype discriminates species found in or specific to any environment.

The first phenotype is the probability of observing a microbial species  $m$  in an environment  $e_x$ , that is, its overall prevalence  $P(m|e_x)$ . Both genes relating to survival in the GI tract and genes relating to survival, persistence, and dispersal in the outside environment are expected to correlate with overall prevalence. Prevalence can be estimated by the frequency with which the species is observed in a sample from the environment, for example, using a logit transform to enable linear modeling and pseudocounts to avoid estimates of 0 or 1 (see Methods, “Estimating the prevalence phenotype”).

The second type of quantitative phenotype is the *environmental specificity* of a microbial species, which we define as the conditional probability that a sample is derived from one



**Fig 1. Failing to account for tree structure results in an elevated false positive rate.** Continuous phenotypes and binary genotypes were simulated across the trees for the four phyla under consideration. A-D show results for the null of no true phenotype-genotype correlation. A-B) Histogram of  $p$ -values for simulated phenotypes and genotypes on the Bacteroidetes tree, using (A) phylogenetic or (B) standard linear models. The phylogenetic model distribution was similar to a uniform distribution, while the standard model was very anticonservative, having an excess of small  $p$ -

values. C-D) False positive rates (Type I error rates) at  $p = 0.05$  for the C) phylogenetic and D) standard models, across varying levels of true phylogenetic signal (Ives-Garland  $\alpha$ ). E) Traits with varying levels of “true” association spanning values we observed in real data were simulated, and power ( $y$ -axis) was computed using phylogenetic linear models.

<https://doi.org/10.1371/journal.pcbi.1006242.g001>

environment in a set of environments, given that the species is present in the sample. This parameter captures the power of a given microbe as a marker to discriminate between two or more different environments, such as different body sites or types of hosts (see [Methods](#), “Estimating environmental specificity scores”). This is distinct from its overall prevalence in the environment.

We developed an estimator for environmental specificity and applied it to two separate gut microbial phenotypes. First, we considered a phenotype defined as the conditional probability that a given body site is the gut and not another body site, given that a particular species is present. The physical distance between body sites is much smaller than the distance between hosts, and microbes from one body site are likely to be transiently introduced to others. Hence, enrichment of a species in one body site over others is stronger evidence for selection (versus dispersal) than is overall prevalence in that body site alone. We estimate this parameter with a *body-site specificity score* that uses metagenomics data to measure how predictive a particular microbe is for the gut versus other body sites (e.g., skin, urogenital tract, oropharynx, or lung).

The second type of environmental specificity we considered is the conditional probability that a host has a disease given that a particular species is present. This *disease-specific specificity score* is estimated in a similar way to the body-site specificity score (see [Methods](#), “Estimating environmental specificity scores”). We focus on Crohn’s disease, a type of inflammatory bowel disease known to be associated with dramatic shifts in the gut microbiota and in gut-immune interactions [36]. Genes associated with this disease-specific prevalence could illuminate differences in selective pressures between healthy vs. diseased gut environments. Both scores are based on maximum *a posteriori* (MAP) estimates of the conditional probability of a sample being from the gut given that a microbe is observed in the sample. To account for sampling noise, we use a shrunken estimate with a Laplace prior (see [Methods](#), “Estimating environmental specificity scores”).

## Genes associated with species prevalence in healthy human gut metagenomes

We assembled a compendium of published DNA sequencing data from healthy human stool microbiomes across five studies in North America, Europe, and China (426 subjects total). Using the MIDAS database and pipeline [32], we mapped metagenomic sequencing reads from each run to a panel of phylogenetic marker genes, and from these, estimated the relative abundances of species. Multiple runs corresponding to the same individual were averaged. We then estimated the prevalence (probability of non-zero abundance) of each species across these subjects, weighting each study equally and adding pseudocounts to avoid probabilities of exactly 0 or 1 (see [Methods](#), “Estimating the prevalence phenotype”). Finally, we determined whether genes (here, we take “genes” to mean members of a FIGfam protein family, which are designed to approximate “isofunctional homologs” [37]) were present or absent in the pangenomes of each species, based on sequenced genomes included in the MIDAS database, such that any FIGfam annotated in at least one sequenced isolate was considered to be present in the pangenome. This approach to genotyping could be extended to additionally include single-amplified genomes and metagenome-assembled genomes (see [Discussion](#)). Our analysis framework can also be applied to genotypes other than gene presence-absence

(e.g., nucleotide or amino acid changes). We note that while the FIGfam database does include many hypothetical protein families of unknown function, many bacterial genes lack even this level of annotation, so a more comprehensive grouping of genes into orthologous or functionally homologous groups could reveal yet more novel associations.

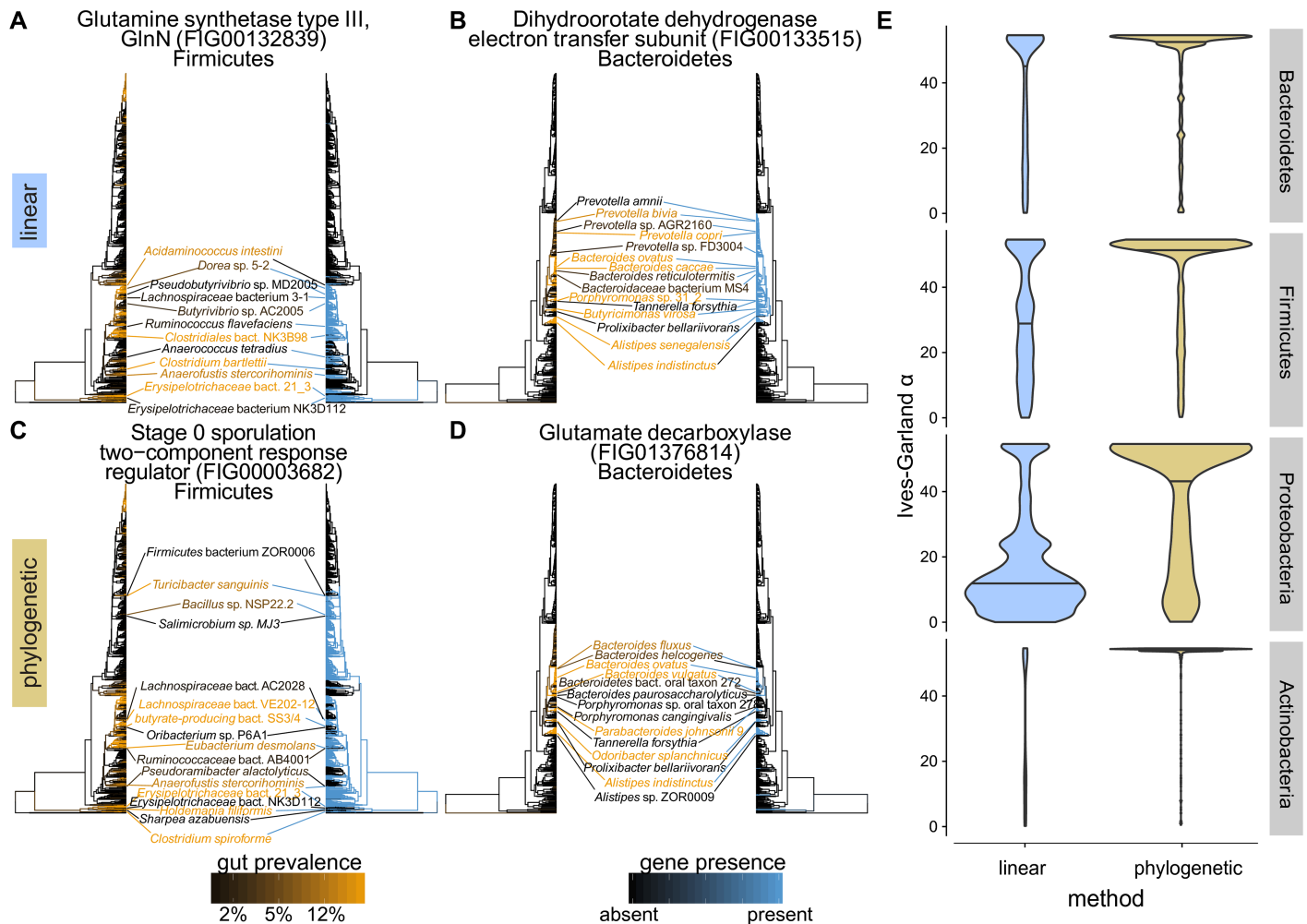
As expected, the most prevalent species overall included *Bacteroides vulgatus*, *Bacteroides ovatus*, and *Faecalibacterium prausnitzii*, while the least prevalent included halophiles and thermophiles (S1 Table). Gut prevalence had a strong phylogenetic signal (Pagel's  $\lambda = 0.97$ , likelihood-ratio  $p < 10^{-22}$ ), meaning that it was strongly correlated with the evolutionary relatedness of species. This emphasizes the need for phylogeny-aware modeling so that signal linking genes to prevalence will not be drowned out by shared variation in gene content between closely-related species.

To demonstrate the effect of phylogenetic correlation empirically, we fit both a standard linear model and a phylogenetic linear model for each of the four common gut phyla and all genes present in that phylum. These models relate logit-transformed estimates of the prevalence of different species in a phylum to a gene's presence-absence in those species' pangenomes. Recall that the residual variation in logit-prevalence is independent and normally distributed in the standard linear model, but has a distribution encoding correlations proportional to species relatedness in the phylogenetic linear model (see Methods, "Fitting linear vs. phylogenetic models"). For both standard and phylogenetic linear models, separate models were fit for each phylum. While this means that genes weakly-associated across the entire tree of life may have been missed by this approach, it has the advantage of both reducing the memory needed to store the gene presence-absence matrix and allowing for phylum-specific rates of evolution for our phenotype of interest. We modeled associations for 144,651 genes total across the four phyla, fitting 190,923 models total (since some genes are present in multiple phyla).

We used the parameter estimates and their standard errors from fitted models to test null hypotheses of the form  $H_0: \beta_{1,g} = 0$ , meaning gene  $g$  is not associated with gut prevalence of species in a particular phylum. The  $p$ -values were adjusted for multiple testing using the false discovery rate (FDR) (see Methods, "Fitting linear vs. phylogenetic models"). We found 9,455 FIGfam gene families positively associated with logit-prevalence within at least one phylum (FDR  $q \leq 0.05$ ) using phylogenetic linear models, 47% of which had no annotated function. We observed that 75% of the significant genes from these tests had effect sizes larger than (Bacteroidetes) 0.95, (Firmicutes) 1.09, (Proteobacteria) 0.35, and (Actinobacteria) 2.08, which are within the range of effect sizes for which phylogenetic linear models showed good performance in simulations (see above).

With standard linear models our tests identified 25,226 genes positively associated with gut prevalence, substantially more than with phylogenetic linear models (17.4% versus 6.5% of total). Based on our simulations, these likely included many false positives. The top results of phylogenetic versus standard linear models (Fig 2) illustrate the pitfalls of not correcting for phylogenetic correlation. Using the standard model, we recover associations such as those seen in Fig 2A and 2B: a subunit of dihydroorotate dehydrogenase in Bacteroidetes (Fig 2B) and in Firmicutes, a particular type of glutamine synthetase (Fig 2A). While these associations might look reasonable at a first glance, on closer inspection, they depend on the fact that these genes are near-uniformly present in entire clades of bacteria. These clades are, in general, more prevalent in the gut compared to the rest of the species in the tree. However, any finer structure relating to differences between close neighbors is lost. For example, dihydroorotate dehydrogenase (Fig 2B) is found not only in the human gut commensal *Bacteroides caccae*, but also in its relative *Bacteroides reticulotermitis*, which was not only low-prevalence in our samples but was indeed isolated from the gut of a subterranean termite [38].





**Fig 2. Examples of hits from standard linear (blue highlights) and phylogenetic (orange highlights) models.** In each panel, the tree on the left is colored by species prevalence (black to orange), while the tree on the right is colored by gene presence-absence (blue to black). Selected species are displayed in the middle; lines link species with the leaves to which they refer. The color of the line matches the color of the leaf. A-B) The standard model recovered hits that matched large clades but without recapitulating fine structure. C-D) The phylogenetic model recovered associations for which more of the fine structure was mirrored between the left-hand and right-hand trees, as exemplified by the species labeled in the middle. E) Violin plots of Ives-Garland  $\alpha$ , a summary of the rate of gain and loss of a binary trait across a tree, for genes significantly associated with prevalence in the standard (left, blue) and phylogenetic (right, orange) linear models. Horizontal lines mark the median of the distributions. The phylogenetic (orange) and standard linear (blue) models were significantly different for each phylum (Wilcox test for Bacteroidetes:  $8.2 \times 10^{-41}$ ; Firmicutes:  $7.6 \times 10^{-279}$ ; Proteobacteria:  $1.8 \times 10^{-235}$ ; Actinobacteria:  $9.0 \times 10^{-133}$ ).

<https://doi.org/10.1371/journal.pcbi.1006242.g002>

While this alone does not necessarily constitute evidence *against* these genes having adaptive functions in the human gut, we do expect that matched pheno- and genotypic differences between close phylogenetic neighbors offer stronger evidence for an association. An analogy can be drawn with genome-wide association mapping in humans: models that do not account for correlations between sites caused by population structure, as opposed to selective pressure, will tend to identify more spurious associations. In contrast, because the phylogenetic null model “expects” phenotypic correlations to scale with the evolutionary distance between species, this approach will tend to upweight cases where phylogenetically close relatives have different phenotypes and where distant relatives have similar phenotypes. This leads to the identification of candidate genes that capture more variation between close neighbors

(Fig 2C and 2D). Thus, phylogenetic linear models will identify genes whose presence in genomes is more frequently changing between sister taxa in association with a phenotype.

We provided further evidence that this trend is true in general by calculating the phylogenetic signal of the top hits from each model using Ives and Garland's  $\alpha$  [34]. This statistic captures the rate of transitions between having and not having a binary trait (here, a gene) across a tree; higher values therefore correspond to more disagreement between closely related species and lower values correspond to more agreement. Indeed, across all four phyla, the linear model identified gene families with significantly lower Ives-Garland  $\alpha$  than the phylogenetic model (Fig 2E, linear model  $p < 10^{-16}$ ), indicating that these genes' presence versus absence tended to be driven more by clade-to-clade differences (i.e., shared evolution).

These results suggest that standard linear models can identify genes that are truly important for colonizing an environment, such as the healthy human gut, but in addition will identify other genes that may simply be common in clades associated with that environment. The latter set will likely include many false positive associations from the perspective of understanding functions necessary for living in the environment. Phylogenetic linear models overcome this problem by adding the expectation that closely-related species will have similar phenotypes and distantly-related species will have less similar phenotypes, effectively upweighting instances where this is not the case. These conclusions are supported by our simulations and by an *in vivo* functional screen (see Results, "Deletion of gut-specific genes lowers fitness in the mouse microbiome").

### Gene families associated with gut prevalence provide insight into colonization biology

Several of the gene families that we observed to be associated with gut prevalence have previously been linked to gut colonization efficiency. For example, in Firmicutes, we noticed that several top hits were annotated as sporulation proteins (e.g., "Stage 0 sporulation two-component response regulator", Fig 2C). Sporulation is known to be a strategy for surviving harsh environments (such as acid, alcohol, and oxygen exposure) that is used by many, but not all, members of Firmicutes. Resistance to oxygen (aerotolerance) is particularly important because many gut Firmicutes are strict anaerobes [39], sporulation is known to be an important mechanism of transmission and survival in the environment (reviewed in Swick et al. [40]), and sporulation ability has been linked to transmission patterns of gut microbes [32]. Our result associating sporulation proteins to gut prevalence provides further evidence for sporulation as a strategy that is generally important for the propagation and fitness of gut microbes.

In Bacteroidetes, we observed an association between gut prevalence and the presence of a pair of gene families putatively assigned to the GAD operon, namely, the glutamate decarboxylase *gadB* and the glutamate/gamma-aminobutyric acid (GABA) antiporter *gadC*. These genes show a complex pattern of presence that is strongly correlated with gut prevalence (Fig 2D, S5 Fig). Results from research in Proteobacteria, where these genes were first described, shows that their products participate in acid tolerance. L-glutamate must be protonated in order to be decarboxylated to GABA; export of GABA coupled to import of fresh L-glutamate therefore allows the net export of protons, raising intracellular pH [41]. It was previously hypothesized that this acid tolerance mechanism allowed bacteria to survive the harshly acidic conditions in the stomach: indeed, if disrupted in the pathogen *Edwardsiella tarda*, gut colonization in a fish model is impaired [42]. *Listeria monocytogenes* with disrupted GAD systems also become sensitive to porcine gastric fluid [43]. However, while it has previously been shown that gut *Bacteroides* do contain homologs for at least one of these genes [41], their

functional importance has not yet been demonstrated in this phylum. Our results provide preliminary evidence that this system may be important in Bacteroidetes as well as in Proteobacteria.

### Using body sites as a control allows us to differentiate general dispersal from a specific gut advantage

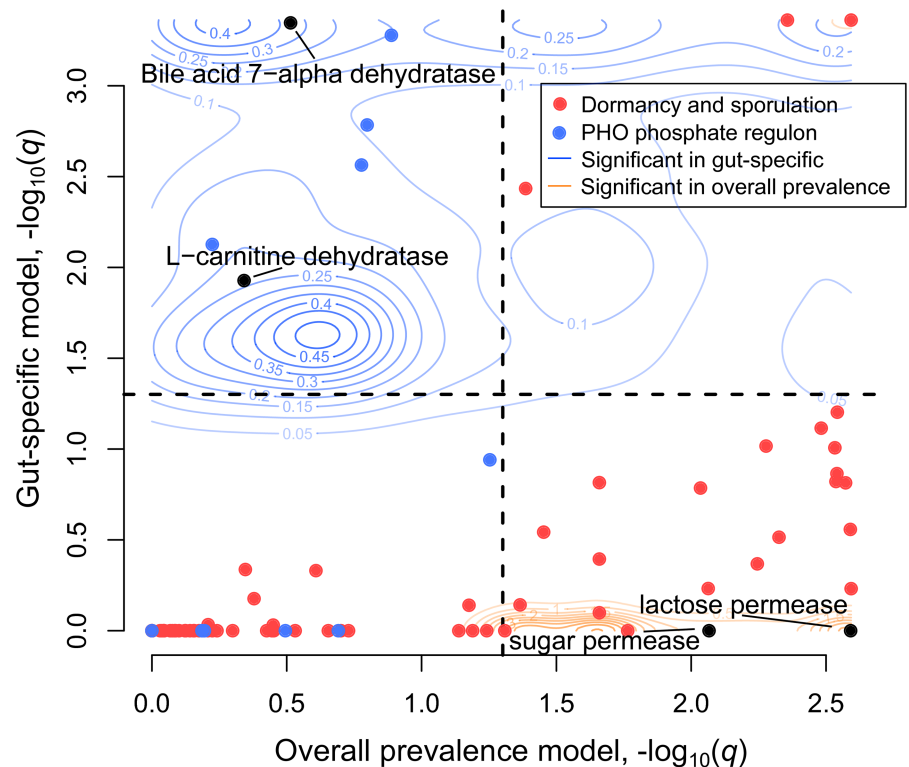
The previous analyses have focused on modeling the phenotype of overall prevalence in the human gut. However, microbes could be prevalent in the gut for at least two main reasons. First, they could be specifically well-adapted to the human gut; second, they could simply be very common in the environment (i.e., highly dispersed). The presence or absence of a gene family could enhance either of these properties. Some genes might, for example, confer improved stress tolerance that was adaptive across a range of harsh conditions, while others might allow, for example, uptake and catabolism of metabolic substrates that were more common in the human gut than in other environments.

With this in mind, we analyzed the relative enrichment of microbes in the gut over other human body sites in 127 individuals from the Human Microbiome Project (HMP) study [35]. We chose other body sites as a control because the physical distance between sites within a host is much smaller than the distance between people, and microbes from one body site are likely to be commonly, if transiently, introduced to other body sites (e.g., skin to oral cavity). To find specifically gut-associated genes, we used the phylogenetic linear model to regress gene presence-absence on the logit-transformed conditional probability  $P(e_{\text{Gut}}|m)$ , i.e., the probability that a body site was the human gut given that a particular species  $m$  was observed, which we estimated using Laplace regularization (see [Methods](#), “Estimating-environmental-specificity”). We identified 4,274 genes whose presence in bacterial genomes was positively associated with those species being present in the gut versus other body sites in at least one phylum (374 in Bacteroidetes, 1,515 in Firmicutes, 1,194 in Proteobacteria, and 1,255 in Actinobacteria).

Overall, the effect sizes for genes learned from this body site-specific model correlated only moderately with those learned from the “gut prevalence” models (median  $r = 0.26$ , range 0.14–0.35), indicating that these two quantitative phenotypes describe distinct phenomena. Additionally, the overlap between significant ( $q \leq 0.05$ ) hits for both models was usually small (median Jaccard index 0.056, range 0.035–0.343). These results are not surprising given that our regularized estimates of gut specificity were only moderately correlated with overall gut prevalence (Spearman’s  $\rho = -0.06$ , Pearson’s  $r = 0.31$ , [S1 Fig](#)), even when prevalence was calculated only from HMP gut samples (Spearman’s  $\rho = 0.05$ , Pearson’s  $r = 0.40$ ). This may arise from different genes being involved in dispersal or adaptation to many different environments versus those involved in adaptation specifically to the gut.

Indeed, when we compare the enrichments for genes significant in either the body site or overall prevalence models, we observe large functional shifts ([Fig 3](#)). For example, in the gut prevalence model, but not the body site-specific model, Firmicutes were strongly enriched for “dormancy and sporulation” ( $q = 4.4 \times 10^{-5}$ ). Because sporulation is likely useful in a wide range of environments beyond the gut, this result seems intuitive. Body site-specific results for Firmicutes were instead enriched for genes involved in “phosphorus metabolism” ( $q = 0.14$ ) and in particular the term “high affinity phosphate transporter and control of PHO regulon” ( $q = 0.058$ ).

We also observed biologically-justified individual gene families that were significant in the body site-specific model ( $q \leq 0.05$ ) but not the overall gut prevalence model ( $q > 0.5$  and/or wrong sign of effect size). In Firmicutes, for example, carnitine dehydratase and bile acid



**Fig 3. Comparison of results from the overall prevalence and body-site specific models for Firmicutes.** FDR-corrected significance (as  $-\log_{10}(q)$ ) of the overall model is plotted on the horizontal axis, whereas the same quantity for the body-site-specific model is plotted on the vertical axis. All FIGfams significant ( $q \leq 0.05$ ) in at least one of the two models are plotted as contour lines: FIGfams significant in the overall prevalence model (and possibly also the gut specific model) are plotted in orange, while FIGfams significant in the gut specific model (and possibly also the overall prevalence model) are plotted in blue. Selected SEED subsystems are displayed as colored points (legend), and selected individual genes are plotted as black points.

<https://doi.org/10.1371/journal.pcbi.1006242.g003>

7-alpha dehydratase were both significant only in the body site-specific model, suggesting a specific role for these genes within the gut environment. Indeed, bile acids are metabolites of cholesterol that are produced by vertebrates and thus unlikely to be encountered outside of the host. While the metabolite L-carnitine is made and used in organisms spanning the tree of life, it is particularly concentrated in animal tissue and especially red meat, and cannot be further catabolized by humans [44], making it available to intestinal microbes. Bile acid transformation by gut commensals is a well-established function of the gut microbiome, with complex influences on health (reviewed in Staley et al. [45]).

In Bacteroidetes, we found that a homolog of the autoinducer 2 aldolase *lsrF* was significant only in the body site-specific model. Autoinducer 2 is a small signaling molecule produced by a wide range of bacteria that is involved in interspecies quorum sensing. The protein *lsrF*, specifically, is part of an operon whose function in *Escherichia coli* is to “quench” or destroy the AI-2 signal [46]. Further, an increase of the AI-2 signal has been shown to decrease the Bacteroidetes/Firmicutes ratio *in vivo* in the intestines of streptomycin-treated mice [47]. Degrading this molecule is therefore a plausible gut-specific colonization strategy for gut Bacteroidetes. These discovered associations make the genes involved, including many genes without known functions or roles in gut biology, excellent candidates for understanding how bacteria adapt to the gut environment.

## Deletion of gut-specific genes lowers fitness in the mouse microbiome

Beyond finding evidence for the plausibility of individual genes based on the literature, we were interested in whether more high-throughput experimental evidence supported the associations we found between gut colonization and gene presence. To interrogate this, we used results from an *in vivo* transposon-insertion screen of four strains of *Bacteroides*. This screen identified many genes whose disruption caused a competitive disadvantage in gnotobiotic mice, as revealed by time-course high-throughput sequencing; 79 gene families significantly affected microbial fitness across all four strains tested [48]. Determining agreement with this screen is somewhat complicated by the fact that we associated gene presence to gut specificity across all members of the phylum Bacteroidetes, and not only within the *Bacteroides* genus. Significance of overlap therefore depends on what we take as the null “background” set, the cutoff used for significance, and the set of results from the screen we choose as true positives (S3 Table).

Despite these complications, this analysis clearly showed that the 79 genes whose disruption led to lower fitness in the murine gut across all four *Bacteroides* species were over-represented among our predictions for gut-specific genes (odds ratio = 4.67,  $q = 6.7 \times 10^{-3}$ ), and remained so if we only considered the gene families that were present in all *Bacteroides* species (odds ratio = 7.02,  $q = 3.4 \times 10^{-3}$ ) (Table 1). Interestingly, we observed the opposite pattern for the overall prevalence model: the prevalence-associated genes we identified were actually depleted for genes found to be important *in vivo* (odds ratio = 0.18,  $q = 1.1 \times 10^{-2}$ ). We believe that this

**Table 1. Assessment of agreement between the *in vivo* results from Wu et al. [48] and gut-specific (“bodysite”) vs. gut prevalence (“overall”) phylogenetic models.** The background sets for enrichment tests were defined as follows: “all tested” (all gene families for which a phylogenetic model was fit), “Bacteroides (core or variable)” (all gene families with at least one representative in *Bacteroides* genome cluster pangenomes), “Bacteroides (core only)” (gene families that were present in all *Bacteroides* genome cluster pangenomes), “Bacteroides (variable only)” (gene families present in some but not all *Bacteroides* genomes clusters), and “Bacteroides thetaiotaomicron only” (only gene families present in *Bacteroides thetaiotaomicron*). The “all tested” background sets are further classified as “all tested (overall)” and “all tested (body site)” since the set of genes tested differed slightly. The *p*-values are from Fisher’s exact tests. These comparisons have been excerpted from the full set, which can be seen in S3 Table; *q*-values were calculated based on this full set of tests using the Benjamini-Hochberg method [89].

Background set	FDR	MODEL	<i>p</i> -value	odds ratio	<i>q</i> -value	significant
<b>All tested</b>	<b>5%</b>	<b>overall</b>	$4.86 \times 10^{-3}$	<b>0.18</b>	$1.12 \times 10^{-2}$	<b>TRUE</b>
<b>Bacteroides (core or variable)</b>	<b>5%</b>	<b>overall</b>	$5.32 \times 10^{-12}$	<b>0.05</b>	$5.32 \times 10^{-11}$	<b>TRUE</b>
Bacteroides (core only)	5%	overall	1.00	0.00	1.00	FALSE
<b>Bacteroides (variable only)</b>	<b>5%</b>	<b>overall</b>	$8.03 \times 10^{-4}$	<b>0.14</b>	$2.54 \times 10^{-3}$	<b>TRUE</b>
<b>Bacteroides thetaiotaomicron only</b>	<b>5%</b>	<b>overall</b>	$4.89 \times 10^{-6}$	<b>0.10</b>	$2.10 \times 10^{-5}$	<b>TRUE</b>
<b>All tested</b>	<b>25%</b>	<b>overall</b>	$2.30 \times 10^{-3}$	<b>0.25</b>	$6.28 \times 10^{-3}$	<b>TRUE</b>
<b>Bacteroides (core or variable)</b>	<b>25%</b>	<b>overall</b>	$2.85 \times 10^{-15}$	<b>0.06</b>	$3.43 \times 10^{-14}$	<b>TRUE</b>
Bacteroides (core only)	25%	overall	$2.49 \times 10^{-1}$	0.00	$3.65 \times 10^{-1}$	FALSE
<b>Bacteroides (variable only)</b>	<b>25%</b>	<b>overall</b>	$4.96 \times 10^{-4}$	<b>0.19</b>	$1.65 \times 10^{-3}$	<b>TRUE</b>
<b>Bacteroides thetaiotaomicron only</b>	<b>25%</b>	<b>overall</b>	$4.20 \times 10^{-8}$	<b>0.12</b>	$2.80 \times 10^{-7}$	<b>TRUE</b>
<b>All tested</b>	<b>5%</b>	<b>bodysite</b>	$2.67 \times 10^{-3}$	<b>4.67</b>	$6.68 \times 10^{-3}$	<b>TRUE</b>
Bacteroides (core or variable)	5%	bodysite	$6.43 \times 10^{-2}$	2.23	$1.20 \times 10^{-1}$	FALSE
<b>Bacteroides (core only)</b>	<b>5%</b>	<b>bodysite</b>	$1.14 \times 10^{-3}$	<b>7.02</b>	$3.43 \times 10^{-3}$	<b>TRUE</b>
Bacteroides (variable only)	5%	bodysite	1.00	0.00	1.00	FALSE
Bacteroides thetaiotaomicron only	5%	bodysite	$1.67 \times 10^{-1}$	1.79	$2.69 \times 10^{-1}$	FALSE
<b>All tested</b>	<b>25%</b>	<b>bodysite</b>	$2.67 \times 10^{-3}$	<b>3.36</b>	$6.68 \times 10^{-3}$	<b>TRUE</b>
Bacteroides (core or variable)	25%	bodysite	$1.80 \times 10^{-1}$	1.69	$2.77 \times 10^{-1}$	FALSE
<b>Bacteroides (core only)</b>	<b>25%</b>	<b>bodysite</b>	$1.09 \times 10^{-2}$	<b>3.47</b>	$2.41 \times 10^{-2}$	<b>TRUE</b>
Bacteroides (variable only)	25%	bodysite	$7.14 \times 10^{-1}$	1.06	$8.40 \times 10^{-1}$	FALSE
Bacteroides thetaiotaomicron only	25%	bodysite	$5.56 \times 10^{-1}$	1.23	$7.41 \times 10^{-1}$	FALSE

<https://doi.org/10.1371/journal.pcbi.1006242.t001>

is because the body-site-specific model, like the experiment, focused specifically on colonization efficiency, while the overall gut prevalence model would have included genes involved in persistence and dispersal in the environment and transfer between hosts. This experimental evidence supports the idea that environment-specific phylogenetic linear models truly identify genes that are important for bacteria to colonize an environment.

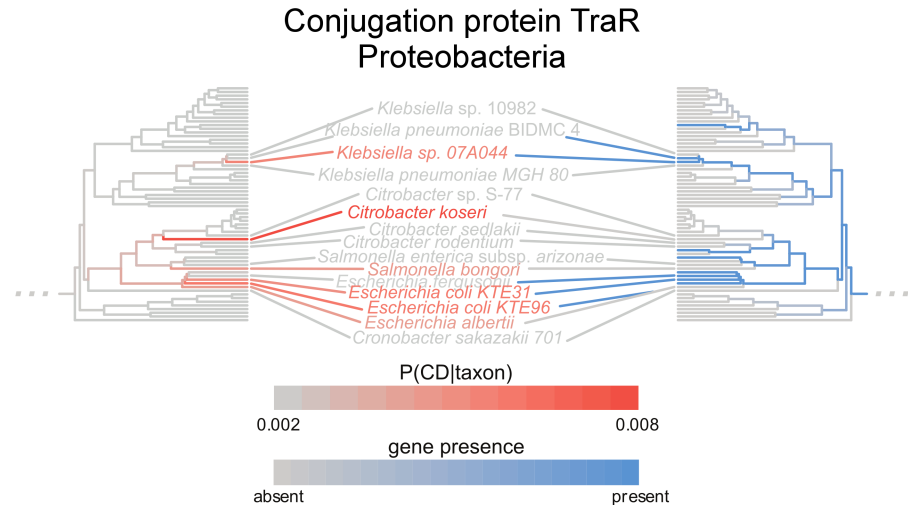
## We identify Proteobacterial gene families associated with microbes that are more prevalent in Crohn's disease

The above analyses were performed with respect to the gut of healthy individuals from the mainly post-industrial populations of North America, Europe and China. However, we also know that taxonomic shifts are common between healthy guts versus the guts of individuals from the same population with diseases such as type 2 diabetes, colorectal cancer, rheumatoid arthritis, and inflammatory bowel disease (reviewed in Wang et al. [49]). One explanation for these results is that sick hosts select for specific microbial taxa, as with the links previously observed between Proteobacteria and the inflammation that accompanies many disease states [50]. Since gut microbes have also been implicated in altering disease progression (reviewed in Lynch and Pedersen [51]), identifying genes associated with colonizing diseased individuals may afford us new opportunities for intervention.

To identify microbiome functions that could be involved in disease-specific adaptation to the gut, we looked for genes that were present more often in microbes that discriminated case from control subjects. Specifically, we compared  $n = 38$  healthy controls from the MetaHIT consortium to  $n = 13$  individuals with Crohn's disease [52, 53]. Similar to our analysis of gut versus other body sites, we used the conditional probability that a subject had Crohn's disease given that we observed a particular microbe in their gut microbiome  $P(e_{CD}|M)$  (see Methods). We identified 1,841 genes whose presence in bacterial genomes was associated with Crohn's after correcting for phylogenetic relationships in at least one phylum (796 in Bacteroidetes, 278 in Firmicutes, 467 in Proteobacteria, and 340 in Actinobacteria).

Eleven of our top Proteobacterial associations were annotated as fimbrial proteins, including one predicted to be involved specifically in the biosynthesis of type 1 fimbriae, or pili (FimI, association  $q = 2.8 \times 10^{-3}$ ), cell surface structures involved in attachment and invasion. Crohn's pathology has been linked to an immune response to invasive bacteria, and adherent-invasive *E. coli* (AIEC) appear to be overrepresented in ileal Crohn's [54]. In an AIEC *E. coli* strain isolated from the ileum of a Crohn's patient, type 1 pili were required for this adherent-invasive phenotype [55]. Chronic infection by AIEC strains was also observed to lead to chronic inflammation, and to an increase in Th17 cells and a decrease in CD8<sup>+</sup> T cells similar to that observed in Crohn's patients [56].

An additional striking feature of the results was the number of Proteobacterial proteins associated with greater risk of Crohn's that were annotated as being involved in type IV secretion and conjugative transfer systems (Fisher's test  $q = 4.9 \times 10^{-4}$ ). Conjugative transfer is a process by which gram-negative bacteria in direct physical contact share genetic material. More specifically, many of these genes, including one annotated as TraR (Fig 4), were homologs of those involved in an "F-type" conjugal system for transferring IncF plasmids, which can be classified as a variety of type IV secretion system [57]. Previously, in a mouse model, gut inflammation was shown to stimulate efficient horizontal gene transfer in Proteobacteria by promoting blooms of *Enterobacteriaceae* and thus facilitating cell-to-cell contact [58]. Future work will be required to determine whether this increased conjugation is a neutral consequence of inflammation, a causative factor, or provides a selective advantage in the inflamed gut.



**Fig 4. Genes involved in conjugative transfer are associated with Crohn’s disease-enriched species.** The conjugation transcriptional regulator *traR* is plotted as an example. The left-hand tree is colored by each species’ disease specificity score, i.e., the conditional probability of Crohn’s given the observation of a given species (grey, which represents the prior, to red, which represents a higher conditional probability). The right-hand tree is colored by gene presence-absence (grey, meaning absent, or blue, meaning present). The mirrored patterns drive the phylogeny-corrected correlation.

<https://doi.org/10.1371/journal.pcbi.1006242.g004>

## Discussion

The present analyses represent a first look into what can be learned by combining shotgun metagenomics with phylogenetically-aware models. Several extensions to our work could be made in the future. First, in addition to modeling prevalence, for instance, we could model abundance using a phylogenetic linear model with random effects [59], potentially allowing us to learn what controls the steady-state abundance of species in the gut. Additionally, we could also use these models to screen for epistatic interactions, which would be near-intractable even in systems with well-characterized genetic tools, but for which a subset of hypotheses could be validated by, e.g., comparing the fitness of wild-type microbes with double knockouts. While controlling the total number of tests would still be important to preserve power, an automated, computational approach to detecting gene interactions would still offer important savings in time and expense over developing a genome-wide experimental library of multiple knockouts per organism under investigation.

Currently, these analyses estimate species abundance and gene presence-absence from available sequenced isolate genomes. However, it has been estimated that on average 51% of genomes in the gut are from novel species [32]. Especially for case/control comparisons, using information from metagenomic assemblies could enable quantification of species with no sequenced representatives, and would yield a more accurate estimate of the complement of genes in the pangenome for species that do have sequenced representatives. This would be particularly helpful in gut communities from individuals in non-industrialized societies that are enriched for novel microbial species [32]. In fact, genes then could be treated as quantitative variables (e.g., coverage or prevalence) rather than binary, which is possible for covariates in phylogenetic linear models and simply changes the interpretation of the association coefficient  $\beta_{1,g}$ .

Another potential extension would be to model prevalence and environment-specific prevalence for taxa other than the species clusters analyzed in this study. We focused on four prevalent and abundant phyla of bacteria, but our methods could be applied more broadly as long

as quantitative phenotypes and genotypes could be accurately estimated. Phylogenetic linear modeling could also be applied directly to genera or higher taxonomic groups, although both phenotypes and genotypes would be averaged over more diverse sets of genomes, which could result in associations with different signs canceling out. As more genome and metagenome data is generated for microbial populations over time, extensions of phylogenetic linear modeling (e.g., with random effects [59]) may also be useful for studying associations between phenotypes and evolving gene copy number and single nucleotide variants at the strain level. This application would require accurate trees with strains as leaves, each with estimates of a phenotype and genotype. Additionally, our current definition of species approximates a 95% average nucleotide identity (ANI) cutoff; while this approach is a standard bioinformatic approach [60], and appears to be a “natural boundary” in analyses of genome compendia [61], the precise definition of a bacterial species remains a matter of active debate, and in the future may include phenotypic information [62] or information about gene flow [63]. Beyond prevalence, other phenotypes will also be interesting to investigate, especially experimentally measured phenotypes from high throughput screens and other techniques that complement genomics.

In summary, using phylogenetic linear models, we were able to discover thousands of specific gene families associated with quantitative phenotypes calculated directly from data: overall gut prevalence, a specificity score for the gut over other body sites, and a specificity score for the gut in Crohn’s disease versus health. Importantly, we have shown through simulation and real examples that standard linear models are inadequate for this task because of an unacceptably high false-positive rate under realistic conditions. Furthermore, many of the results we found also have biological plausibility, both from the literature on specific microbial pathways and from a high-throughput *in vivo* screen directly measuring colonization efficiency. In addition to these expected discoveries, we also found thousands of novel candidates for understanding and potentially manipulating gut colonization. These results illustrate the potential of integrating phylogeny with shotgun metagenomic data to deepen our understanding of the factors determining which microbes come to constitute our gut microbiota in health and disease.

## Methods

A graphical overview of our statistical methods can be found in [S2 Fig](#). Additionally, we provide a glossary of mathematical notation used in this section in [S1 Appendix](#).

### Species definition

We utilized the previously published clustering of 31,007 high-quality bacterial genomes into 5,952 species from the MIDAS 1.0 database [32] ([http://lighthouse.ucsf.edu/MIDAS/midas\\_db\\_v1.0.tar.gz](http://lighthouse.ucsf.edu/MIDAS/midas_db_v1.0.tar.gz)). These species clusters are sets of genomes with high pairwise sequence similarity across a panel of 30 universal, single-copy genes. The genomes in each species clustering have approximately 95% average genome-wide nucleotide identity, a common “gold-standard” definition of bacterial and archaeal species [64]. These species-level taxonomic units are similar to, but can differ from, operational taxonomic units (OTUs) defined solely on the basis of the 16S rRNA gene.

Taxonomic annotations for each species were drawn from the MIDAS 1.0 database. Some taxonomic annotations of species in the MIDAS database were incomplete; these were fixed by searching the NCBI Taxonomy database using their web API via the `rentrez` package [65] and retrieving the full set of taxonomic annotations.



## Pangenomes

Pangenomes for all species used in this study were downloaded from the MIDAS 1.0 database. As previously described [32], pangenomes were constructed by clustering the DNA sequences of the genes found across all strains of each species at 95% sequence identity using UCLUST [66]. Pangenomes were functionally annotated based on the FIGfams [37] which were included in the MIDAS databases and originally obtained from the PATRIC [67] database. Thus, each pangenome represents the set of known, non-redundant genes from each bacterial species with at least one sequenced isolate.

## Phylogenetic tree construction

The tree used for phylogenetic analyses was based on the tree from Nayfach et al. [32] based on an approximate maximum likelihood using FastTree 2 [68] on a concatenated alignment (using MUSCLE [69]) of thirty universal genes. Thus, each tip in the tree represents the phylogenetic placement for one bacterial species. For the current analyses, the tree was rooted using the cyanobacterium *Prochlorococcus marinus* as an outgroup, and the tree was then divided by phylum, retaining the four most prevalent phyla in the human gut (Bacteroidetes, Firmicutes, Actinobacteria, and Proteobacteria). One Actinobacterial species cluster, the radiation-resistant bacterium *Kineococcus radiotolerans*, was dropped from the tree because it had an extremely long branch length, indicating an unusual degree of divergence. Finally, phylum-specific trees were made ultrametric using the `chronos` function in the R package `ape` [70], assuming the default “correlated rates” model of substitution rate variation. We performed this step because first, our taxa were contemporaneously sampled, and second, we assumed that our phenotypes of interest varied with divergence time, as opposed to the number of substitutions per site separating marker gene sequences [71].

## Estimating species abundance across human associated metagenomes

Metagenome samples were drawn from subjects in the Human Microbiome Project (HMP) [35], the MetaHIT consortium [52, 53], a study of glucose control [72], and a study of type 2 diabetes [73]. Accession numbers were identified using the aid of SRADB [74] and downloaded from the Sequence Read Archive (SRA) [75]. The relative abundance of bacterial species in the metagenomes was estimated using MIDAS v1.0 [32], which maps reads to a panel of 15 phylogenetic marker genes. Species relative abundances are computed as previously described [32] (“Species abundance estimation”): essentially, they are normalized counts of reads mapping to bacterial species, with non-uniquely mapped reads assigned probabilistically.

Accession IDs used can be found in S4 Table. For prevalence estimates, we used healthy subjects from all four cohorts; for body site comparisons, we used only healthy subjects from HMP [35], and for the Crohn’s case-control comparison, we used only subjects from the MetaHIT consortium [52, 53].

## Modeling gene-phenotype associations

The basic design is the same for all models that we fit: we model the effect of a categorical variable, gene (specifically, FIGfam family) presence vs. absence, on a particular phenotype estimated for many microbes from data.

Here, let  $\vec{\phi}_{x,E[,D]}(A)$  refer to a vector whose elements  $\phi_{m,x,E[,D]}(A)$  refer to an estimate of the phenotype  $\phi$  for a microbe  $m$ , in an environment  $e_x$  from a set of  $k$  environments  $E = \{e_1, \dots, e_k\}$ , optionally also adjusting for potential dataset effects  $D$ , based on a matrix of microbial presence-absence data  $A$ . We then model the effect on this phenotype of having vs.

lacking each particular gene  $g$ , fitting one model per gene:

$$\vec{\phi}_{x,E[D]}(A) = \beta_0 + \beta_{1,g} \vec{I}_g + \vec{\epsilon}_g$$

where  $\beta_0$  is a baseline intercept value,  $\beta_{1,g}$  is the effect size of gene  $g$ ,  $\vec{I}_g$  is a binary vector whose elements  $I_{g,m}$  are 1 when microbe  $m$ 's pangenome contains the gene  $g$  and 0 otherwise, and  $\vec{\epsilon}_g$  are the residuals. We then test the null hypothesis  $H_0: \beta_{1,g} = 0$ , yielding one  $p$ -value per gene; the resulting genewise  $p$ -values are finally corrected for multiple testing using an adaptive false discovery rate approach ( $q$ -value estimation).

The differences in the models we fit concern only how we obtain phenotype estimates  $\vec{\phi}_{x,E[D]}(A)$ , and our assumptions about how the residuals  $\vec{\epsilon}_g$  are distributed.

### Fitting linear vs. phylogenetic models

The phylogenetic and standard linear models are very similar, except for the assumptions about the distribution of the residuals. In the standard linear model, the residuals are assumed to be independently and identically distributed as a normal distribution, i.e.,  $\epsilon_{g,m} \sim \mathcal{N}(0, \sigma^2)$  or using multivariate notation  $\vec{\epsilon}_g \sim \mathcal{N}(0, \sigma^2 I)$ . In the phylogenetic model, in contrast, the residuals are not independent: rather, they are correlated based on the phylogenetic relatedness of the species. They are therefore distributed  $\epsilon \sim \mathcal{N}(0, \Sigma)$ , with the following covariance matrix:

$$\Sigma = \begin{bmatrix} \sigma^2 & \sigma_{1,2} & \cdots & \sigma_{1,i} \\ \sigma_{2,1} & \sigma^2 & & \\ \vdots & & \ddots & \\ \sigma_{i,1} & & & \sigma^2 \end{bmatrix}$$

where  $i$  is the number of species,  $\sigma^2$  is the overall variance, and  $\sigma_{1,2}$  is the covariance between species 1 and species 2. Under the assumption of the phylogenetic model (evolution of a continuous phenotype according to Brownian motion), this covariance is proportional to the distance between the last common ancestor of species 1 and 2 and the root of the tree. Thus, very closely-related species have a common ancestor that is far from the root, while the last common ancestor of two unrelated species is the root node itself. This method was first described in Grafen [26]; for this study, we use the implementation in the phylolm R package [76].

$\beta_1$  parameters were tested for a significant difference from 0 and the resulting  $p$ -values were converted to  $q$ -values using Storey and Tibshirani's FDR correction procedure [77, 78].

### Metagenomic presence-absence data

We use binary presence-absence data to calculate the phenotypes of interest. More formally, we conceptualize the metagenomic data as a matrix  $A$  of microbial presence-absence with dimensions  $i \times j$ , where  $i$  is the number of microbes and  $j$  is the number of samples, and  $a_{m,n}$  is 1 if the relative abundance of microbe  $m$  (calculated using MIDAS's taxonomic profiling [32])

is greater than 0, and 0 otherwise:

$$A = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,j} \\ a_{2,1} & a_{2,2} & & a_{2,j} \\ \vdots & & \ddots & \vdots \\ a_{i,1} & a_{i,2} & \cdots & a_{i,j} \end{bmatrix}$$

We conceptualize each  $e_1, e_2, \dots, e_k \in E$  as a set of indices, referring to samples collected from that environment, e.g., the oropharynx in healthy subjects, or the gut in subjects with Crohn’s disease, such that for all  $e_x \in E, e_x \subseteq \{1, 2, \dots, j\}$ . Because one environment may be tested in multiple studies, for our prevalence estimates, we also define a similar mapping of samples to datasets  $d_1, d_2, \dots, d_l \in D$  such that for all  $d_y \in D, d_y \subseteq \{1, 2, \dots, j\}$ . (For calculating environmental specificity scores, to avoid having to correct for unbalanced designs, we only use single datasets that measured all environments to be compared.) We also assume that  $E$  and  $D$  are partitions of  $\{1, 2, \dots, j\}$ , such that every sample is covered and no sample belongs to multiple  $e_x$  or  $d_y$ .

### Estimating the prevalence phenotype

The first phenotype we consider is prevalence,  $p$ . Prevalence is usually defined as the fraction of samples in which a particular taxon is observed. Using the formulation above, the prevalence of microbe  $m$  in environment  $e_x$  and study  $d_y$  would be equal to:

$$\frac{\sum_{n \in (e_x \cap d_y)} (a_{m,n})}{\|e_x \cap d_y\|}$$

where we denote the quantity  $\sum_{n \in (e_x \cap d_y)} (1)$ , yielding the number of samples in the intersection of environment  $e_x$  and study  $d_y$ , as  $\|e_x \cap d_y\|$ .

We now take a slightly more general definition, such that a particular taxon’s true prevalence  $p_{m,(e_x \cap d_y)}$  is the probability of observing a particular microbe  $m$  in a set of samples  $(e_x \cap d_y)$ ,  $P(m|(e_x \cap d_y))$ . More specifically, we can say that  $p_{m,(e_x \cap d_y)}$  is the probability parameter of the Bernoulli random variable  $a_{m,n}, n \in (e_x \cap d_y)$ , which generates the values in our data matrix  $A$ :

$$a_{m,n} \sim \text{Bernoulli}(p_{m,(e_x \cap d_y)}), n \in e_x \cap d_y$$

$$p_{m,(e_x \cap d_y)} = P(m|e_x, d_y)$$

The maximum likelihood estimator of  $p_{m,(e_x \cap d_y)}$  given the data matrix  $A$  is then, as above, the fraction of subjects in  $F$  in which microbe  $m$  was observed:

$$\hat{p}_{m,(e_x \cap d_y)}^{\text{MLE}}(A) = \frac{\sum_{n \in (e_x \cap d_y)} (a_{m,n})}{\|e_x \cap d_y\|}$$

Because  $p$  is a proportion or probability, it is bounded between 0 and 1. The distribution of  $p$  is therefore highly non-normal, potentially violating assumptions of our regression model. We will therefore use  $\text{logit}(p)$  as our phenotype. However, this now introduces a problem because  $\hat{p}_{m,(e_x \cap d_y)}^{\text{MLE}}$  can take the values 0 and 1, leading to infinite estimates of  $\text{logit}(\hat{p}_{m,(e_x \cap d_y)})$ . We therefore instead use a shrunken estimate of  $\hat{p}_{m,(e_x \cap d_y)}$ .

Shrinkage estimators reduce the variance in the estimate of a parameter by combining it with prior information. These priors can be estimated from data (as in empirical Bayes approaches), estimated from independent information about the distribution of the parameter, or selected to be uninformative. Here, we use an uninformative prior, in this case a uniform distribution:

$$a_{m,n} \sim \text{Bernoulli}(p_{m,(e_x \cap d_y)}), \quad n \in e_x \cap d_y$$

$$p_{m,(e_x \cap d_y)} \sim \text{Beta}(1, 1)$$

Mechanistically, this is equivalent to performing additive smoothing, which effectively adds one pseudocount to the numbers of absences and presences:

$$\hat{p}_{m,(e_x \cap d_y)}^{\text{ADD}}(A) = \frac{1 + \sum_{n \in (e_x \cap d_y)} a_{m,n}}{2 + \|e_x \cap d_y\|}$$

Finally, we note that this estimate of prevalence is only valid within a single study  $d_y$ . However, what we really want is an estimator of prevalence that depends only on the environment  $e_x$ . We therefore marginalize out the effect of  $d_y$ :

$$P(m|e_x) = \sum_y P(m|e_x, d_y)P(d_y) = \sum_y P(m|(e_x \cap d_y))P(d_y)$$

where we let the prior probability  $P(d_y)$  simply be the inverse of the number of datasets present:

$$P(d_y) = \frac{1}{\|D\|}$$

Effectively, this weights each dataset inverse-proportionally to the number of samples, so that the study with the largest number of samples does not dominate our estimates of prevalence. (An alternative weighting could use, for example, the square root of the sample size, as is sometimes performed in fixed-effect meta-analyses.)

To avoid effects from additive smoothing dominating our estimates (as might happen if the same smoothing were applied to samples with different numbers of samples), we first obtain a marginalized version of the maximum-likelihood estimator, then perform additive smoothing on these marginalized estimates:

$$\hat{p}_{m,x,E,D}^{\text{ADDW}}(A) = \frac{\left( \sum_y \hat{p}_{m,(e_x \cap d_y)}^{\text{ML}} \frac{1}{\|D\|} \right) \left( \sum_y \|e_x \cap d_y\| \right) + 1}{\left( \sum_y \|e_x \cap d_y\| \right) + 2}$$

Finally, we use the logit of this estimate, i.e.,  $\text{logit}(\hat{p}_{m,x,E,D}^{\text{ADDW}}(A))$ , as the elements  $\phi_{m,x,E,D}^{\text{Prev}}(A)$  of our first phenotype  $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ :

$$\phi_{m,x,E,D}^{\text{Prev}}(A) = \text{logit} \left( \frac{\left( \sum_y \hat{p}_{m,(e_x \cap d_y)}^{\text{ML}} \frac{1}{\|D\|} \right) \left( \sum_y \|e_x \cap d_y\| \right) + 1}{\left( \sum_y \|e_x \cap d_y\| \right) + 2} \right)$$

To recapitulate, we use a logit-transformed, shrunken estimate of prevalence in a given environment, weighted so that each study contributes equally.

### Estimating environmental specificity scores

**Formulating the specificity score.** Prevalence gives us information about how commonly a microbe is seen in a particular environment. While useful, this concept does not address the difference between microbes that are specific for a given environment and those that have a cosmopolitan distribution. We therefore wanted to design a statistic capturing this environmental specificity. We define this statistic in terms of how *predictive* a particular microbe is for one out of a set of possible environments. (For simplicity, we only consider cases in which all environments were measured within a single study, and therefore drop  $D$  from these equations. This estimator could be extended in the future to account for study effects as above.)

Recall that prevalence can be defined as the probability  $P(m|e_x)$ , i.e., the probability of observing microbe  $m$  in environment  $e_x$ . To avoid potential sources of confounding error, we only consider environments both measured within the same study. We therefore let the environmental specificity score  $s_{m,x,E}$  equal the probability of observing a particular environment  $e_x$  out of a set of  $k$  environments  $E = \{e_1, e_2, \dots, e_k\}$ , given that we observe microbe  $m$ :

$$s_{m,x,E} = P(e_x|m)$$

which, by application of Bayes' rule and then marginalization, becomes:

$$P(e_x|m) = \frac{P(m|e_x)P(e_x)}{P(m)} = \frac{P(m|e_x)P(e_x)}{\sum_{y=1}^k P(m|e_y)P(e_y)} = \frac{p_{m,e_x}P(e_x)}{\sum_{y=1}^k p_{m,e_y}P(e_y)}$$

where  $P(m|e_x)$  is the prevalence  $p_{m,e_x}$  of microbe  $m$  in environment  $e_x \in E$ , and  $P(e_x)$  is the prior probability of observing environment  $e_x$ . The priors  $P(e)$  can be uninformative, in which case  $P(e_x) = 1/k$  for all  $x$ , meaning that all environments are equally likely. This is the approach we take for body site comparisons. Alternatively, for a disease state, they could be drawn from actual epidemiological data about the frequency of that disease in the population of interest. This is the approach we take for the Crohn's disease comparisons, taking  $P(e_{CD}) = 0.002$  [79], since in a Crohn's case-control study, the fraction of individuals with Crohn's will be much higher than the true prevalence of this disease in the population. In either case, the  $P(e_x)$  values do not depend on the values in the dataset  $A$ . An intriguing third possibility that would depend on  $A$  would be to estimate the priors  $P(e_x)$  based on the average observed  $\alpha$ -diversity within environment  $e_x$ , such that more diverse environments would be modeled as *a priori* more likely to contain any particular microbe.

**Motivating a shrunken estimator of  $s_{m,x,E}(A)$ .** One simple way to estimate  $\hat{s}_{m,x,E}(A)$  would be to simply plug in our estimates of  $\hat{p}_{m,e_x}^{ADD}(A)$ , yielding:

$$\hat{s}_{m,x,E}^{ADD}(A) = \frac{\hat{p}_{m,e_x}^{ADD}(A) \cdot P(e_x)}{\sum_{y=1}^k \hat{p}_{m,e_y}^{ADD}(A) \cdot P(e_y)}$$

However, for cases in which the number of total observations of a microbe  $\sum_{n \in e_y} a_{m,n}$  is low (imagine, e.g., a microbe that is observed once in environment  $e_1$  and zero times in  $e_2$ ), even the shrunken estimate  $\hat{p}_{m,e_x}^{ADD}(A)$  will have relatively high variance. This is particularly problematic here because both the numerator and denominator of  $\hat{s}_{m,x,E}^{ADD}(A)$  depend on  $\hat{p}_{m,x,E}^{ADD}(A)$ , so as  $p_{m,1,E} \rightarrow 0$ , the standard error of  $\hat{s}_{m,x,E}^{ADD}(A)$  will tend to increase. This means that the microbes with the least-confidently estimated prevalences will tend to have high leverages in the regression, distorting the results. The confounding between the magnitude of  $\hat{s}_{m,x,E}^{ADD}(A)$  and its standard error also leads to heteroskedasticity, or unequal variance across the residuals, violating one of the main linear model assumptions.

To account for these issues, we construct a more aggressively-shrunk estimator of  $s_{m,x,E}$ . We assume that most microbes do not differ substantially between environments, and therefore shrink estimates of  $s_{m,x,E} = P(e_x|m)$  towards the prior probability  $P(e_x)$ , indicating that that this microbe is uninformative about the environment. To accomplish this, we use a maximum a posteriori (MAP) estimator,  $\hat{s}_{m,x,E}^{\text{MAP}}(A)$ , with a Laplace prior centered on  $P(e_x)$ . Laplace priors are also used in the Bayesian lasso to make parameter estimates more sparse, by shrinking them to zero. However, critically, we are not using the Laplace prior to perform model selection, since the exact same model is fit as in Eq (1); we are only using it to reduce the variance in estimating  $\hat{s}_{m,x,E}$ ; unlike the Bayesian lasso, we therefore use no information about the independent variable (gene presence-absence) in obtaining estimates of  $\hat{s}_{m,x,E}$ .

**An estimator of  $s_{m,x,E}$  using Laplace shrinkage.** We are interested in estimating the specificity score  $s_{m,x,E}$  of a microbe  $m$  in environment  $e_x$ , given observations of microbial presence or absence  $a_{m,n}$  where  $n$  is a sample, a mapping of samples to environments  $E$ , and the prior probability  $P(e_x)$  of choosing the environment  $e_x \in E$ . In order to do this, we construct a model of the data that depends on  $s_{m,x,E}$ , then maximize the posterior probability of observing the data given the model.

Previously, we modeled the presence-absence data  $a_{m,n}$  in terms of prevalence  $p_{m,e_x}$ . It is therefore useful to define a function of our parameter of interest  $s_{m,x,E}$  that yields  $p_{m,e_x}$ . We first restate the definition of the specificity score  $s_{m,x,E}$ , noting that it depends on the prevalences  $p_{m,e_y}$  of microbe  $m$  across all  $k$  environments  $e_y \in E$ ,  $|E| = k$ , and the prior probabilities of these environments  $P(e_y)$ :

$$s_{m,x,E} = \frac{p_{m,e_x} P(e_x)}{\sum_{y=1}^k p_{m,e_y} P(e_y)}$$

Rearranging this equation to get the prevalence in an environment  $p_{m,e_x}$  in terms of specificity  $s_{m,x,E}$  we get:

$$p_{m,e_x} = \frac{(s_{m,x,E}) \sum_{y=1}^k p_{m,e_y} P(e_y)}{P(e_x)}$$

We then define the function  $h(s_{m,x,E})$  as the function that transforms specificity scores into prevalences:

$$h(s_{m,x,E}) = p_{m,e_x} = \frac{(s_{m,x,E}) \sum_{y=1}^k p_{m,e_y} P(e_y)}{P(e_x)}$$

Next we consider the Laplace distribution that we use as the prior for our specificity score  $s_{m,x,E}$ . This distribution has two parameters, a mean and a width. The mean is given by the prior probability of environment  $e_x$ ,  $P(e_x)$ . For the moment, we assume that we have chosen a width  $b$ , which governs the amount of shrinkage. We can then model the environmental specificity score as follows:

$$\begin{aligned} \text{logit}(s_{m,x,E}) &\sim \text{Laplace}(P(e_x), b) \\ a_{m,n} &\sim \text{Bernoulli}(h(s_{m,x,E})), \quad n \in e_x \end{aligned}$$

We can then obtain an estimate of  $s_{m,x,E}$  given the prior probability of environment  $e_x$  ( $P(e_x)$ ), our observed data  $a_{m,n}$ , and a Laplace width  $b$ , by maximizing the posterior probability of our

data  $a_{m,n}$ . This maximum *a posteriori* (MAP) estimate takes the following form:

$$\hat{s}_{m,x,E}^{\text{MAP}}(A; b) = \operatorname{argmax}_{s_{m,x,E}} \mathcal{L}(A|h(s_{m,x,E})) \times f(A|b) \tag{2}$$

where  $\hat{s}$  is the parameter being estimated,  $A$  represents our data matrix with elements  $a_{m,n}$ ,  $\mathcal{L}$  represents the likelihood function of the distribution from which the data is assumed to be drawn,  $f$  represents the density function of our Laplace prior on  $s$  (without which the estimator reduces to the maximum-likelihood estimator), and  $b$  is the given width parameter as above.

In order to calculate  $h(s_{m,x,E})$ , we need to know the prevalences across all environments  $p_{m,e_y}$ ,  $e_y \in E$ , which are not given. However, we can estimate them as before from  $a_{m,n}$ . We let  $h^{\text{ML}}$  be the function  $h$  where the values of  $p_{m,e_y}$  come from the maximum-likelihood estimates  $\hat{p}_{m,e_y}^{\text{ML}}$ , so that:

$$h^{\text{ML}}(s_{m,x,E}, A) = \frac{(s_{m,x,E}) \sum_{y=1}^k \hat{p}_{m,y,E,A}^{\text{ML}} P(e_y)}{P(e_x)}$$

We can now expand this to give the final maximization:

$$\hat{s}_{m,x,E}^{\text{MAP}}(A; b) = \operatorname{argmax}_{s_{m,x,E}} \left[ \left( \binom{\|e_x\|}{n_{\text{obs}}} (h^{\text{ML}}(s_{m,x,E}, A))^{n_{\text{obs}}} (1 - h^{\text{ML}}(s_{m,x,E}, A))^{\|e_x\| - n_{\text{obs}}} \right) \times \left( \frac{1}{2b} \exp\left(-\frac{|\operatorname{logit}(s_{m,x,E}) - \operatorname{logit}(P(e_x))|}{b}\right) \right) \right] \tag{3}$$

where  $n_{\text{obs}}$  refers to  $\sum_{n \in e_x} a_{m,n}$ , i.e., the number of times microbe  $m$  was observed in environment  $e_x$ .

**Choosing the Laplace width parameter.** We have assumed above that we are given the Laplace width parameter  $b$ . To choose appropriate, dataset-specific values of  $b$ , which controls how much  $\hat{s}_{m,x,E}^{\text{MAP}}(A; b)$  is shrunk back to the prior, we performed simulations. We chose a simulation-based approach instead of, for example, cross-validation because we lacked labeled examples of microbes that truly differed between environments. Instead, we constructed datasets  $A'$  with elements  $a'_{m,n}$  where we “knew” that some microbes ( $m \in M_0$ ) were not informative about the environment and others ( $m \notin M_0$ ) had “true” differences, by simulating data with the following model:

$$\begin{aligned} a'_{m,n} &\sim \text{Bernoulli}(p_{m,y,E}), \quad n \in e_y \\ p_{m,y,E} &= \begin{cases} q_m & (y \neq x) \\ \text{logistic}(\operatorname{logit}(q_m + z_m)) & (y = x) \end{cases} \\ z_m &\sim \begin{cases} z \cdot (2 \cdot (\text{Bernoulli}(r)) - 1) & (m \notin M_0) \\ 0 & (m \in M_0) \end{cases} \\ q_m &\sim \text{Beta}(a, b) \end{aligned}$$

In other words, for each species  $m$  in different environments  $e_y \in E$ , presence-absence  $a'_{m,e_y}$  was modeled as a Bernoulli random variable. The success parameter from this Bernoulli was drawn from a Beta distribution with parameters  $a$  and  $b$ , which were fit from a single environment in the corresponding real dataset using maximum-likelihood, thus ensuring that the simulated species had similar baseline prevalences as real species. In species with no difference between environments  $m \in M_0$ , the true prevalence  $p_{m,x,E}$  was set to be equal between the environment of interest  $e_x$  and all other environments; in species with true differences between

environments ( $m \notin M_0$ ), in contrast, the effect size  $z$  was either added or subtracted from the logit-prevalence (with the parameter  $r$  controlling the proportion of positive true effects). The number of null species  $||M_0||$  was set to 25% of the total number of simulated species  $||M||$ , which was matched to the real dataset.

For a given simulated dataset and value of  $b$ , the false positive rate ( $FPR_b$ ) and the true positive rates for  $F > 0$  and  $F < 0$  ( $TPR_{pos}$  and  $TPR_{neg}$ , respectively) were calculated:

$$\begin{aligned} FPR_b &= \#(|P(e_x|(m \in M_0)) - P(e_x)| > \epsilon) / ||M_0|| \\ TPR_{pos_b} &= \#(P(e_x|(z_m > 0)) - P(e_x) > \epsilon) / \#(z_m > 0) \\ TPR_{neg_b} &= \#(P(e_x) - P(e_x|(z_m < 0)) > \epsilon) / \#(z_m < 0) \end{aligned}$$

Since we are using numerical optimization, posterior probabilities are not always shrunk exactly to the prior; we therefore use a tolerance parameter  $\delta$  set at  $P(e_x) \cdot 0.005$  to account for numerical error. The tuning parameter  $b$  was then optimized according to the following piecewise continuous function, which increases from 0 to 1 until the false positive rate drops to 0.05 or lower (in order to guide the optimizer), and then increases above 1 in proportion to the average (geometric mean) of the positive and negative effect true positive rates:

$$b_{\text{optim}} = \operatorname{argmax}_b \begin{cases} 1 - FPR_b & FPR_b > 0.05 \\ 1 + \sqrt{TPR_{pos_b} \times TPR_{neg_b}} & FPR_b \leq 0.05 \end{cases}$$

Given  $z = 2$  and  $r = 0.5$ , for Crohn's disease,  $b_{\text{optim}}$  was estimated at  $b = 0.16$  and for the body site specificity,  $b_{\text{optim}}$  was estimated at  $b = 0.282$ . (Changing  $z$  to 1 or 0.5, or changing  $r$  to 0.1 or 0.9, resulted in very similar estimates of  $b_{\text{optim}}$ . Additionally,  $b_{\text{optim}}$  estimates were consistent across several orders of magnitude of  $\epsilon$ ; see [S4 Fig](#)).

**Worked example.** An example showing the effect of this procedure on real data can be seen in [S3 Fig](#). The microbe *Bacillus subtilis* is detected once in the healthy cohort and once in the Crohn's cohort, while *Bacteroides fragilis* is present in 24/38 healthy subjects but 13/13 Crohn's subjects ([S3A Fig](#)). The maximum-likelihood values of  $\hat{p}_{m,CD,E}(A)$  and  $\hat{s}_{m,CD,E}(A)$  are therefore much higher for *B. subtilis* than for *B. fragilis*, even though the evidence for a difference across environments in the prevalence of *B. subtilis* is much weaker ([S3B and S3C Fig](#)). In contrast, the Laplace prior ([S3D Fig](#)) successfully shrinks the estimate of the *B. subtilis* specificity score back to the baseline ( $P(e_{CD}) = 0.002$ ), while the evidence for *B. fragilis* overcomes this prior and yields an estimate close to the maximum-likelihood value (0.0031; [S3E Fig](#)).

### Alternatives to shrinkage estimation of environmental specificity scores

An alternative to using the Laplace shrinkage estimator would be to allow all taxa to contribute to the regression, but to downweight taxa with less-confidently measured phenotypes. In generalized least squares (GLS), this is typically accomplished by scaling the variance-covariance matrix of the residuals by the variances of the estimators. This in effect says that the residuals are expected to be more dispersed around the regression line when the variance of the estimator is high: equivalently, this procedure weights each point in least-squares by the inverse of the estimator's variance. We represent these variances as  $v_m \equiv (SE(\hat{s}_{m,x,E}(A)))^2$ . The covariance



matrix is then:

$$\Sigma^{\text{WLS}} = \begin{bmatrix} \sigma_1 v_1 & \sigma_{1,2} \sqrt{v_1 v_2} & \cdots & \sigma_{1,m} \sqrt{v_1 v_m} \\ \sigma_{2,1} \sqrt{v_2 v_1} & \sigma_2 v_2 & & \\ \vdots & & \ddots & \\ \sigma_{m,1} \sqrt{v_m v_1} & & & \sigma_m v_m \end{bmatrix}$$

where  $\sigma$  values are as above. (Off-diagonal elements are weighted by the geometric mean of the variances, thus giving the same correlation structure as before.)

Because we have the data matrix  $A$  from which  $\hat{s}_{m,x,E}^{\text{ADD}}$  was estimated, we can estimate the variances of these estimates by bootstrapping. Denoting the  $\hat{s}$  estimates derived from bootstrap sample  $c$  as  $\hat{s}_{m,x,E}^c$ , where  $c \in \{1, C\}$  and  $C$  is the number of bootstraps, and letting the mean across bootstrap samples be  $\bar{s}_{m,x,E}$ , then  $v_m = (\frac{1}{C-1} \sum_{c=1}^C (\hat{s}_{m,x,E}^c - \bar{s}_{m,x,E})^2)^2$ .

Both approaches (Laplace shrinkage estimation and WLS) account for variability in the accuracy of estimating  $s_{m,x,E}$ , but in different ways. We would expect them to agree more when  $s_{m,x,E}$  values were more confidently estimated (meaning the evidence for difference from the prior would be stronger and the extent of downweighting would be lower), and when more of the  $s_{m,x,E}$  values diverged substantially from the prior (leading to less sparsity in  $\hat{s}_{m,x,E}^{\text{Laplace}}$ ). Indeed, when comparing body site specificities, which are based on higher numbers of samples and involve comparisons across highly divergent environments, both approaches yield more similar estimates, and become more concordant when the effect sizes are calculated only over significantly-different gene families. In contrast, for the Crohn’s disease comparison, in which the environments (healthy vs. diseased gut) are more closely related and the number of samples is smaller, the two methods tend to disagree more, especially in phyla where most taxa are shrunk back to the prior (S2 Table). (We implemented WLS using the `gls` function provided by the `nlme` R package [80].)

These results reflect the different underlying assumptions of each estimator: Laplace shrinkage assumes that taxa are not truly varying across environments without strong evidence, while WLS uses information from all taxa but downweights less-confidently-observed species. For the purposes of this manuscript, we focused on the results based on Laplace shrinkage estimation, since we believe the assumption that most taxa do not change is appropriate when comparing the same body site in health and disease. However, either approach may be preferable depending on the precise scenario being studied. It could also be possible to combine the two approaches by, for example, using the full posterior distribution of  $\hat{s}_{m,x,E}$  to derive weights.

### Power analysis

To test the power and false positive rate of our method, we used parametric simulations, either under the null hypothesis in which a gene had no effect on the phenotype, or under the alternative hypothesis in which it had a defined effect. These involved generating one binary genotype and one continuous phenotype per simulation. These were parameterized as follows:

- The continuous phenotype  $\vec{\phi}^{\text{Sim}}$  is simulated according to a Brownian motion model with parameters  $\beta_0$  corresponding to the ancestral state of the phenotype and  $\sigma^2$  corresponding to the phenotype’s overall variance (i.e., the diagonal of  $\Sigma$  in the phylogenetic model).

- The binary genotype  $\vec{I}^{\text{sim}}$  is generated from a Markov process as in Ives and Garland [34], with parameters  $\alpha$  and  $\beta_1$ .  $\alpha$  gives the sum of the transition probabilities going from 0 to 1 and from 1 to 0.  $\beta_1$  gives the effect size: that is, how much the simulated phenotype influences the binary genotype (in logit space).

We perform the following process, given a choice of  $\alpha$  and  $\beta_1$ , for each phylum  $h$ :

1. Estimate the parameters  $\hat{\beta}_0^{\text{Prev}}$  and  $\hat{\sigma}^2$  by fitting the following intercept-only phylogenetic model to the real prevalence phenotype:

$$\vec{\phi}_{x,E}^{\text{Prev}}(A) = \beta_0^{\text{Prev}} + \vec{\epsilon}$$

where  $\epsilon \sim \mathcal{N}(0, \Sigma)$  and  $\text{Diag}(\Sigma) = \sigma^2$ .

2. For each of  $B$  simulations:

- a. Generate a continuous phenotype  $\vec{\phi}^{\text{Sim}}$  according to a Brownian motion process, evolving along the tree of phylum  $h$ , with parameters  $\hat{\beta}_0^{\text{Prev}}$  and  $\hat{\sigma}^2$ .
- b. Generate a binary genotype  $\vec{I}^{\text{sim}}$  according to an Ives-Garland Markov process with parameters  $\alpha$  and  $\beta_1$  and a covariate  $\vec{\phi}_{x,E}^{\text{Sim}}$ .
- c. Fit the following regression equation:

$$\vec{\phi}^{\text{Sim}} = \beta_0^{\text{Test}} + \beta_1^{\text{Test}} \vec{I}^{\text{sim}} + \epsilon$$

using either a standard linear model or the phylogenetic model (as specified in Methods, “Fitting linear vs. phylogenetic models”).

- d. Return the  $p$ -value for the test of the null hypothesis  $\beta_1^{\text{Test}} = 0$ .
3. The fraction of  $p$ -values  $\leq 0.05$  yields the power (when given  $\beta_1 > 0$ ) or the false positive rate (when given  $\beta_1 = 0$ ) of the test.

The binary genotype effect size  $\beta_1$  is not a linear function of the effect of the gene on prevalence. A more intuitive description of the effect size might be the (average) fold-change in prevalence associated with a gene’s presence. Because the parameters  $\beta^{\text{Test}}$  are on a logit scale, this quantity would be equal to:

$$F = \frac{\text{logistic}(\beta_1^{\text{Test}} + \beta_0^{\text{Test}})}{\text{logistic}(\beta_0^{\text{Test}})}$$

This quantity will depend on the amount of phylogenetic signal in the binary genotype  $\alpha$ , the tree along which genotypes and phenotypes are simulated, and the “input” effect size  $\beta_1$ . Accordingly, we simulated sets of 50 binary genotypes with given effect sizes  $\beta_1 \in \{0, 0.5, 0.75, 1.0, 1.25\}$  and values of  $\alpha \in \{0, 25, 50\}$ . While there is substantial variation, in general an “input” effect size (i.e.,  $\beta_1$ ) of 1.0 approximately corresponds to  $F \approx 1.72$ , a 72% increase in prevalence, and an “input” effect size of 0.5 corresponds to  $F \approx 1.5$ , a 50% increase (S8 Fig).

### Assessing the potential impact of sampling with left-censoring

One potential pitfall with applying linear methods occurs when the distribution of the response variable (here, our phenotype) has a minimum value. This arises because, within a particular dataset, the lowest possible value of  $\hat{p}_{m,x,E}^{\text{ADD}}(A)$  is equal to  $(\|e_x\| + 2)^{-1}$ , where  $\|e_x\|$  is the number of samples in environment  $e_x$ . This phenomenon is referred to as “left-censoring.”

(Some may have encountered the term “left-censoring” in the context of participants who join a study having already experienced an event of interest. The time they experienced this event is therefore lower by an unknown amount than the lowest-possible measured value. While the domain and application are different, the statistical phenomenon is the same.) Left-censoring can result in inaccurate  $p$ -values because the variance is mis-estimated for the data points below the limit of detection. We therefore empirically assessed the impact of left-censoring in simulation, and also created extensions of the method to be used when its impact is noticeable.

Empirically, our prevalence phenotype  $\vec{\phi}_{x,E}^{\text{Prev}}(A)$  displays substantial left-censoring (S6A Fig). The distribution of this phenotype fits well to a normal with left-censoring at the limit of detection, in this case approximately  $-0.50$  standard deviations below the mean (AIC using truncated normal and censoring: 20844.82; AIC using standard normal: 21833.38).

We therefore repeated the simulation process above, but after using our continuous phenotype to generate the binary genotype in step 2.b, we truncated the continuous phenotype  $\vec{\phi}^{\text{Sim}}$  artificially at a specified number of standard deviations  $K$  below the mean, yielding  $\vec{\phi}^{\text{Cens}}$ . We then replaced  $\vec{\phi}^{\text{Sim}}$  with  $\vec{\phi}^{\text{Cens}}$  in the regression in step 2.c.

Using this simulation framework, we benchmarked three different ways to test the significance of  $\beta_1^{\text{Test}}$  in the phylogenetic model. We computed  $p$ -values in one of the following three ways:

1. T-STATISTIC: Return the  $p$ -value of a  $t$ -test of the null hypothesis that  $\beta_1^{\text{Test}} = 0$ , as in Methods, “Modeling gene-phenotype association”.
2. PARAMETRIC BOOTSTRAP: Simulate a number  $C$  of null binary genotypes  $\vec{I}^c$  with  $\beta_1 = 0$ , where  $c \in \{1, \dots, C\}$ , and collect the estimates  $\hat{\beta}_1^c$  based on the following phylogenetic linear model:

$$\vec{\phi}^{\text{Sim}} = \beta_0^c + \beta_1^c \vec{I}^c + \epsilon$$

Then compute the fraction that are at least as extreme as the test statistic  $\hat{\beta}_1^{\text{Test}}$  and return this as a  $p$ -value:

$$\frac{1}{C} \sum_{c=1}^C \begin{cases} 1 & ((\hat{\beta}_1^c - \bar{\beta}_1^C)^2) \geq ((\hat{\beta}_1^{\text{Test}} - \bar{\beta}_1^C)^2) \\ 0 & ((\hat{\beta}_1^c - \bar{\beta}_1^C)^2) < ((\hat{\beta}_1^{\text{Test}} - \bar{\beta}_1^C)^2) \end{cases}$$

where  $\bar{\beta}_1^C$  is the mean of  $\beta_1^c$  values. To save computation time on high  $p$ -values, we use early stopping: after every 25 such simulations, if the resulting  $p$ -value would already be guaranteed to exceed 0.05, we stop and return the  $p$ -value based on the current number of simulations.

3. MOCK-UNCENSORED BOOTSTRAP: As #2, simulate a number  $C$  of additional binary genotypes with  $\beta_1 = 0$ . Instead of calculating  $p$ -values as in #2, however:
  - a. Simulate the same number  $C$  of “uncensored” versions of the continuous phenotype  $\vec{\phi}^{\text{Unc}}$ , in effect “filling in” or imputing the censored values with random values from the predicted tail of the distribution (see S6B and S6C Fig for an illustration):
    - i. First, fit a left-truncated normal distribution to the part of  $\vec{\phi}^{\text{Cens}}$  that is above the lowest value (assumed to be the limit of detection) by maximum likelihood (using `fitdistscens` in R package `fitdistrplus` [81]), yielding mean  $\mu^{\text{Trunc}}$ , standard deviation  $\sigma^{2\text{Trunc}}$ , and lower truncation point  $\min(\phi_m^{\text{Cens}})$  statistics (e.g., S6B Fig).

- ii. Next, for  $c \in \{1, \dots, C\}$ , generate a vector  $\vec{T}^c$  whose elements  $T_m^c$  are realizations of the random variable  $T \sim N^{\text{RTrunc}}(\mu^{\text{Trunc}}, \sigma^{2\text{Trunc}}, \min(\phi_m^{\text{Cens}}))$ . Here,  $N^{\text{RTrunc}}$  represents a right-truncated normal distribution, having mean  $\mu^{\text{Trunc}}$ , variance  $\sigma^{2\text{Trunc}}$ , and upper truncation point  $\min(\phi_m^{\text{Cens}})$ . Then generate the “uncensored” vector  $\vec{\phi}^{\text{Unc}^c}$  with elements  $\phi_m^{\text{Unc}^c}$  as follows:

$$\phi_m^{\text{Unc}^c} = \begin{cases} \phi_m^{\text{Cens}} & \phi_m^{\text{Cens}} > \min(\phi_m^{\text{Cens}}) \\ T_m^c & \phi_m^{\text{Cens}} = \min(\phi_m^{\text{Cens}}) \end{cases}$$

(An example can be seen in [S6C Fig](#)).

- b. For each  $c \in \{1, \dots, C\}$ :

- i. estimate a test statistic  $\hat{\beta}_1^{\text{Test}^c}$  by fitting the following phylogenetic linear model:

$$\vec{\phi}^{\text{Unc}^c} = \beta_0^{\text{Test}^c} + \beta_1^{\text{Test}^c} \vec{T}^{\text{Sim}} + \epsilon$$

- ii. estimate a null test statistic  $\hat{\beta}_1^c$  by fitting the following phylogenetic linear model:

$$\vec{\phi}^{\text{Unc}^c} = \beta_0^{\text{Test}^c} + \beta_1^c \vec{T}^c + \epsilon$$

- c. Calculate  $p$ -values as follows:

$$\frac{1}{C} \sum_{c=1}^C \begin{cases} 1 & ((\hat{\beta}_1^c - \bar{\beta}_1^C)^2) \geq ((\bar{\beta}_1^{\text{Test}^C} - \bar{\beta}_1^C)^2) \\ 0 & ((\hat{\beta}_1^c - \bar{\beta}_1^C)^2) < ((\bar{\beta}_1^{\text{Test}^C} - \bar{\beta}_1^C)^2) \end{cases}$$

where  $\bar{\beta}_1^C$  is the mean of the  $\beta_1^c$  values and  $\bar{\beta}_1^{\text{Test}^C}$  is the mean of the  $\beta_1^{\text{Test}^c}$  values.

Intuitively, the second method generates a null distribution via simulation, while the third method additionally reduces the impact of data points at the limit of detection, by randomly imputing them from the best-fit normal distribution. We then calculated power and FPR for each of these three methods, varying the amount of censoring  $K$  and the effect size  $\beta_1$  (see [S7 Fig](#)). Interestingly, for the level of censoring in our data ( $K = -0.50$ ), the false-positive rate in all three methods remained well-controlled, although power dropped. The mock-uncensored bootstrap had lower power overall and became more conservative, especially in the case where the phylogenetic signal was highest ( $\alpha = 0$ ) and where the level of censorship  $K$  was highest.

Another common approach to this problem is the tobit model: the true value of the response variable is treated as a hidden variable, and expectation-maximization is used to fit the regression parameters based on the observed censored values. Given that the degree of censoring we observed did not appear to inflate the false positive rate in any of the methods we tested, we opted not to construct a phylogenetic tobit model; however, this could be an interesting area of future research.

### Assessing the impact of compositionality

Because relative abundances are compositional (i.e., sum to 1), changes in highly abundant taxa combined with read sampling can lead to skewed estimates of relative abundance. For

example, if a very abundant microbe exhibits large changes in relative abundance across samples, other microbes will appear to become less abundant simply because they make up a smaller proportion of the total reads, regardless of whether their level has actually changed [82]. This necessitates the use of compositional data analysis methods such as fitting intrinsically compositional distributions to the data (e.g., multinomial [83]) or transforming them such that they no longer lie on a simplex (e.g., the clr-transform [84]). However, the impact of compositionality on prevalence was *a priori* less clear, because while prevalences are based on presence versus absence, which could be affected by sampling, prevalences themselves do not have to sum to any particular value. Nonetheless, since we compute prevalence from relative abundance, we undertook a simulation based investigation of prevalence estimation accuracy as a function of various study design and biological parameters.

To directly compare the degree to which spurious correlation in prevalence and relative abundance could be induced by compositionality, we simulated correlated absolute abundances of microbes with no true correlations between them, then converted these uncorrelated absolute abundances to relative abundances. To look for spurious correlation, we plotted the distribution of pairwise microbial correlations in absolute and relative abundances. We next used relative abundance to compute prevalence. Since there is only one prevalence measurement per microbe per data matrix, we repeated the simulation multiple times with the same parameters, then assessed microbial correlations across simulation runs, again starting from both absolute and relative abundances.

**Simulating data.** For each of 25 simulations  $i \in I$ , we assumed each of 1000 simulated microbes  $m \in M$  had absolute abundances  $b_{(m,n)_i}^{Abs}$ , over a set of 100 simulated samples  $n \in N$ . These were independently log-normally distributed, centered on microbe-specific average values  $\mu_m$  that were themselves log-normally distributed. We also assumed a microbe-specific prevalence  $z_m$  drawn from a logistic-normal distribution:

$$\begin{aligned} b_{(m,n)_i}^{Abs} &\sim \log \mathcal{N}(\mu_m) \cdot c_{(m,n)_i} \\ c_{(m,n)_i} &\sim \text{Bernoulli}(1 - z_m) \\ \log(\mu_m) &\sim \mathcal{N}(1, 3) \\ \text{logit}(z_m) &\sim \mathcal{N}(-4, 2) \end{aligned}$$

We converted the resulting absolute abundance matrices  $B_i^{Abs}$  to relative abundances in two different ways. First, we simply column-normalized, yielding a matrix  $B_i^{relCN}$  with elements  $b_{(m,n)_i}^{relCN} = b_{(m,n)_i}^{Abs} / \sum_{m \in M} b_{(m,n)_i}^{Abs}$ . Second, we first generated compositional count data  $B_i^{Counts}$ , with column vectors  $\vec{b}_{n_i}^{Counts} = [b_{1,n_i}^{Counts}, b_{2,n_i}^{Counts}, \dots, b_{||M||,n_i}^{Counts}]^t$  simulated as follows:

$$\begin{aligned} \vec{b}_{n_i}^{Counts} &\sim \text{Multinomial}(r_n, \theta_{n_i}) \\ \theta_{n_i} &\sim \text{Dirichlet}(b_{(1,n)_i}^{Abs}, b_{(2,n)_i}^{Abs}, \dots, b_{(||M||,n)_i}^{Abs}) \\ r_n &\sim \text{DiscreteUniform}(10^5, 10^6) \end{aligned}$$

We then converted these counts to relative abundances  $B_i^{relDM}$  such that  $b_{(m,n)_i}^{relDM} = b_{(m,n)_i}^{Counts} / \sum_{m \in M} b_{(m,n)_i}^{Counts}$ .

**Evaluating extent of spurious correlation for relative abundance.** For each  $||M|| \times ||N||$  matrix  $B$ , we obtained a set of correlations  $L(B) = \{\rho_{(x,y)_B} : x, y \in M, y < x\}$ , where  $\rho_{(x,y)_B}$  is the Pearson correlation between the row-vectors  $\vec{b}_{(x,n \in N)}$  and  $\vec{b}_{(y,n \in N)}$ .  $L(A)$  can also be thought of as the entries in the lower-triangular part of the correlation matrix of a given  $B$ .

To show that compositionality is enough to induce spurious correlation in abundances, we computed the correlations from absolute abundances  $L(B_i^{\text{Abs}})$  and compared them with the correlation matrices derived from relative abundances,  $L(B_i^{\text{relCN}})$  and  $L(B_i^{\text{relDM}})$ . Because discrete sampling introduces many more zeros into the matrix  $B_i^{\text{relDM}}$ , we also made the further comparison of pairwise correlations calculated using only non-zero elements of  $B_i^{\text{relDM}}$ . We denote this non-zero correlation as  $LZ(B) = \{\rho_{(x,y)_B}^z : x, y \in M, y < x\}$ . Here,  $\rho_{(x,y)_B}^z$  is the Pearson correlation between the row-vectors  $\vec{b}_{x,v \in V}$  and  $\vec{b}_{y,v \in V}$ , and  $V$  is the set of samples in which both microbes  $x$  and  $y$  have non-zero abundance, i.e.,  $V = \{v: (b_{x,v} > 0, b_{y,v} > 0)\}$ .

The correlation matrices of simulated absolute abundance data, as expected, are centered on zero. In contrast, the correlation matrix for the relative abundances  $L(B_i^{\text{relCN}})$  is shifted markedly to the right (S9A Fig). While this seems to improve when plotting the results of the Dirichlet-Multinomial simulations,  $L(B_i^{\text{relDM}})$ , when correlations are calculated over only non-zero elements, we see that they once again are shifted to the right by the same degree. These simulations show that relative abundance estimates are biased due to compositionality, consistent with previous work [85, 86], but that a high proportion of zeros (i.e., microbe not detected) can reduce this effect.

**Evaluating extent of spurious correlation for prevalence.** While relative abundance can be measured for each microbe in each sample, prevalence must be measured over a set of samples. To assess the impact of compositionality on prevalence, we therefore repeated the above simulation for 100 simulations  $i \in I$ , but this time held  $\mu_m$  and  $z_m$  constant across runs of the simulation to make them comparable.

Mirroring the procedure we used to estimate prevalence in real data, we binarized our  $B_i$  abundance matrices to yield corresponding presence-absence matrices  $A_i$ , and then calculated estimated prevalence of each microbe  $m$ ,  $\hat{p}_{m,N}^{\text{ADD}}(A_i)$ , from our absolute abundances  $A_i^{\text{Abs}}$  and from the derived relative abundances  $A_i^{\text{relDM}}$ . We performed this for every simulation  $i \in I$ , forming two matrices of prevalences  $\Pi^{\text{Abs}}$  and  $\Pi^{\text{relDM}}$ , each with  $\|M\|$  rows,  $\|I\|$  columns, and elements  $\pi_{m,i} = \hat{p}_{m,N}^{\text{ADD}}(A_i)$ . We also investigated the impact of defining prevalence as absolute abundance over a low but non-zero threshold of detection, which we defined as  $\Pi^{\text{AbsThresh}}$ . We chose a threshold of  $b_{m,n}^{\text{Abs}} > 0.5$ . We then plotted densities of the lower-triangular part of the microbial correlation matrices  $L(\Pi^{\text{Abs}})$ ,  $L(\Pi^{\text{AbsThresh}})$ , and  $L(\Pi^{\text{relDM}})$ , analogously to above. (We did not compute  $L(\Pi^{\text{relCN}})$  because it should always be equal to  $L(\Pi^{\text{Abs}})$ .)

In contrast to what we observed for relative abundance, the distributions of prevalences were almost identical when based on absolute abundances versus resampled relative abundances (S9B Fig). This strongly suggests that compositionality does not induce spurious correlations between taxon prevalences in the same way, or to the same degree, as it does for relative abundances.

**Other sources of error in prevalence.** Another source of error in prevalence estimates, which is related to our discussion of left-censoring, is the distinction between “structural” versus “sampling” zeros in relative abundance data [87]. Structural zeros arise from samples in which a microbe is really not present, whereas sampling zeros arise when a microbe is present but does not happen to be sampled. The fewer the number of samples, and the fewer reads per sample, the more likely sampling zeros are to occur.

To get a better sense of when prevalence estimates are likely to be reliable, we simulated relative abundance data, varying the number of samples and reads, and then calculated the correlation between the true and the estimated prevalences. Our simulation framework can be thought of as a zero-inflated multinomial, with details as follows (using  $\odot$  to mean element-

wise multiplication):

$$\begin{aligned} \vec{b}_n &\sim \text{Multinomial}(r_n, \theta) \odot \vec{c}_{M,n} \\ c_{m,n} &\sim \text{Bernoulli}(1 - z_m) \\ \text{logit}(z_m) &\sim N(-4, 2) \\ \theta_n &\sim \text{Dirichlet}(\alpha_1, \dots, \alpha_{||M||}) \\ \alpha_m &\sim \text{DiscreteUniform}(1, 10^5) \\ \log(r_n) &\sim N(q, 0.5) \end{aligned}$$

In other words, we randomly simulate one Dirichlet “meta-community” with  $||M|| = 6000$  uniformly-distributed  $\alpha$  parameters (matching the approximate number of taxa measured in MIDAS). For each sample  $n$ , we then draw a multinomial  $\theta_n$  from this Dirichlet, then draw reads from that multinomial, given a number of total reads  $r_n$ . These total reads are log-normally distributed around a central value  $q$ . We then draw random prevalences  $z_m$  with a logit-normal distribution, then draw structural zeros  $c_{m,n}$  based on these (sampling zeros arise from the sampling inherent in the multinomial distribution). Element-wise multiplication with the vector of structural zeros therefore leads to counts from a multinomial with excess zeros.

Our simulation varies  $q$ , the average (geometric mean) number of reads across samples, and  $||N||$ , the total number of samples. After constructing a simulated abundance matrix  $B_i^{\text{Counts}}$  from each simulation, we binarize it to obtain a presence-absence matrix  $A_i$ , then compute estimated prevalences for each microbe  $m$ ,  $\hat{p}_{m,N}^{\text{ADD}}(A_i)$ . We finally calculate the correlation of  $\text{logit}(z_m)$  and  $\text{logit}(\hat{p}_{m,N})$ , as well as the fraction of censoring that we observe (i.e., the fraction of values of  $\text{logit}(z_m) \leq \min(\text{logit}(\hat{p}_{m,N}))$ ), and the correlation over only the non-censored values (i.e., for  $\{m : \hat{p}_{m,N} > \frac{1}{N+2}\}$ ).

In these simulations, the factor with the greatest effect on prevalence reliability is the number of samples, not read depth (S10 Fig). While this might initially seem counter-intuitive, we are assuming in this simulation that most taxa have low prevalences (the mean of  $\text{logit}(z_m)$  is  $-4$ , corresponding to  $\sim 2\%$  prevalence), much like what we observed in real data. In this setting, structural zeros will tend to outweigh sampling zeros, so increasing the sampling depth will not necessarily improve estimates of prevalence. In this setting, having fewer than 50 samples increases the level of censoring above approximately 50%; this amount of censoring (zero standard deviations below the mean in S7 Fig) does not increase the false discovery rate in our simulations, but does lower the power to detect associations. However, the taxa that were not left-censored remained reliably-measured even at the shallowest total read depths ( $r > 0.9$ ) for  $||N|| \geq 100$ .

We were curious as to whether read depth became more important when the ratio of structural to sampling zeros was lowered. To this end, we repeated this simulation with a distribution of prevalences centered on 0.5, so that  $\text{logit}(z_m) \sim \mathcal{N}(0, 1)$ . Since it would probably not be realistic for nearly every microbe in the MIDAS database to have a prevalence above 5%, we lowered the number of taxa  $||M||$  to 1000. Indeed, under these assumptions, read depth became significantly more important (S11 Fig). This is not surprising, since with these assumptions true structural zeros should be less common than above. With these assumptions, censoring does not vary substantially across conditions (censored percentage  $\leq 1\%$ ), but at least 50,000 reads per sample, on average, were necessary to reach a correlation of  $r > 0.9$  for those microbes whose prevalences were above the limit of detection.

Taken together, these simulations suggest that read depth, sample size, and true prevalence values affect the accuracy of prevalence estimates. Read depth becomes most important when

most true prevalences are relatively large (5 – 95%). In contrast, when most prevalences are low, sample size is likely the most important consideration. Finally, when the sample size is sufficiently high ( $|N| > 100$ ) and most prevalences are low, microbes with prevalences above the limit of detection tend to be estimated accurately. These results provide guidelines for when the methods in this manuscript will be most accurate if applied to metagenomes from environments with different characteristics from human stool samples.

### Enrichment analysis

Enrichment analysis was performed using SEED subsystem annotations for FIGfams [88, 37]. Each subsystem was tested for a significant overlap with significant hits from the linear models ( $q \leq 0.05$ ), given the set of FIGfams tested, by Fisher’s exact test. For each gene set, a  $2 \times 2$  contingency table was constructed with the following form:

$$\begin{bmatrix} \|( \text{subsys} \cap \text{signif} ) \cap \text{BG} \| & \|( \text{subsys} \setminus \text{signif} ) \cap \text{BG} \| \\ \|( \text{signif} \setminus \text{subsys} ) \cap \text{BG} \| & \|\text{BG} \setminus ( \text{subsys} \cup \text{signif} ) \| \end{bmatrix}$$

where “subsys” is the set of FIGfams in a given SEED subsystem, “signif” is the set of FIGfams in a particular phylum that were significant hits, and “BG” is the set of all FIGfams tested in that phylum. Two-tailed  $p$ -values were corrected using the Benjamini-Hochberg procedure [89] and an FDR of 25% was set for detecting significant enrichment and depletion (only enrichment is reported). We used this significance threshold in accordance with accepted practice for gene set enrichment analysis [90]. We used the Benjamini-Hochberg procedure since unlike the  $q$ -value method it does not require the estimation of the proportion of true nulls, which is more difficult with small numbers of tests.

### Overlap with *in vivo* results

Results of the screen were obtained from the Supplemental Material of Wu et al. (downloaded on 2017 May 3) [48]. Genes were mapped to FIGfams by matching identifiers in the Supplemental Material to genome annotations from PATRIC [67]. Significance of overlap between these genes and the results for the Bacteroidetes phylum from the body-site-specific or overall models was determined by Fisher’s exact test.

This test depends both on how we determine which genes from the *in vivo* screen count as true positives, and on the choice of the “background set,” i.e., which genes would be possible to find in the *in vivo* study. Rather than committing to one method of picking the “true positive” and “background” sets, we instead enumerated several possibilities, performed all possible combinations (S3 Table), and corrected for multiple comparisons. The options we tested for true positive sets were 1. genes in the screen that were significantly associated with fitness in all four *Bacteroides* strains tested, 2. *Bacteroides thetaiotaomicron* genes significantly associated with diet-independent fitness effects, and 3. *B. thetaiotaomicron* genes associated with either diet-dependent or -independent effects. The background sets we tested were 1. all gene families for which a phylogenetic model was fit, 2. all gene families appearing at least once in a *Bacteroides* genome cluster pangenome, 3. all gene families present in all *Bacteroides* pangenomes, 4. gene families present in some but not all *Bacteroides* pangenomes, and 5. gene families present in *Bacteroides thetaiotaomicron*. Similarly to our approach to gene set enrichment analysis, for each test we assembled a  $2 \times 2$  contingency table as follows:

$$\begin{bmatrix} \|\{ \text{pos} \cap \text{signif} \} \cap \text{BG} \| & \|\{ \text{pos} \setminus \text{signif} \} \cap \text{BG} \| \\ \|\{ \text{signif} \setminus \text{pos} \} \cap \text{BG} \| & \|\text{BG} \setminus \{ \text{pos} \cup \text{signif} \} \| \end{bmatrix}$$



where “pos” refers to the true positive FIGfam set, “signif” refers to the set of significant FIGfam hits from the phylogenetic model, and “BG” refers to the background FIGfam set. The full results are depicted in [S3 Table](#), with the results for true positive set 1 excerpted from this full comparison in [Table 1](#).

## Codebase

The analyses were carried out using R [91]. The code used to perform these analyses is available at <http://www.bitbucket.com/pbradz/plr> in the form of an Rmarkdown notebook [92].

Other R packages used include `phylo1m` [76], `phyloseq` [93], `rentrez` [65], `pbapply` [94], `XML` [95], `qvalue` [77], `magrittr` [96], `compiler` [91], `phytools` [97], `ggtree` [98], `gridExtra` [99], `MASS` [100], `nlme` [80], `fitdistrplus` [81], `truncnorm` [101], `ape` [102], `geiger` [103], `MCMCpack` [104], `knitr` [105], `data.table` [106], `parallel` [91], and `dplyr` [107].

## Supporting information

**S1 Appendix. Glossary of terms.** Reference material giving definitions of mathematical symbols used in the Methods.

(PDF)

**S1 Table. Species prevalences, gut specificities, and Crohn’s disease specificities for all genome clusters (species) tested.** `logit.Prevalence`, `logit.BodySite`, and `logit.Crohns` column titles refer to our estimates of  $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ ,  $\vec{\phi}_{\text{Gut},E}^{\text{Spec}}(A)$ , and  $\vec{\phi}_{\text{CD},E}^{\text{Spec}}(A)$ , respectively. Row labels (5-digit numbers) correspond to MIDAS taxon IDs.

(CSV)

**S2 Table. Concordance of  $\beta_{1,g}$  estimates for weighted least squares vs. Laplace shrinkage procedures.** Pearson’s correlation coefficient  $r$  comparing estimates of  $\beta_{1,g}$  for the weighted phylogenetic least squares procedure with the unweighted phylogenetic least squares using Laplace-shrunk estimates of specificity scores were computed across: all tested genes (all), all genes significant in either the weighted or shrunk phylogenetic model (one significant), or all genes significant in both models (both significant). Comparisons were performed for both body site environmental specificity scores (bodysite) and Crohn’s disease (Crohn’s). Additionally, the number of taxa not shrunk back to the prior by Laplace shrinkage for each environmental specificity score are given (non-shrunk).

(XLSX)

**S3 Table. Full assessment of whether genes linked to microbial fitness in an *in vivo* experiment [48] were enriched for significant hits of the body site-specific and overall gut prevalence models.** The different sets of true positives were defined as: “Bacteroides” (genes in the screen significantly associated with fitness in all four strains), “BthetaDietIndep” (genes present in *Bacteroides thetaiotaomicron* that had diet-independent fitness effects in the screen), and “BthetaAny” (same, but for diet-dependent as well as -independent effects). The “background sets” were defined as follows: “all tested” (all gene families for which a phylogenetic model was fit), “Bacteroides (core or variable)” (all gene families with at least one representative in *Bacteroides* genome cluster pangenomes), “Bacteroides (core only)” (gene families that were present in all *Bacteroides* genome cluster pangenomes), “Bacteroides (variable only)” (gene families present in some but not all *Bacteroides* genomes clusters), and “Bacteroides thetaiotaomicron only” (only gene families present in *Bacteroides thetaiotaomicron*). Two false discovery rates for each model were tested (5% and 25%). Fisher tests yielded  $p$ -values that

were then converted to  $q$ -values using the Benjamini-Hochberg approach [89].  
(XLSX)

**S4 Table. SRA accession IDs used to estimate prevalence and environmental specificity scores.**

(CSV)

**S1 Fig. Estimates of logit-gut prevalence (x-axis) vs. logit-gut environmental specificity score (y-axis), showing only modest correlation.**

(PDF)

**S2 Fig. Method overview.** Using MIDAS, we calculate species relative abundances from shotgun sequencing data. These are binarized to yield a matrix of microbial presence/absence, with rows corresponding to microbes and columns corresponding to samples. Samples are organized into environments (i.e., the environments from which the sample was collected) and into datasets (corresponding to samples collected as part of the same project). Using the presence/absence matrix together with these metadata, we estimate phenotype vectors  $\vec{\phi}(A)$ , whose elements are estimates of microbial phenotypes. These phenotypes fall into two groups: prevalence ( $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ ) and environmental specificity scores ( $\vec{\phi}_{x,E}^{\text{Spec}}(A)$ ). Separately, we use the whole genomes incorporated into the MIDAS database to assemble a matrix of gene presence/absence in the pangenome of microbes, and to construct a phylogenetic species tree based on previously-validated single-copy marker genes. We subset this tree to yield four phylum-specific trees. The inputs to our phylogenetic models are a phenotype vector, a gene presence-absence vector, and a phylogenetic tree. Based on these models, we estimate  $p$ -values for a non-zero effect of the gene on the phenotype, then convert these  $p$ -values into  $q$ -values to obtain predicted gene-phenotype interactions at a given false discovery rate (here, 5%).

(PDF)

**S3 Fig. Laplacian regularization reduces noise in estimating  $\hat{s}_{m,CD,E}(A)$ .** Two species are compared, one that was infrequently observed in both Crohn's disease cases and controls (*Bacillus subtilis*, left) and one with a significant bias for Crohn's disease cases (*Bacteroides fragilis*, right). A) Total counts across subjects for *Bacillus subtilis* and *Bacteroides fragilis*. B) Likelihood function for  $\hat{s}_{m,CD,E}(A)$ , or prevalence in Crohn's disease. The maximum-likelihood value is given in the inset. C) Unregularized likelihood for  $\text{logit}(\hat{s}_{m,CD,E}(A))$ , or the environmental specificity of the microbe. Note that the maximum-likelihood value (inset) was actually almost twice as large for *Bacillus subtilis* as for *Bacteroides fragilis* despite the relative paucity of data for *B. subtilis* (compare Y-axes, which show that the distribution for *B. subtilis* is flatter). D) Laplace prior around  $P(e_{CD}) = 0.002$  with width parameter  $b = 0.16$  (optimized using simulation). E) Log-likelihood plot for the posterior  $P(e_{CD}|m) = \hat{s}_{m,CD,E}^{\text{MAP}}(A)$ , obtained by taking the product of the prior distribution and the unregularized distribution. The maximum *a posteriori* (MAP) estimates are the modes of these distributions (inset).

(PDF)

**S4 Fig. Sensitivity plot for  $\epsilon$  tolerance parameter in Laplace shrinkage.** Y-axis gives the best  $b_{\text{optim}}$  value obtained given a particular  $\log_{10}(\epsilon)$  selected when performing Laplace shrinkage of  $\hat{s}_{m,CD,E}^{\text{MAP}}(A)$  estimates. The value of  $\epsilon$  used in the manuscript (0.005) is highlighted with a vertical dashed line.

(PDF)

**S5 Fig. Illustration showing logit-prevalence vs. the pattern of glutamate-GABA decarboxylase (*gadC*) inheritance in Bacteroidetes.** As in Fig 2, the tree on the left is colored by species

prevalence (black to orange), while the tree on the right is colored by gene presence-absence (blue to black), with selected species called out in the middle, and lines linking species labels to leaves that match leaf color.

(PDF)

**S6 Fig. Distribution of estimated logit-prevalence  $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ , showing impact of censoring.**

A) Density of estimated logit-prevalence distribution,  $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ , showing pile-up of values at the limit of detection  $K$ . B) Density of a normal distribution with mean and standard deviation obtained from best-fit of truncated normal to  $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ . C) “Uncensored” version of  $\vec{\phi}_{x,E,D}^{\text{Prev}}(A)$ . Data points at or below  $K$  have been replaced by random sampling from a truncated normal, with mean and standard deviation as in B and with  $K$  as upper truncation point.

(PDF)

**S7 Fig. Impact of left-censoring on the false positive rate and power of phylogenetic tests.**

Bars give replicate measurements of false positive rate (left, effect size of 0) and power (right, effect size of 0.75) across the different phyla (colors), based on simulating binary genotypes and continuous phenotypes as in Methods, “Power analysis”, with varying levels of left-censoring (“censoring point”), and obtaining  $p$ -values with the three methods described in Methods, “Assessing the potential impact of sampling with left-censoring.” Horizontal dashed lines give a rate of 0.05. Binary genotypes had varying levels of Ives-Garland  $\alpha$  (0, 25, 50), representing high to low phylogenetic signal.

(PDF)

**S8 Fig. Simulations showing the prevalence ratios  $F$  corresponding to various effect sizes.**

To give a more intuitive sense of scale for simulated effect sizes, simulations were performed as in Methods, “Power analysis”, with effect sizes  $\beta_1$  ranging from 0 to 1.25. After fitting phylogenetic models to the simulated phenotypes and genotypes, the average prevalences with the simulated gene,  $\text{logistic}(\beta_{1,g} + \beta_{0,g})$ , and without,  $\text{logistic}(\beta_{0,g})$ , were computed, and their ratio  $F$  was taken.  $\log_2(F)$  is plotted here, such that a value of 1 means the gene conferred (on average) a 2-fold change in prevalence. Violin plots were made of 50 simulations.

(PDF)

**S9 Fig. Results of simulations illustrating compositionality artifacts in relative abundances, but not prevalence estimates.**

Lines are density plots of all pairwise correlations (either over all samples or over only non-zero values) for (A), the relative abundance of microbes within a simulation, and (B), the prevalence of simulated microbes across simulations with the same  $\mu_m$  (mean) and  $z_m$  (zero-inflation) values. A) Pairwise correlations for absolute abundances  $L(B_i^{\text{Abs}})$  (black), column-normalized relative abundances  $L(B_i^{\text{relCN}})$  (blue), and relative abundances obtained through Dirichlet-Multinomial sampling  $L(B_i^{\text{relDM}})$  are shown, along with pairwise correlations over non-zero elements of relative abundances obtained through Dirichlet-Multinomial sampling  $LZ(B_i^{\text{relDM}})$ . B) Pairwise correlations for true prevalences  $L(\Pi_i^{\text{Abs}})$ , for true prevalences above a sampling floor  $L(\Pi_i^{\text{AbsThresh}})$ , and for prevalences calculated from relative abundances obtained through Dirichlet-Multinomial sampling  $L(\Pi_i^{\text{relDM}})$  are compared.

(PDF)

**S10 Fig. Reliability of prevalence estimates, assuming that most microbes are low-prevalence in the given environment.**

We simulated microbiome data using Dirichlet-Multinomial sampling combined with zero inflation. Here, we assumed that true prevalences were distributed  $\text{logit}(z_m) \sim \mathcal{N}(-4, 2)$ . We then plotted (blue contours) the true logit-prevalence  $z_m$  (x-

axis) versus the calculated logit-prevalence from the simulation  $\text{logit}(\hat{p}_{m,N})$ , for varying geometric-mean read depths  $r_n$  (columns) and sample sizes  $||N||$  (rows). The values given for each sample size and read depth are the Pearson correlation  $r$ , the Pearson correlation over non-censored microbes only  $r_{\text{unc}}$ , and the percent of microbes with censored prevalences (“cens”). (PDF)

**S11 Fig. Reliability of prevalence estimates, assuming that microbial prevalence is centered on 50% in the given environment.** As in S10 Fig, but with  $\text{logit}(z_m) \sim \mathcal{N}(0, 1)$ . (PDF)

## Acknowledgments

The authors would like to thank Joshua Ladau, Nandita Garud, and other members of the Pollard and Turnbaugh groups, attendees of the 2017 Keystone meeting on the Microbiome in Health and Disease, and attendees of the Second Workshop in Statistics and Algorithmic Challenges in Microbiome Data Analysis (SACMDA2) for helpful suggestions and discussions.

## Author Contributions

**Conceptualization:** Patrick H. Bradley, Katherine S. Pollard.

**Data curation:** Patrick H. Bradley, Stephen Nayfach.

**Formal analysis:** Patrick H. Bradley.

**Funding acquisition:** Katherine S. Pollard.

**Investigation:** Patrick H. Bradley.

**Methodology:** Patrick H. Bradley, Katherine S. Pollard.

**Project administration:** Katherine S. Pollard.

**Resources:** Stephen Nayfach.

**Software:** Patrick H. Bradley, Stephen Nayfach.

**Supervision:** Katherine S. Pollard.

**Validation:** Patrick H. Bradley.

**Visualization:** Patrick H. Bradley.

**Writing – original draft:** Patrick H. Bradley, Katherine S. Pollard.

**Writing – review & editing:** Patrick H. Bradley, Stephen Nayfach, Katherine S. Pollard.

## References

- Slack E, Hapfelmeier S, Stecher B, Velykoredko Y, Stoel M, Lawson MAE, et al. Innate and adaptive immunity cooperate flexibly to maintain host-microbiota mutualism. *Science*. 2009; 325(5940):617–620. <https://doi.org/10.1126/science.1172747> PMID: 19644121
- Atarashi K, Tanoue T, Shima T, Imaoka A, Kuwahara T, Momose Y, et al. Induction of colonic regulatory T cells by indigenous *Clostridium* species. *Science*. 2011; 331(6015):337–341. <https://doi.org/10.1126/science.1198469> PMID: 21205640
- Mazmanian SK, Round JL, Kasper DL. A microbial symbiosis factor prevents intestinal inflammatory disease. *Nature*. 2008; 453(7195):620–625. <https://doi.org/10.1038/nature07008> PMID: 18509436
- Sassone-Corsi M, Raffatellu M. No vacancy: how beneficial microbes cooperate with immunity to provide colonization resistance to pathogens. *Journal of Immunology*. 2015; 194(9):4081–7. <https://doi.org/10.4049/jimmunol.1403169>

5. Yano JM, Yu K, Donaldson GP, Shastri GG, Ann P, Ma L, et al. Indigenous bacteria from the gut microbiota regulate host serotonin biosynthesis. *Cell*. 2015; 161(2):264–76. <https://doi.org/10.1016/j.cell.2015.02.047> PMID: 25860609
6. Peng L, Li ZR, Green RS, Holzman IR, Lin J. Butyrate enhances the intestinal barrier by facilitating tight junction assembly via activation of AMP-activated protein kinase in Caco-2 cell monolayers. *Journal of Nutrition*. 2009; 139(9):1619–1625. <https://doi.org/10.3945/jn.109.104638> PMID: 19625695
7. Reber SO, Siebler PH, Donner NC, Morton JT, Smith DG, Kopelman JM, et al. Immunization with a heat-killed preparation of the environmental bacterium *Mycobacterium vaccae* promotes stress resilience in mice. *Proceedings of the National Academy of Sciences*. 2016; 113(22):E3130–E3139. <https://doi.org/10.1073/pnas.1600324113>
8. Garrett WS, Gallini CA, Yatsunenko T, Michaud M, DuBois A, Delaney ML, et al. *Enterobacteriaceae* act in concert with the gut microbiota to induce spontaneous and maternally transmitted colitis. *Cell Host & Microbe*. 2010; 8(3):292–300. <https://doi.org/10.1016/j.chom.2010.08.004>
9. Kullberg MC, Ward JM, Gorelick PL, Caspar P, Hieny S, Cheever A, et al. *Helicobacter hepaticus* triggers colitis in specific-pathogen-free interleukin-10 (IL-10)-deficient mice through an IL-12- and gamma interferon-dependent mechanism. *Infection and Immunity*. 1998; 66(11):5157–66. PMID: 9784517
10. Kostic A, Chun E, Robertson L, Glickman J, Gallini C, Michaud M, et al. *Fusobacterium nucleatum* potentiates intestinal tumorigenesis and modulates the tumor-immune microenvironment. *Cell Host & Microbe*. 2013; 14(2):207–215. <https://doi.org/10.1016/j.chom.2013.07.007>
11. Bartlett JG. *Clostridium difficile*-associated enteric disease. *Current Infectious Disease Reports*. 2002; 4(6):477–483. <https://doi.org/10.1007/s11908-002-0032-0> PMID: 12433321
12. van Nood E, Vrieze A, Nieuwdorp M, Fuentes S, Zoetendal EG, de Vos WM, et al. Duodenal infusion of donor feces for recurrent *Clostridium difficile*. *New England Journal of Medicine*. 2013; 368(5):407–415. <https://doi.org/10.1056/NEJMoa1205037> PMID: 23323867
13. Weingarden A, González A, Vázquez-Baeza Y, Weiss S, Humphry G, Berg-Lyons D, et al. Dynamic changes in short- and long-term bacterial composition following fecal microbiota transplantation for recurrent *Clostridium difficile* infection. *Microbiome*. 2015; 3(1):10. <https://doi.org/10.1186/s40168-015-0070-0> PMID: 25825673
14. Khanna S, Vazquez-Baeza Y, González A, Weiss S, Schmidt B, Muñoz-Pedrogo DA, et al. Changes in microbial ecology after fecal microbiota transplantation for recurrent *C. difficile* infection affected by underlying inflammatory bowel disease. *Microbiome*. 2017; 5(1):55. <https://doi.org/10.1186/s40168-017-0269-3> PMID: 28506317
15. Carvalho F, Koren O, Goodrich J, Johansson MV, Nalbantoglu I, Aitken J, et al. Transient inability to manage Proteobacteria promotes chronic gut inflammation in TLR5-deficient mice. *Cell Host & Microbe*. 2012; 12(2):139–152. <https://doi.org/10.1016/j.chom.2012.07.004>
16. Chassaing B, Koren O, Carvalho FA, Ley RE, Gewirtz AT. AIEC pathobiont instigates chronic colitis in susceptible hosts by altering microbiota composition. *Gut*. 2014; 63(7):1069–1080. <https://doi.org/10.1136/gutjnl-2013-304909> PMID: 23896971
17. Vilhjálmsson BJ, Nordborg M. The nature of confounding in genome-wide association studies. *Nature Reviews Genetics*. 2012; 14(1):1–2. <https://doi.org/10.1038/nrg3382> PMID: 23165185
18. Kim PJ, Price ND. Genetic co-occurrence network across sequenced microbes. *PLoS Computational Biology*. 2011; 7(12):e1002340. <https://doi.org/10.1371/journal.pcbi.1002340> PMID: 22219725
19. Porter SS, Chang PL, Conow CA, Dunham JP, Friesen ML. Association mapping reveals novel serpentine adaptation gene clusters in a population of symbiotic *Mesorhizobium*. *The ISME Journal*. 2017; 11(1):248–262. <https://doi.org/10.1038/ismej.2016.88> PMID: 27420027
20. Collins C, Didelot X. A phylogenetic method to perform genome-wide association studies in microbes that accounts for population structure and recombination. *PLoS Computational Biology*. 2018; 14(2):e1005958. <https://doi.org/10.1371/journal.pcbi.1005958> PMID: 29401456
21. Zhao N, Chen J, Carroll I, Ringel-Kulka T, Epstein M, Zhou H, et al. Testing in microbiome-profiling studies with MiRKAT, the microbiome regression-based kernel association test. *The American Journal of Human Genetics*. 2015; 96(5):797–807. <https://doi.org/10.1016/j.ajhg.2015.04.003> PMID: 25957468
22. Silverman JD, Washburne AD, Mukherjee S, David LA. A phylogenetic transform enhances analysis of compositional microbiota data. *eLife*. 2017; 6. <https://doi.org/10.7554/eLife.21887>
23. Tang ZZ, Chen G, Alekseyenko AV, Li H. A general framework for association analysis of microbial communities on a taxonomic tree. *Bioinformatics*. 2016; 33(9):btw804. <https://doi.org/10.1093/bioinformatics/btw804>

24. Lozupone CA, Hamady M, Cantarel BL, Coutinho PM, Henrissat B, Gordon JL, et al. The convergence of carbohydrate active gene repertoires in human gut microbes. *Proceedings of the National Academy of Sciences*. 2008; 105(39):15076–15081. <https://doi.org/10.1073/pnas.0807339105>
25. Felsenstein J. Phylogenies and the comparative method. *The American Naturalist*. 1985;. <https://doi.org/10.1086/284325>
26. Grafen A. The phylogenetic regression. *Philosophical Transactions of the Royal Society of London Series B, Biological Sciences*. 1989; 326(1233):119–57. <https://doi.org/10.1098/rstb.1989.0106> PMID: 2575770
27. Dunn CW, Zapata F, Munro C, Siebert S, Hejnal A. Pairwise comparisons across species are problematic when analyzing functional genomic data. *Proceedings of the National Academy of Sciences*. 2018; 115(3):E409–E417. <https://doi.org/10.1073/pnas.1707515115>
28. Levy A, Salas Gonzalez I, Mittelviehhaus M, Clingenpeel S, Herrera Paredes S, Miao J, et al. Genomic features of bacterial adaptation to plants. *Nature Genetics*. 2018; 50(1):138–150. <https://doi.org/10.1038/s41588-017-0012-9> PMID: 29255260
29. Ord TJ, Martins EP. Tracing the origins of signal diversity in anole lizards: phylogenetic approaches to inferring the evolution of complex behaviour. *Animal Behaviour*. 2006; 71(6):1411–1429. <https://doi.org/10.1016/j.anbehav.2005.12.003>
30. Zaneveld JRR, Parfrey LW, Van Treuren W, Lozupone C, Clemente JC, Knights D, et al. Combined phylogenetic and genomic approaches for the high-throughput study of microbial habitat adaptation. *Trends in Microbiology*. 2011; 19(10):472–82. <https://doi.org/10.1016/j.tim.2011.07.006> PMID: 21872475
31. Washburne AD, Morton JT, Sanders J, McDonald D, Zhu Q, Oliverio AM, et al. Methods for phylogenetic analysis of microbiome data. *Nature Microbiology*. 2018; 3(6):652–661. <https://doi.org/10.1038/s41564-018-0156-0> PMID: 29795540
32. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome Research*. 2016; 26(11):1612–1625. <https://doi.org/10.1101/gr.201863.115> PMID: 27803195
33. Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, et al. A core gut microbiome in obese and lean twins. *Nature*. 2009; 457(7228):480–4. <https://doi.org/10.1038/nature07540> PMID: 19043404
34. Ives AR, Garland T. Phylogenetic logistic regression for binary dependent variables. *Systematic Biology*. 2010; 59(1):9–26. <https://doi.org/10.1093/sysbio/syp074> PMID: 20525617
35. Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nature*. 2012; 486(7402):207–14. <https://doi.org/10.1038/nature11234> PMID: 22699609
36. Gevers D, Kugathasan S, Denson L, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host & Microbe*. 2014; 15(3):382–392. <https://doi.org/10.1016/j.chom.2014.02.005>
37. Meyer F, Overbeek R, Rodriguez A. FIGfams: yet another set of protein families. *Nucleic Acids Research*. 2009; 37(20):6643–54. <https://doi.org/10.1093/nar/gkp698> PMID: 19762480
38. Sakamoto M, Ohkuma M. *Bacteroides reticulotermis* sp. nov., isolated from the gut of a subterranean termite (*Reticulitermes speratus*). *International Journal of Systematic and Evolutionary Microbiology*. 2013; 63(Pt 2):691–695. <https://doi.org/10.1099/ijs.0.040931-0> PMID: 22544795
39. Browne HP, Forster SC, Anonye BO, Kumar N, Neville BA, Stares MD, et al. Culturing of “unculturable” human microbiota reveals novel taxa and extensive sporulation. *Nature*. 2016; 533(7604):543–546. <https://doi.org/10.1038/nature17645> PMID: 27144353
40. Swick MC, Koehler TM, Driks A. Surviving between hosts: sporulation and transmission. *Microbiology Spectrum*. 2016; 4(4). <https://doi.org/10.1128/microbiolspec.VMBF-0029-2015> PMID: 27726794
41. De Biase D, Pennacchiotti E. Glutamate decarboxylase-dependent acid resistance in orally acquired bacteria: function, distribution and biomedical implications of the *gadBC* operon. *Molecular Microbiology*. 2012; 86(4):770–786. <https://doi.org/10.1111/mmi.12020> PMID: 22995042
42. Srinivasa Rao PS, Lim TM, Leung KY. Functional genomics approach to the identification of virulence genes involved in *Edwardsiella tarda* pathogenesis. *Infection and Immunity*. 2003; 71(3):1343–51. <https://doi.org/10.1128/IAI.71.3.1343-1351.2003> PMID: 12595451
43. Cotter PD, Gahan CG, Hill C. A glutamate decarboxylase system protects *Listeria monocytogenes* in gastric fluid. *Molecular Microbiology*. 2001; 40(2):465–75. <https://doi.org/10.1046/j.1365-2958.2001.02398.x> PMID: 11309128
44. Wargo MJ, Meadows JA. Carnitine in bacterial physiology and metabolism. *Microbiology*. 2015; 161(6):1161–1174. <https://doi.org/10.1099/mic.0.000080> PMID: 25787873

45. Staley C, Weingarden AR, Khoruts A, Sadowsky MJ. Interaction of gut microbiota with bile acid metabolism and its influence on disease states. *Applied Microbiology and Biotechnology*. 2017; 101(1):47–64. <https://doi.org/10.1007/s00253-016-8006-6> PMID: 27888332
46. Marques JC, Oh IK, Ly DC, Lamosa P, Ventura MR, Miller ST, et al. LsrF, a coenzyme A-dependent thiolase, catalyzes the terminal step in processing the quorum sensing signal autoinducer-2. *Proceedings of the National Academy of Sciences*. 2014; 111(39):14235–14240. <https://doi.org/10.1073/pnas.1408691111>
47. Thompson J, Oliveira R, Djukovic A, Ubeda C, Xavier K. Manipulation of the quorum sensing signal ai-2 affects the antibiotic-treated gut microbiota. *Cell Reports*. 2015; 10(11):1861–1871. <https://doi.org/10.1016/j.celrep.2015.02.049> PMID: 25801025
48. Wu M, McNulty NP, Rodionov DA, Khoroshkin MS, Griffin NW, Cheng J, et al. Genetic determinants of in vivo fitness and diet responsiveness in multiple human gut *Bacteroides*. *Science*. 2015; 350(6256):aac5992–aac5992. <https://doi.org/10.1126/science.aac5992> PMID: 26430127
49. Wang J, Jia H. Metagenome-wide association studies: fine-mining the microbiome. *Nature Reviews Microbiology*. 2016; 14(8):508–522. <https://doi.org/10.1038/nrmicro.2016.83> PMID: 27396567
50. Shin NR, Whon TW, Bae JW. Proteobacteria: microbial signature of dysbiosis in gut microbiota. *Trends in Biotechnology*. 2015; 33(9):496–503. <https://doi.org/10.1016/j.tibtech.2015.06.011> PMID: 26210164
51. Lynch SV, Pedersen O. The human intestinal microbiome in health and disease. *The New England Journal of Medicine*. 2016; 375(24):2369–2379. <https://doi.org/10.1056/NEJMra1600266> PMID: 27974040
52. Nielsen HB, Almeida M, Juncker AS, Rasmussen S, Li J, Sunagawa S, et al. Identification and assembly of genomes and genetic elements in complex metagenomic samples without using reference genomes. *Nature Biotechnology*. 2014; 32(8):822–828. <https://doi.org/10.1038/nbt.2939> PMID: 24997787
53. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. An integrated catalog of reference genes in the human gut microbiome. *Nature Biotechnology*. 2014; 32(8):834–841. <https://doi.org/10.1038/nbt.2942> PMID: 24997786
54. Glasser AL, Boudeau J, Barnich N, Perruchot MH, Colombel JF, Darfeuille-Michaud A. Adherent invasive *Escherichia coli* strains from patients with Crohn's disease survive and replicate within macrophages without inducing host cell death. *Infection and Immunity*. 2001; 69(9):5529–37. <https://doi.org/10.1128/IAI.69.9.5529-5537.2001> PMID: 11500426
55. Barnich N, Boudeau J, Claret L, Darfeuille-Michaud A. Regulatory and functional co-operation of flagella and type 1 pili in adhesive and invasive abilities of AIEC strain LF82 isolated from a patient with Crohn's disease. *Molecular Microbiology*. 2003; 48(3):781–794. <https://doi.org/10.1046/j.1365-2958.2003.03468.x> PMID: 12694621
56. Small CLN, Reid-Yu SA, McPhee JB, Coombes BK. Persistent infection with Crohn's disease-associated adherent-invasive *Escherichia coli* leads to chronic inflammation and intestinal fibrosis. *Nature Communications*. 2013; 4:1957. <https://doi.org/10.1038/ncomms2957> PMID: 23748852
57. Lawley TD, Klimke WA, Gubbins MJ, Frost LS. F factor conjugation is a true type IV secretion system. *FEMS Microbiology Letters*. 2003; 224(1):1–15. [https://doi.org/10.1016/S0378-1097\(03\)00430-0](https://doi.org/10.1016/S0378-1097(03)00430-0) PMID: 12855161
58. Stecher B, Denzler R, Maier L, Bernet F, Sanders MJ, Pickard DJ, et al. Gut inflammation can boost horizontal gene transfer between pathogenic and commensal *Enterobacteriaceae*. *Proceedings of the National Academy of Sciences*. 2012; 109(4):1269–1274. <https://doi.org/10.1073/pnas.1113246109>
59. Ives AR, Helmus MR. Generalized linear mixed models for phylogenetic analyses of community structure. *Ecological Monographs*. 2011; 81(3):511–525. <https://doi.org/10.1890/10-1264.1>
60. Konstantinidis KT, Tiedje JM. Genomic insights that advance the species definition for prokaryotes. *Proceedings of the National Academy of Sciences*. 2005; 102(7):2567–72. <https://doi.org/10.1073/pnas.0409727102>
61. Jain C, Rodríguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High-throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *bioRxiv*. 2017; p. 225342.
62. Rosselló-Móra R, Amann R. Past and future species definitions for Bacteria and Archaea. *Systematic and Applied Microbiology*. 2015; 38(4):209–216. <https://doi.org/10.1016/j.syapm.2015.02.001> PMID: 25747618
63. Cadillo-Quiroz H, Didelot X, Held NL, Herrera A, Darling A, Reno ML, et al. Patterns of gene flow define species of thermophilic archaea. *PLoS Biology*. 2012; 10(2):e1001265. <https://doi.org/10.1371/journal.pbio.1001265> PMID: 22363207

64. Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition. *Proceedings of the National Academy of Sciences*. 2009; 106(45):19126–31. <https://doi.org/10.1073/pnas.0906412106>
65. Winter DJ. rentrez: An R package for the NCBI eUtils API. *PeerJ Preprints*. 2017
66. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*. 2010; 26(19):2460–2461. <https://doi.org/10.1093/bioinformatics/btq461> PMID: 20709691
67. Wattam AR, Abraham D, Dalay O, Disz TL, Driscoll T, Gabbard JL, et al. PATRIC, the bacterial bioinformatics database and analysis resource. *Nucleic Acids Research*. 2014; 42(Database issue):D581–91. <https://doi.org/10.1093/nar/gkt1099> PMID: 24225323
68. Price MN, Dehal PS, Arkin AP. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS ONE*. 2010; 5(3):e9490. <https://doi.org/10.1371/journal.pone.0009490> PMID: 20224823
69. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*. 2004; 32(5):1792–1797. <https://doi.org/10.1093/nar/gkh340> PMID: 15034147
70. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics*. 2004; 20(2):289–290. <https://doi.org/10.1093/bioinformatics/btg412> PMID: 14734327
71. Garamszegi LZ, Gonzalez-Voyer A. Working with the tree of life in comparative studies: how to build and tailor phylogenies to interspecific datasets. In: *Modern phylogenetic comparative methods and their application in evolutionary biology*. Berlin, Heidelberg: Springer Berlin Heidelberg; 2014. p. 19–48. Available from: [http://link.springer.com/10.1007/978-3-662-43550-2\\_2](http://link.springer.com/10.1007/978-3-662-43550-2_2).
72. Karlsson FH, Tremaroli V, Nookaew I, Bergström G, Behre CJ, Fagerberg B, et al. Gut metagenome in European women with normal, impaired and diabetic glucose control. *Nature*. 2013; 498(7452):99–103. <https://doi.org/10.1038/nature12198> PMID: 23719380
73. Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F, et al. A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*. 2012; 490(7418):55–60. <https://doi.org/10.1038/nature11450> PMID: 23023125
74. Zhu Y, Stephens RM, Meltzer PS, Davis SR. SRADB: query and use public next-generation sequencing data from within R. *BMC Bioinformatics*. 2013; 14(1):19. <https://doi.org/10.1186/1471-2105-14-19> PMID: 23323543
75. Leinonen R, Akhtar R, Birney E, Bower L, Cerdeno-Tarraga A, Cheng Y, et al. The European Nucleotide Archive. *Nucleic Acids Research*. 2011; 39(Database):D28–D31. <https://doi.org/10.1093/nar/gkq967> PMID: 20972220
76. si Tung Ho L, Ané C. A linear-time algorithm for Gaussian and non-Gaussian trait evolution models. *Systematic Biology*. 2014; 63(3):397–408. <https://doi.org/10.1093/sysbio/syu005>
77. Storey JD, Bass AJ, Dabney A, Robinson D. qvalue: Q-value estimation for false discovery rate control; 2015. Available from: <http://github.com/jdstorey/qvalue>.
78. Storey JD, Tibshirani R. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences*. 2003; 100(16):9440–5. <https://doi.org/10.1073/pnas.1530509100>
79. Kappelman MD, Rifas-Shiman SL, Kleinman K, Ollendorf D, Bousvaros A, Grand RJ, et al. The prevalence and geographic distribution of Crohn's disease and ulcerative colitis in the United States. *Clinical Gastroenterology and Hepatology*. 2007; 5(12):1424–1429. <https://doi.org/10.1016/j.cgh.2007.07.012> PMID: 17904915
80. Pinheiro J, Bates D, DebRoy S, Sarkar D, R Core Team. nlme: Linear and Nonlinear Mixed Effects Models; 2018. Available from: <https://CRAN.R-project.org/package=nlme>.
81. Delignette-Muller ML, Dutang C. fitdistrplus: An R Package for Fitting Distributions. *Journal of Statistical Software*. 2015; 64(4):1–34. <https://doi.org/10.18637/jss.v064.i04>
82. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are compositional: and this is not optional. *Frontiers in Microbiology*. 2017; 8:2224. <https://doi.org/10.3389/fmicb.2017.02224> PMID: 29187837
83. Fernandes AD, Reid JN, Macklaim JM, McMurrough TA, Edgell DR, Gloor GB. Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis. *Microbiome*. 2014; 2(1):15. <https://doi.org/10.1186/2049-2618-2-15> PMID: 24910773
84. Aitchison J. The statistical analysis of compositional data. *Journal of the Royal Statistical Society, Series B (Methodological)*. 1982; 44(2):139–177.
85. Friedman J, Alm EJ. Inferring correlation networks from genomic survey data. *PLoS Comput Biol*. 2012; 8. <https://doi.org/10.1371/journal.pcbi.1002687>



86. Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, et al. Correlation detection strategies in microbial data sets vary widely in sensitivity and precision. *ISME J*. 2016; 10(7). <https://doi.org/10.1038/ismej.2015.235> PMID: 26905627
87. Xu L, Paterson AD, Turpin W, Xu W. Assessment and selection of competing models for zero-inflated microbiome data. *PLoS ONE*. 2015; 10(7):e0129606. <https://doi.org/10.1371/journal.pone.0129606> PMID: 26148172
88. Overbeek R, Begley T, Butler RM, Choudhuri JV, Chuang HY, Cohoon M, et al. The subsystems approach to genome annotation and its use in the project to annotate 1000 genomes. *Nucleic Acids Research*. 2005; 33(17):5691–5702. <https://doi.org/10.1093/nar/gki866> PMID: 16214803
89. Hochberg Y, Benjamini Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B (Methodological)*. 1995; 1:289–300.
90. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America*. 2005; 102(43):15545–50. <https://doi.org/10.1073/pnas.0506580102> PMID: 16199517
91. R Core Team. R: A language and environment for statistical computing; 2017. Available from: <https://www.R-project.org/>.
92. Allaire J, Xie Y, McPherson J, Luraschi J, Ushey K, Atkins A, et al. rmarkdown: Dynamic Documents for R; 2018. Available from: <https://CRAN.R-project.org/package=rmarkdown>.
93. McMurdie PJ, Holmes S. phyloseq: An R package for reproducible interactive analysis and graphics of microbiome census data. *PLoS ONE*. 2013; 8(4):e61217. <https://doi.org/10.1371/journal.pone.0061217> PMID: 23630581
94. Solymos P, Zawadzki Z. pbapply: Adding Progress Bar to '\*apply' Functions; 2018. Available from: <https://CRAN.R-project.org/package=pbapply>.
95. Lang DT, the CRAN Team. XML: Tools for Parsing and Generating XML Within R and S-Plus; 2018. Available from: <https://CRAN.R-project.org/package=XML>.
96. Bache SM, Wickham H. magrittr: A Forward-Pipe Operator for R; 2014. Available from: <https://CRAN.R-project.org/package=magrittr>.
97. Revell LJ. phytools: An R package for phylogenetic comparative biology (and other things). *Methods in Ecology and Evolution*. 2012; 3:217–223. <https://doi.org/10.1111/j.2041-210X.2011.00169.x>
98. Yu G, Smith D, Zhu H, Guan Y, Lam TTY. ggtree: an R package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods in Ecology and Evolution*. 2017; 8:28–36. <https://doi.org/10.1111/2041-210X.12628>
99. Auguie B. gridExtra: miscellaneous functions for "grid" graphics; 2017. Available from: <https://CRAN.R-project.org/package=gridExtra>.
100. Venables WN, Ripley BD. *Modern Applied Statistics with S*. 4th ed. New York: Springer; 2002. Available from: <http://www.stats.ox.ac.uk/pub/MASS4>.
101. Mersmann O, Trautmann H, Steuer D, Bornkamp B. truncnorm: truncated normal distribution; 2018. Available from: <https://CRAN.R-project.org/package=truncnorm>.
102. Paradis E, Claude J, Strimmer K. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics*. 2004; 20:289–290. <https://doi.org/10.1093/bioinformatics/btg412> PMID: 14734327
103. Harmon LJ, Weir JT, Brock CD, Glor RE, Challenger W. GEIGER: investigating evolutionary radiations. *Bioinformatics*. 2008; 24:129–131. <https://doi.org/10.1093/bioinformatics/btm538> PMID: 18006550
104. Martin AD, Quinn KM, Park JH. MCMCpack: Markov Chain Monte Carlo in R. *Journal of Statistical Software*. 2011; 42(9):22. <https://doi.org/10.18637/jss.v042.i09>
105. Xie Y. knitr: A comprehensive tool for reproducible research in R. In: Stodden V, Leisch F, Peng RD, editors. *Implementing reproducible computational research*. Chapman and Hall/CRC; 2014. Available from: <http://www.crcpress.com/product/isbn/9781466561595>.
106. Dowle M, Srinivasan A. data.table: Extension of 'data.frame'; 2018. Available from: <https://CRAN.R-project.org/package=data.table>.
107. Wickham H, Francois R, Henry L, Müller K. dplyr: a grammar of data manipulation; 2017. Available from: <https://CRAN.R-project.org/package=dplyr>.