

Phylogenetic factorization of compositional data

September 8, 2016

Authors: Alex D Washburne^{1,2,†}, Justin D Silverman^{3,4,5,6}, Jonathan W Leff², Dominic J Bennett^{7,8}, John L. Darcy⁹, Sayan Mukherjee^{3,10}, Noah Fierer², and Lawrence A David^{3,5,6}

¹ Nicholas School of the Environment, Duke University, Durham, NC 27708

² Cooperative Institute for Research in Environmental Sciences (CIRES), University of Colorado, Boulder, Boulder CO 80309

³ Program in Computational Biology and Bioinformatics, Duke University, Durham, NC 27708

⁴ Medical Scientist Training Program, Duke University, Durham, NC 27708

⁵ Center for Genomic and Computational Biology, Duke University, Durham, NC 27708

⁶ Department of Molecular Genetics and Microbiology, Duke University, Durham, NC 27708

⁷Department of Earth Science and Engineering, Imperial College London, London UK

⁸Institute of Zoology, Zoological Society of London, London UK

⁹Department of Ecology and Evolutionary Biology, University of Colorado, Boulder, CO 80309

¹⁰Departments of Statistical Science, Mathematics, and Computer Science, Duke University, Durham, NC 27708

†Corresponding author; contact at alex.d.washburne@gmail.com

1 Abstract

2 Marker gene sequencing of microbial communities has generated big datasets
3 of microbial relative abundances varying across environmental conditions, sam-
4 ple sites and treatments. These data often come with putative phylogenies,
5 providing unique opportunities to investigate how shared evolutionary history
6 affects microbial abundance patterns. Here, we present a method to identify the
7 phylogenetic factors driving patterns in microbial community composition. We
8 use the method, “phylofactorization”, to re-analyze datasets from human body
9 and soil microbial communities, demonstrating how phylofactorization can be
10 a dimensionality-reducing tool, an ordination-visualization tool, and also mass-
11 produce inferences on the edges in the phylogeny in which meaningful differences
12 arose.

13 Background

14 Microbial communities play important roles in human [6], livestock [16] and
15 plant [3] health, biogeochemical cycles [2, 12], the maintenance of ecosystem pro-
16 ductivity, bioremediation, and other ecosystem services. Given the importance
17 of microbial communities and the vast number of uncultured and undescribed
18 microbes associated with animal and plant hosts and in natural and engineered
19 systems, understanding the factors determining microbial community structure
20 and function is major challenge for modern biology.

21 Marker gene sequencing (e.g. 16S rRNA gene sequencing to assess bacterial and
22 archaeal diversity and 18S markers for Eukaryotic diversity) is now one of the
23 most commonly used approaches for describing microbial communities, quanti-
24 fying the relative abundances of individual microbial taxa, and characterizing
25 how microbial communities change across space, time, or in response to known
26 biotic or abiotic gradients.

27 Analyzing these data is challenging due to the peculiar noise structure of sequence-
28 count data [30], the inherently compositional nature of the data [15], deciding
29 the taxonomic scale of investigation [7, 8, 31], and the high-dimensionality of
30 species-rich microbial communities [13]. There is a great need and opportunity
31 to develop tools to more efficiently analyze these datasets and leverage infor-
32 mation on the phylogenetic relationships among taxa to better identify which

33 clades are driving differences in microbial community composition across sample
34 categories or measured biotic or abiotic gradients [24]. In this paper, we take on
35 these challenges by developing a means to perform regression of biotic/abiotic
36 gradients on branches in the phylogenetic tree, allowing dimensionality reduc-
37 tion to a series of branches in the phylogeny in a manner consistent with the
38 compositional nature of the data.

39 Many of these challenges can be resolved by performing regression on clades
40 identified in the phylogeny. Consider a study on the effect of oxazolidinones,
41 which affect gram-positive bacteria, on microbial community composition. Rather
42 than regression of antibiotic treatment on abundance at numerous taxonomic
43 levels, statistical analysis of bacterial communities treated with an oxazolidinone
44 should instantly identify the split between gram-positive and gram-negative bac-
45 teria as the most important phylogenetic factor determining response to oxazo-
46 lidinones. Subsequent factors should then be identified by comparing bacteria
47 within the previously-identified groups: identify clades within gram-positives
48 which may be more resistant or susceptible than the remaining gram-positives.
49 Splitting the phylogeny at each inference and making comparisons within the
50 split groups ensures that subsequent inferences are independent of the gram pos-
51 itive - gram negative split which we have already obtained. All of this analysis
52 must be done consistent with the compositional nature of sequence count data.

53 Here, we provide a method to analyze phylogenetically-structured compositional
54 data. The algorithm, referred to as “phylofactorization”, iteratively identifies
55 the most important clades driving variation in the data through their associa-
56 tions with independent variables. Clades are chosen based on some metric of
57 the strength or importance of their regressions with meta-data, and subsequent
58 clades are chosen by comparison of sub-clades within the previously-identified
59 bins of phylogenetic groups. Each “factor” identified corresponds to an edge in
60 the phylogeny, and phylofactorization builds on literature from compositional
61 data analysis to construct a set of orthogonal axes corresponding to those edges;
62 the output orthonormal basis allows the projection of sequence-count relative
63 abundances onto these phylogenetic axes for dimensionality reduction, visual-
64 ization, and standard multivariate statistical analyses. The visualizations and
65 inferences drawn from phylofactorization can be tied back to splits in a given
66 phylogenetic tree and thereby allow researchers to annotate the microbial phy-
67 logeny from the results of microbiome datasets.

68 We show with simulations that phylofactor is able to correctly identify affected

69 clades. We then phylofactor a dataset of human oral and fecal microbiomes
70 to determine the phylogenetic factors driving variation in human body site [4],
71 and a dataset of soil microbes using a multiple regression of pH, carbon concen-
72 tration and nitrogen concentration [28]. In the human microbiome dataset, we
73 find three splits in the phylogeny that together capture 17.6% of the variation
74 community composition across two body sites. Phylofactorization reveals splits
75 between unclassified OTUs not identifiable by taxonomic grouping, important
76 clades of monophyletic yet para-taxonomic OTUs, and a spectrum of taxo-
77 nomic scales for binning and analyzing taxonomic units that varies across taxa
78 - all features that would be missed by standard taxonomy-based analysis. In
79 the soil microbiome dataset, we use phylofactor-based dimensionality reduction
80 and ordination-visualization - either using the orthogonal axes corresponding to
81 splits in the phylogeny, or binning OTUs based on their inferred phylogenetic
82 factors - to find that pH drives most of the variation in the dominant clades in
83 the soil dataset, and confirm this finding by dominance analysis on the under-
84 lying regressions in phylofactorization, indicating that >90% of the explained
85 variation in the first three factors is explained by pH. The axes in our ordination-
86 visualization plots correspond to identifiable edges on the phylogeny that have
87 clear biological interpretations and can be used and tested across studies. User-
88 friendly code for implementing, summarizing and visualizing phylofactorization
89 is provided in an R package - 'phylofactor'.

90 Results

91 We find three main results. First, we find that our algorithm out-performs a
92 standard tool for analyzing compositions of parts related by a tree - what we
93 refer to as the "rooted ILR" transform - and that we can obtain a conservative
94 estimate of the number of phylogenetic factors in simulated datasets with a
95 known number of affected clades. Second, we phylofactor a dataset of the human
96 oral and fecal microbiomes and find three edges in the phylogeny that account
97 for 17.6% of the variation in microbial communities across these sample sites,
98 edges that are not assigned a unique taxonomic label and are thus invisible to
99 taxonomic-based analyses. Third, we show that phylofactorization can be com-
100 bined with multiple regression to reveal that pH drives the main phylogenetic
101 patterns of community composition in soil microbiomes, and show that in four

102 factors we split the Acidobacteria three times - including one split that identi-
103 fies a monophyletic clade of Acidobacteria that consists of alkaliphiles. Finally,
104 using the soil dataset, we demonstrate how phylofactorization yields two compli-
105 mentary methods for dimensionality reduction and ordination-visualization that
106 tell a simplified story of how the major phylogenetic groups of OTUs change
107 with pH.

108 **Power Analysis and Conservative Stopping of Phylofactor-** 109 **ization**

110 Phylofactorization remedies the structured residuals from the rooted ILR re-
111 gression on data with fold-changes in abundances within clades. Phylofac-
112 torization also remedies the problem of high false-positive rates arising from
113 the nested-dependence and correlated coordinates of the rooted ILR transform,
114 as sequential inferences in phylofactorization are independent. Phylofactoriza-
115 tion out-performs the rooted ILR in identifying the correct clades with a given
116 fold-change in abundance (Figs. 1a and 1b), and can be paired with other
117 algorithms assessing residual structure to stop factorization when there is no
118 residual structure and thus accurately identify the number of affected clades
119 (Fig. 1c). Finally, by focusing the inferences on edges instead of nodes in the
120 phylogeny, this algorithm works on trees with polytomies and doesn't require a
121 forced resolution of polytomies to construct a sequential binary partition of the
122 OTUs. Since edges are the locations of the phylogeny where functional traits
123 arise, the identification of edges that drive variation yields a clear, biological
124 interpretation.

125 **Oral-Fecal Microbiome**

126 Phylofactorization of the oral-fecal microbiome dataset, with 290 OTUs and 40
127 samples, yields three factors that explain 17.6% of the variation in the dataset,
128 factors which correspond to clearly visible blocks in phylogenetic heatmaps of
129 the OTU table (Fig. 1). The factors span a range of taxonomic scales and all of
130 them would be invisible to taxonomic-based analyses. Below, we summarize the
131 factors - the P-values from regression, the taxa split at each factor, the body
132 site associations predicted by generalized linear modeling of the ILR coordinate
133 against body site, and finer detail about the taxonomic identities and known
134 ecology of monophyletic taxa being split. Phylofactorization of these data indi-

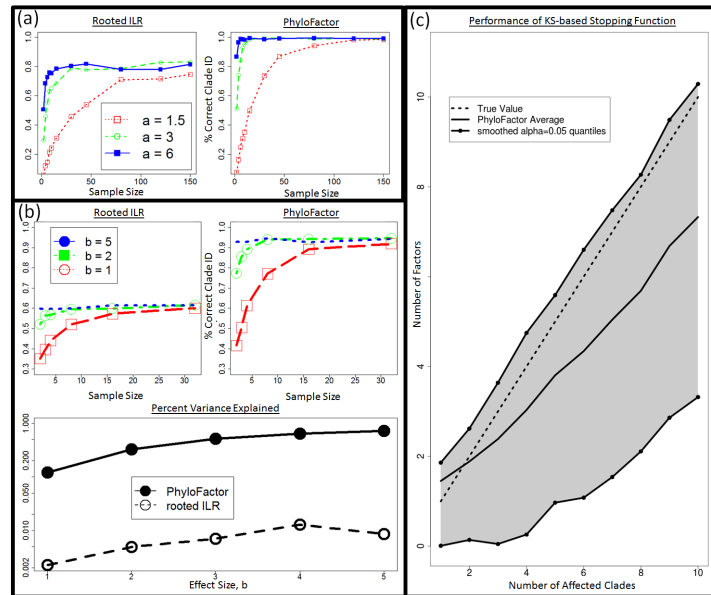


Figure 1: (a) **Power Analysis - 1 Clade.** The rooted ILR transform that minimizes residual variance when regressed against sample site is less able to identify the correct clade compared to phylofactorization for a variety of effect sizes, a , and sample sizes. (b) **Three Significant Clades:** When three significant clades are chosen and given a set of effects increasing in intensity with the parameter b , choosing the top rooted ILR coordinates under performs phylofactorization in correctly identifying the affected clades. Phylofactorization also explains more variation in the data: across effect sizes, phylofactorization explains 2 orders of magnitude more of the variance in the dataset than the sequential rooted ILR. (c) **Stopping Phylofactorization:** Plots of the true number of affected clades in simulated datasets against the number of clades identified by the R package 'phylofactor'. One can terminate phylofactorization when the true number of affected clade is unknown by choosing a stopping function aimed at stopping when there is no evidence of a remaining signal. By stopping the iteration when the distribution of P-values from analyses of variance of regression on candidate ILR basis elements is uniform (specifically, stopping when a KS test against a uniform distribution yields $P > 0.05$), we obtain a conservative estimate of the number of phylogenetic factors in the data.

135 cates that a few clades explain a large fraction of the variation in the data, and
136 many more clades can be identified as containing the same intricate detail as the
137 phylogenetic factors presented below. The biology of microbial human-body-site
138 association can focus on these dominant factors - which traits and evolutionary
139 history drive these monophyletic groups' strong, common association with body
140 sites?

141 The first factor ($P = 4.90 \times 10^{-30}$) split Actinobacteria and Alpha-, Beta-
142 , Gamma-, and Delta-proteobacteria from Epsilonproteobacteria and the rest
143 (Fig. S4). The underlying generalized linear model predicts the Actinobacteria
144 and non-Epsilon-proteobacteria to be 0.4x as abundant as the rest in the gut and
145 3.7x as abundant as the rest in the tongue. The Actinobacteria identified as more
146 abundant in the tongue include four members of the plaque-associated family
147 Actinomycetaceae, one unclassified species of *Cornybacterium*, three members
148 of the mouth-associated genus *Rothia* [20], and one unclassified species of the
149 vaginal-associated genus *Atopobium* [9]. With a standard multivariate analysis
150 of the CLR-transformed data, all nine of these Actinobacteria were identified
151 as significantly more abundant in the tongue from regression of the individual
152 OTUs when using either a 1% false-discovery rate or a Bonferonni correction -
153 these monophyletic taxa all individually show a strong preference for the same
154 body site, and their basal branch was identified as our first phylogenetic fac-
155 tor. The remaining Alpha-, Beta-, Gamma- and Delta-proteobacteria grouped
156 with the Actinobacteria consisted of 31 OTUs, and the Epsilonproteobacteira
157 split from the rest were three unclassified species of the genus *Campylobacter*.
158 The grouping of Actinobacteria with the non-Epsilon Proteobacteria motivates
159 the need for accurate phylogenies in phylofactorization, but also illustrates the
160 promise of identifying clades of interest where the phylogeny is correct and the
161 taxonomy is not.

162 The second factor ($P = 1.15 \times 10^{-31}$) splits 16 Firmicutes of the class Bacilli
163 from the obligately anaerobic Firmicutes class Clostridia and the remaining
164 paraphyletic group containing Epsilonproteobacteria and the rest. The Bacilli
165 are, on average, 0.3x as abundant in the gut as the paraphyletic remaining OTUs
166 and 3.9x as abundant in the tongue. The 16 Bacilli OTUs factored here contain
167 12 unclassified species of the genus *Streptococcus*, well known for its association
168 with the mouth [18], one member of the genus *Lactococcus*, one unclassified
169 species of the mucosal-associated genus *Gemella*, and two members the family
170 Carnobacteriaceae often associated with fish and meat products [22].

171 The third factor ($P = 1.37 \times 10^{-28}$) separated 15 members of the Bacteroidetes
172 family Prevotellaceae from all other Bacteroidetes and the remaining para-
173 phyletic group of OTUs not split by previous factors. The Prevotellaceae split in
174 the third factor were all of the genus *Prevotella*, including the species *Prevotella*
175 *melaninogenica* and *Prevotella nanceiensis* found to have abundances 0.3x as
176 abundant in the gut and 4.0x as abundant in the tongue relative to the other
177 taxa from which they were split.

178 These first three factors capture major blocks visible in the dataset (Fig. 1) can
179 be used as dimensionality reduction tool with a phylogenetic interpretation (Fig.
180 1). While traditional ordination-visualization tools may capture larger fractions
181 of variation of the data, phylogenetic factorization yields a few variables - ratios
182 of clades - which capture large blocks of variation in the data and can be traced
183 to single edges in the phylogeny corresponding to meaningful splits between
184 taxa, edges where traits likely arose which govern the differential abundances
185 across sample sites and environmental gradients or responses to treatments (Fig.
186 1b, supplemental Figs. S4-S8).

187 Using the KS-test stopping criterion, phylofactorization was terminated at 142
188 factors, each corresponding to a branch in the phylogenetic tree separating two
189 groups of OTUs based on their differential abundances in the tongue and fe-
190 ces. These 142 factors define 143 groups, or what we call 'bins', of taxa which
191 remain unsplit by the phylofactorization. The bins vary in size; 112 bins con-
192 tained only single OTUs, whereas 8 were monophyletic clades and the rest are
193 paraphyletic groups of OTUs, the result of taxa within a monophyletic group
194 being factored, yielding one monophyletic group and one paraphyletic group. Of
195 the 112 single-OTU bins extracted from phylofactorization, 78 were also iden-
196 tified as significant at a false-discovery rate of 1%. Some monophyletic bins
197 included groups of unclassified genera that would not be grouped at the genus
198 level under standard taxonomy-based analyses. For instance, two monophyletic
199 clades of the Firmicutes family Lachnospiraceae were identified as having dif-
200 ferent preferred body sites, yet both clades were unclassified at the genus level.
201 Taxonomic-based analyses would either omit these unclassified genera, or group
202 them together and make it difficult to observe a signal due to the two sub-groups
203 having different responses to body site.

204 Performing regression on centered log-ratio (CLR) transformed OTU tables
205 yielded 236 significant OTUs at a false-discovery rate of 1%, and the phylo-
206 genetic signal of these OTUs may be difficult to parse out. However, three

207 iterations of phylofactorization yielded the three major splits in the phylogeny,
208 all of which are consistent with known distributions of taxa. Algorithms such as
209 phylofactor [19], which track P-values up the tree, identify clades with common
210 significance, yet not necessarily clades with common signal - it is a common
211 signal, not a common significance, which better indicates a putative trait driv-
212 ing predictable responses in microbes. In the 142 factors above, phylofactor
213 identified numerous clades with common significance yet different signals.

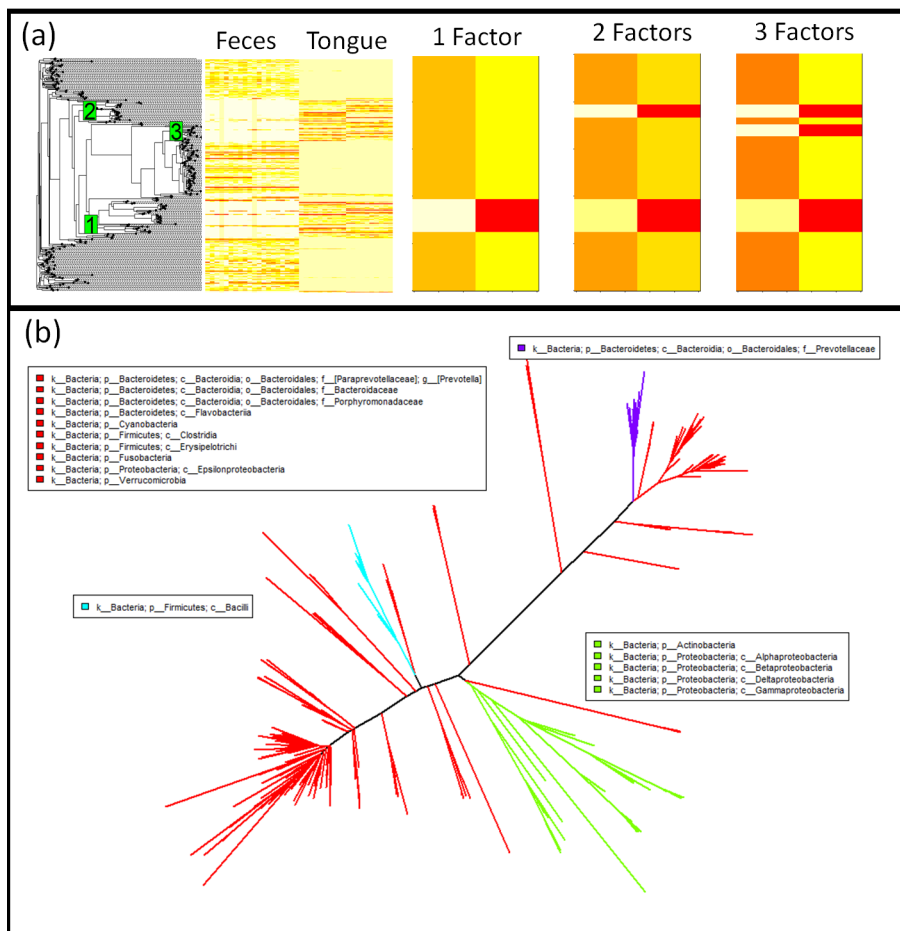


Figure 2: Phylofactorization of human feces/tongue dataset identifies clades differentiating sites. (a) Phylogenetic structure is visible as blocks using a phylogenetic heatmap from the R package 'phytools' [29]. The first factor separates Actinobacteria and some Proteobacteria from the rest, the second factor separates the class Bacilli from the remaining non-Proteobacteria and non-Actinobacteria, the third factor pulls out the genus *Prevotella* from Bacteroidetes and indicates that it, unlike many other taxa in Bacteroidetes, is unrepresented in the tongue. Each factor captures a major block of variation in the data, and the orthogonality of the ILR coordinates from each factor allow multiple factors to be combined easily for estimates of community composition. (b) These three factors splits the phylogeny into four bins. Three of those bins are monophyletic and the final bin is a "remainder" bin, containing taxa split off by the previous monophyletic bins. The three factors are identifiable edges between nodes that can be mapped to an online database containing those nodes.

214 Soil Microbiome

215 The soil microbiome dataset was much larger - 3,379 OTUs and 580 samples -
216 and a much smaller fraction of the variation could be explained by phylofactor-
217 ization. Phylofactorization allows meaningful dimensionality reduction by both
218 factors - plots of the ILR coordinates for the dominant factors - and by bins of
219 taxa that remain un-split at a given level of factorization. Phylofactorization
220 confirmed that the pH of the environment plays a dominant role in the micro-
221 bial community composition, consistent with previous analyses based on Mantel
222 tests [28]. Dominance analysis of the generalized linear models associated with
223 each factor determined pH to account for approximately 92.87%, 89.78%, and
224 92.94% of the explained variance in the first, second, and third factor, respec-
225 tively. C and N were relatively unimportant, and the dominance of pH in the
226 first three factors can be visualized by ordination-visualization plots of the ILR
227 coordinates of the first three factors (Fig. 2a).

228 The first factor splits a group of 206 OTUs in two classes of Acidobacteria from
229 all other bacteria: class Acidobacteriia and class DA052 are shown to decrease in
230 relative abundance with increasing pH. The second factor split 31 OTUs in the
231 order Actinomycetales (some from the family Thermomonosporaceae and the
232 rest unclassified at the family level) from the remainder of all other bacteria,
233 and these monophyletic Actinomycetales also decrease in relative abundance
234 with increasing pH. The third factor identified another clade within the phylum
235 Acidobacteria to decrease with pH: 115 bacteria from the classes Solibacteres
236 and TM1.

237 Interestingly, the fourth factor identifies a large collection of 193 OTUs in the re-
238 mainder of phylum Acidobacteria (i.e. those Acidobacteria not mentioned above
239 in factors 1 and 3) as having relative abundances that increase with pH (dom-
240 inance analysis: 94.79% of explained variance attributable to pH). Unlike the
241 previous three factors above which were acidophiles, this monophyletic group of
242 Acidobacteria consists of alkaliphiles, which includes the classes Acidobacteria-6,
243 Chloracidobacteria, S053 and three OTUs unclassified at the class level.

244 The first four factors can be used to define 5 bins of OTUs that we refer to as
245 “binned phylogenetic units” or BPUs: a monophyletic group of Acidobacteria
246 (classes Chloracidobacteria, Acidobacteria-6, and S035), another monophyletic
247 group of Acidobacteria (classes Solibacteres and TM1), a monophyletic group of
248 several families of the order Actinomycetales, a monophyletic group of Acidobac-
249 teria (classes Acidobacteriia and DA0522), and a paraphyletic amalgamation

250 of the remaining taxa. Binning the OTUs based on these BPUs tells a simplified
251 story of how pH drives microbial community composition (Fig. 2b).

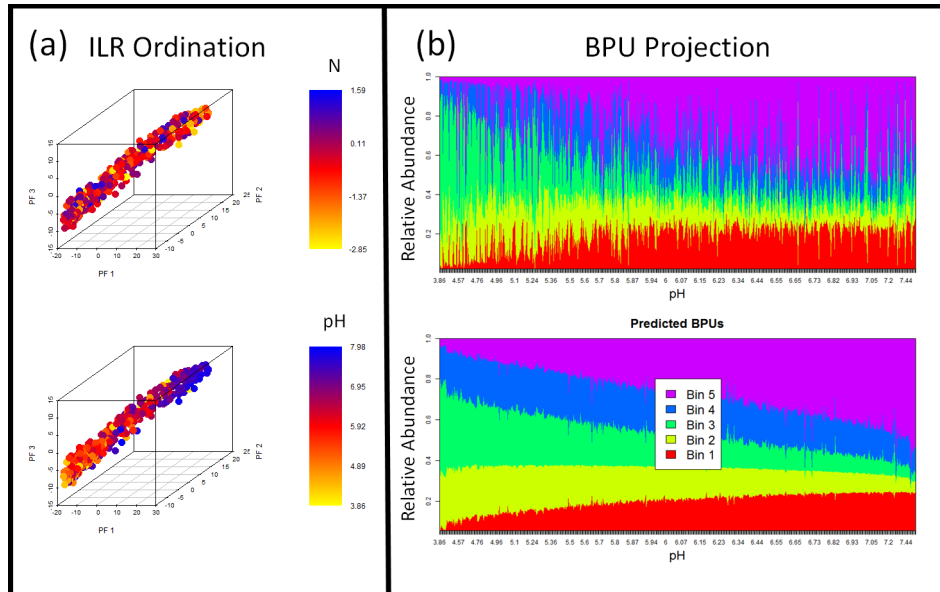


Figure 3: Dimensionality Reduction and Ordination-Visualization of soil microbiome dataset. Phylofactor presents two complementary methods for projecting and visualizing the high-dimensional phylogenetically-structured compositional data. (a) The ILR coordinates have asymptotic normality properties and provide biologically informative ordination-visualization plots. Here, we see that pH is a much better predictor than N of the major phylogenetic factors in Central Park soils. Dominance analysis indicated that pH accounts for approximately 92.87%, 89.78%, and 92.94% of the explained variance in the first, second, and third factor, respectively, consistent with previous results based on Bray-Curtis distances and Mantel tests showing the dominance of pH in structuring soil microbiomes [28]. (b) Every edge separates one group of taxa into two, and those split groups of taxa - what we refer to as bins - can be used to amalgamate taxa and construct a lower-dimensional, compositional dataset of “binned phylogenetic units” (BPUs). Bin 5 is an amalgamation of a monophyletic group of Acidobacteria (classes Chloracidobacteria, Acidobacteria-6, and S035) that increase in relative abundance with pH. Bin 4 is a monophyletic group of Acidobacteria (classes Solibacteres and TM1), Bin 3 is a monophyletic group of several families of the order Actinomycetales, Bin 2 is a monophyletic group of Acidobacteria (classes Acidobacteriia and DA0522), and Bin 1 is a paraphyletic amalgamation of the remaining taxa.

252 Discussion

253 Overview

254 We have introduced a simple and generalizable exploratory data analysis al-
255 gorithm, phylofactorization, to identify clades driving variation in microbiome
256 datasets. Phylofactorization integrates both the compositional and phylogenetic
257 structure of microbiome datasets and produces outputs that contain biological
258 information: effects of independent variables on edges in the phylogeny, includ-
259 ing the tips of the tree traditionally analyzed. The output of phylofactorization
260 contains a sequence of “factors”, or splits in the tree identifying sub-groups of
261 taxa which respond differently to treatment relative to one-another. The splits
262 identified in phylofactorization need not be splits in the Linear taxonomy but
263 can identify strong responses in clades of unclassified taxa. The researcher
264 does not need to choose a taxonomic level at which to perform analysis - those
265 taxonomic levels are output based on whichever clades maximize the objective
266 function, and so researchers will be able to identify multiple taxonomic scales
267 of importance.

268 Phylofactorization outputs an isometric log-ratio transform of the data with
269 known asymptotic normality properties, coordinates that can be analyzed with
270 standard multivariate methods [25]. The resulting coordinates correspond to
271 particular edges between clearly identifiable nodes in the tree of life, allow-
272 ing researchers to annotate a given phylogenetic tree with correlations between
273 clades and various environmental meta-data, sample categories, or experimental
274 treatments.

275 Future Work

276 The generality of phylofactorization opens the door to future work employing
277 phylofactorization with other objective functions. As we showed with the human
278 oral/fecal microbiomes, phylofactorization is not restricted to basal clades, but
279 includes the tips as possible clades of interest, but the objective function we used
280 minimized residual variance in the whole community and thereby may prioritize
281 deeply rooted edges or abundant taxa with weaker effects over individual OTUs
282 with stronger effects. Other objective functions could be constructed to meet
283 the needs of the researcher. If researchers are interested in identifying basal
284 lineages, their objective function can weight edges based on distance from the

285 tips. If researchers are interested in identifying putative traits, they may be
286 interested in an objective function weighting edges based on edge length under
287 an assumption that the probability of a trait arising increases with the amount
288 of time elapsed.

289 Each edge identified in phylofactorization corresponds to two bins of taxa on
290 each side of the edge, and consequently phylofactorization brings in two com-
291plementary perspectives for analyzing the data: factor-based analysis and bin-
292based analysis. Factor-based analysis looks at the each factor as an inference
293on an edge in the phylogeny, conditioned on the previous inferences already
294made, and indicating that taxa on one side of an edge respond differently to
295the independent variable compared to taxa on the other side of the edge. Bin-
296based analysis, on the other hand, looks at the set of clades resulting from a
297certain number of factors - what we call a “binned phylogenetic unit” (BPU).
298These bins will create a lower-dimensional, compositional dataset and can be
299freed from the underlying ILR coordinates for different analyses on these amal-
300gamated clades. While factor-based analysis provides inferences about the splits
301in the phylogeny, BPU-based analysis conditions on the factors and bins OTUs
302based on which factors they contain. BPU-based analysis can inform sequence
303binning in future research aimed at controlling for previously-identified phyloge-
304netic causes of variation, and combine the effects of multiple up-stream factors
305for predictions of OTU abundance. See the supplementary text for a more
306detailed discussion of factor-based and bin-based analyses.

307 Phylofactorization will benefit from community discussion and further research
308overcoming general statistical challenges common to greedy algorithms and anal-
309ysis of phylogenetically-structured compositional data. For instance, the log-
310ratio transform at the heart of phylofactorization requires researchers deal with
311zeros in compositional datasets. While there are many methods for dealing with
312zeros [1, 23, 25], it’s unclear which method is most robust for downstream phylo-
313factorization of sparse OTU tables. Second, phylofactorization as presented here
314does not allow for multiple regression of ILR basis elements - the set of factors
315identified after n iterations may explain less variation combined than an al-
316ternative set of factors that did not maximize the explained variance at each
317iteration. This limitation may be overcome by running many replicates of a
318stochastic greedy algorithm and choosing that which maximizes the explained
319variance after n factors. Third, the researcher must choose an objective function
320which matches her question, and future research can map out which objective

321 functions are appropriate for which questions in microbial ecology. Fourth, like
322 any method performing inference based on phylogenetic structure, phylofactor-
323 ization assumes an accurate phylogeny. Accurate statistical statements about a
324 researcher's confidence in phylofactors must incorporate the uncertainty in our
325 constructed phylogeny. Finally, future research can investigate the unique kinds
326 of errors in phylofactorization: in addition to the multiple-hypothesis testing of
327 edges, phylofactorization may propagate errors in the greedy algorithm, and,
328 even when taxa are correctly factored into the appropriate functional bins, the
329 presence of multiple factors in the same region of the tree can lead to uncer-
330 tainty about the exact edge along which a putative trait arose (see supplement
331 for more discussion on the uncertainty of which edge to annotate).

332 Incorporating that phylogenetic structure into the analysis of microbiome datasets
333 has been a major challenge [24], and now phylofactorization provides a general
334 framework for rigorous exploration of phylogenetically-structured compositional
335 datasets. The soil dataset analyzed above, for instance, contains 3,379 OTUs
336 and 580 samples, and phylofactorization of the clades affected by pH in the
337 soil dataset yielded not just the three dominant factors used for ordination-
338 visualization, but 2,091 factors in all, each with an intricate phylogenetic story.
339 Many Acidobacteria are acidophiles, but some - Chloracidobacteria, Acidobacteria-
340 6, S035, and some undescribed classes of bacteria factored here - appear to be al-
341 kaliphiles. By incorporating the phylogenetic structure of microbiome datasets,
342 the big data of the modern sequence-count boom just got bigger, and future
343 research will need to consider how to organize, analyze and visualize the large
344 amounts of phylogenetic detail that can now be obtained from the analysis of
345 microbiome datasets.

346 Conclusions

347 Phylofactorization is a robust tool for analyzing marker gene sequence-count
348 datasets for phylogenetic patterns underlying microbial community responses
349 to independent variables. Phylofactorization accounts for the compositional na-
350 ture of the data and the underlying phylogeny and produces inferences that are
351 independent and more powerful than application of the ILR transform to the
352 rooted phylogeny. The R package 'phylofactor' has built-in parallelization that
353 can be used to analyze large microbiome datasets, and allows generalized linear

354 modeling to identify clades which respond to treatments or multiple environ-
355 mental gradients.

356 Phylofactorization can connect the pipeline of microbiome studies to focused
357 studies of microbial physiology. As researchers identify lineages with putative
358 functional ecological responses, taxa within those lineages - even if they are not
359 the same OTUs - can be cultivated and their genomes screened to uncover the
360 physiological mechanisms underlying the lineages' shared response.

361 Phylofactorization improves the pipeline for analyzing microbiome datasets by
362 allowing researchers to objectively determine the appropriate phylogenetic scales
363 for analyzing microbiome datasets - a family here, an unclassified split there -
364 instead of performing multiple comparisons at each taxonomic level. Instead of
365 principle components analysis or principle coordinates analysis, phylofactoriza-
366 tion can be used as for exploratory data analysis and dimensionality reduction
367 tool in which the "components" are identifiable clades in the tree of life, a far
368 more intuitive and informative component for biological variation than multi-
369 species loadings.

370 Phylofactorization can allow researchers to annotate online databases of the
371 microbial tree of life, permitting predictions about the physiology of unclassi-
372 fied and uncharacterized life forms based on previous phylogenetic inferences in
373 sequence-count data. By allowing researchers to make inferences on the same
374 tree and potentially annotate an online tree of life, phylofactorization may bring
375 on a new era of characterizing high-throughput phylogenetic annotations, filling
376 in the gaps the microbial tree of life.

377 An R package for phylofactorization with user-friendly parallelization is now
378 available online at <https://github.com/reptalex/phylofactor>.

379 Methods

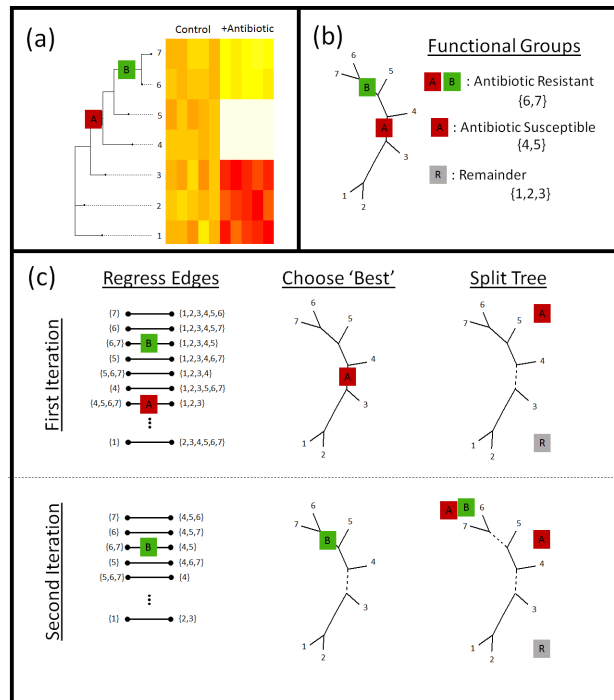


Figure 4: Phylofactorization: (a) Phylofactorization changes variables from tips of the phylogeny (OTUs used in analysis of microbiome datasets) to edges of the phylogeny with the largest predictable differences between taxa on each side of the edge. To illustrate this method, we consider the treatment of a bacterial community with an oxazolidinone. Oxazolidinones target gram-positive bacteria and will likely lead to a decrease in the relative abundances of gram-positive bacteria (antibiotic susceptible clade, A, having the antibiotic target). Among the antibiotic susceptible bacteria, phylofactor can identify monophyletic clades that are resistant relative to other antibiotic-susceptible bacteria due to a vertically-transmitted trait (B) such as the loss of the antibiotic target or enzymes that break down the antibiotic. (b) The two phylogenetic factors produce three meaningful bins of taxa - those susceptible to antibiotics (A), those within the susceptible clade that are resistant to antibiotics (A+B), and a potentially paraphyletic remainder. (c) Phylofactorization is a greedy algorithm to extract the edges which capture the most predictable differences in the response of relative abundances among taxa on the two sides of each edge. (c, top row) For the first iteration, all edges are considered - an ILR coordinate is created for each edge using equation (1) and the ILR coordinate is regressed against the independent variable. The edge which maximizes the objective function is chosen. Depicted above, the first factor corresponds to the edge separating antibiotic susceptible bacteria from the rest. Then, the tree is split - all subsequent comparisons along edges will be contained within the sub-trees. The conceptual justification for limiting comparisons within sub-trees is to prevent over-lapping comparisons: once we identify the antibiotic susceptible clade, we want to look at which taxa within that clade behave differently from other taxa within that clade. (c, bottom row) For the second iteration, the remaining edges are considered, ILR coordinates within sub-trees are constructed. The edge maximizing the objective function is selected and the tree is split at that edge. For more details, see the section “PhyloFactor” in the supplemental info.

380 **Phylogenetically-Structured compositional data**

381 Microbiome datasets are “phylogenetically-structured compositional data”, com-
382 positions of parts linked together by a phylogeny for which only inferences on
383 relative abundances can be drawn. The phylogeny is the scaffolding for the
384 evolution of vertically-transmitted traits, and vertically-transmitted traits may
385 underlie an organism’s functional ecology and response to perturbations or envi-
386 ronmental gradients. Performing inference on the edges in a phylogeny driving
387 variation in the data can be useful for identifying clades with putative traits
388 causing related taxa to respond similarly to treatments, but such inferences
389 must account for the compositional nature of the sequence-count data.

390 A standard analysis of microbiome datasets uses only the distal edges of the
391 tree - the OTUs - and a few edges within the tree separating Linear taxonomic
392 groups. However, a phylogeny of D taxa and no polytomies is composed of
393 $2D - 3$ edges, each connecting two disjoint sets of taxa in the tree with no
394 guarantee that splits in Linear taxonomy corresponds to phylogenetic splits
395 driving variation in our dataset. Thus, instead of analyzing just the tips and
396 a series of Linear splits in the tree, a more robust analysis of phylogenetically-
397 structured compositional data should analyze all of the edges in the tree. To
398 do that, we draw on the isometric log-ratio transform from compositional data
399 analysis, which has been used to search for a taxonomic signature of obesity
400 in the human gut flora [14] and incorporated into packages for downstream
401 principal components analysis [21]. However, to the best of our knowledge,
402 the previous literature using the isometric log-ratio transform in microbiome
403 datasets has used random or standard sequential binary partitions, and not
404 explicitly incorporated the phylogeny as their sequential binary partition.

405 **The Isometric Log-Ratio Transform of a rooted phylogeny**

406 The isometric log-ratio (ILR) transform was developed as a way to transform
407 compositional data from the simplex into real space where standard statistical
408 tools can be applied [11, 10]. A sequential binary partition is used to construct a
409 new set of coordinates, and the phylogeny is a natural choice for the sequential
410 binary partition in microbiome datasets. Instead of analyzing relative abun-
411 dances, y_i , of D different OTUs, the ILR transform produces $D - 1$ coordinates,
412 x_i^* (called “balances”). Each balance corresponds to a single internal node of the
413 tree and represents the averaged difference in relative abundance between the

414 taxa in the two sister clades descending from that node (the difference being
415 appropriately measured as a log-ratio due to the compositional nature of the
416 data; see SI for more detailed description of the ILR transform). For an arbitrary
417 node indicating the split of a group, R with r elements from the group, S
418 with s elements, the ILR balance can be written as

$$x_{\{R,S\}}^* = \sqrt{\frac{rs}{r+s}} \log \left(\frac{g(\mathbf{y}_R)}{g(\mathbf{y}_S)} \right) \quad (1)$$

419 where $g(\mathbf{y}_R)$ is the geometric mean of all y_i for $i \in R$.

420 We refer to the ILR transform corresponding to a rooted phylogeny as the
421 “rooted ILR”. The rooted ILR creates a set of ILR coordinates, $\{x_i^*\}$, where
422 each coordinate corresponds to the “balance” between sister clades at each split
423 in the phylogenetic tree. The balances in a rooted ILR transform in equation (1)
424 can be intuited as the average difference between taxa in two groups, and splits
425 in the tree which meaningfully differentiate taxa will be those splits in which
426 the average difference between taxa in two groups changes predictably with an
427 independent variable. Inferences on ILR coordinates, then, map to inferences
428 on lineages in the phylogenetic tree.

429 The rooted ILR coordinates provide a natural way to analyze microbiota data as
430 they measure the difference in the relative abundances of sister clades and may
431 be useful in identifying effects contained within clades such as zero-sum competition
432 of close relatives or the substitution of one relative for another across
433 environments. However, if we desire to link the effect of an external covariate
434 (e.g. antibiotics vs. no antibiotic treatment) to clades within the phylogeny,
435 the best comparison may not be between sister clades, but instead between all
436 other clades, controlling for any other phylogenetic splits or factors we may
437 know of (e.g. we may compare a lineage within gram-positives with all other
438 gram-positives, once we’ve identified the gram-positive vs. gram-negative split
439 as an important factor for antibiotic susceptibility). We refer to this unrooted
440 approach as ‘phylofactorization’.

441 For the task of linking an external covariate to individual clades in the phylogeny,
442 we examine three features of the rooted ILR that can be improved on
443 by phylofactorization by considering a treatment that decreases the abundance
444 of one and only one clade, B , whose closest relative is clade A . Regression on
445 the rooted ILR coordinates may identify the balance $x_{\{A,B\}}^*$ corresponding to
446 the most recent common ancestor of clades A and B as having that strongest

447 response to the treatment, but regression on this coordinate will suggest that
448 clade B decreases relative to A , leading to structured residuals in the original
449 dataset due to an inability to account for the increase in clade B relative to the
450 rest of the OTUs in the data (Fig. 4a). Second, all partitions between affected
451 clade and the root will be affected. If each balance is tested independently,
452 the rooted ILR may identify many clades that are affected by antibiotics; the
453 correlations between coordinates can yield a high false-positive rate if just one
454 clade is affected (Fig. 4b). Finally, the ILR transformation does not work with
455 polytomies common in real, unresolved phylogenies. Any polytomy will produce
456 a split in the phylogeny between three or more taxa, and there is no general
457 way to describe the balance of relative abundances of three or more parts using
458 only one coordinate.

459 Nonetheless, the simplicity and theoretical foundations underlying the ILR, and
460 the instant appeal of applying it to the sequential-binary partition of the phy-
461 logeny, motivate the rooted ILR as a simple tool for analysis of the phylogenetic
462 structure in compositional data. For that reason, we use the rooted ILR as a
463 baseline for comparison of our more complicated method of phylogenetic factor-
464 ization.

465 **Phylofactorization**

466 The shortcomings of the rooted ILR can be remedied by modifying the ILR
467 transform to apply not to the nodes or splits in a phylogeny, but to the edges in
468 an unrooted phylogeny. While ILR coordinates of nodes allow a comparison of
469 sister clades, ILR coordinates along edges allow comparison of taxa with putative
470 traits that arose along the edge against all taxa without those putative traits.
471 Traits arise along edges of the phylogeny and so, for annotation of online trees
472 of life, effects in a clade are best mapped to a chain of edges in the phylogeny.

473 However, the ILR transform requires a sequential binary partition, and the edges
474 don't immediately provide a clear candidate for a sequential binary partition. In
475 what we refer to as "phylofactorization", one can iteratively construct a sequen-
476 tial binary partition from the unrooted phylogeny by using a greedy algorithm
477 by sequentially choosing edges which maximize a researcher's objective function.
478 Phylofactorization consists of 3 steps (Box 1): (1) Consider the set of possible
479 primary ILR basis elements corresponding to a partition along any edge in the
480 tree (including the tips). (2) Choose the edge whose corresponding ILR basis

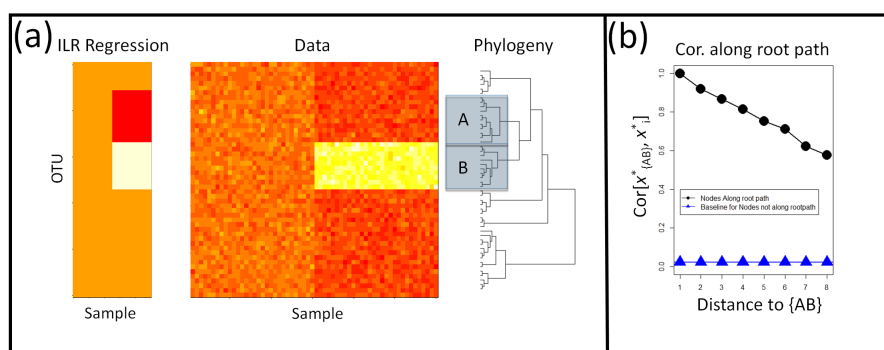


Figure 5: Shortcomings of Rooted ILR (a) The isometric log-ratio transform corresponding to a phylogeny rooted at the common ancestor is inaccurate for geometric changes within clades. Here, absolute abundances of 50 taxa in 30 samples per site were simulated across two sites. An affected clade, B , is up-represented in the second site. Regression on the rooted ILR coordinates, x_i^* , against the sample site indicated that the partition separating clade A, B , referred to as $x_{\{A,B\}}^*$, had the highest test-statistic, but the rooted ILR predicts fold-changes in B relative to A , not fold changes in B relative to the rest of the taxa. (b) Consequently, when one clade increase in abundance while the rest remain unaffected, partitions between the affected clade and the root will also have a signal leading to a correlation in the coordinates along the path from B to the root. The correlation plotted here is the absolute value of the correlation coefficient, and the baseline correlation was estimated as the average absolute value of the correlation coefficient between ILR coordinates not along the root-path of the affected clade.

481 element maximizes some objective function - such as the test-statistic from re-
482 gression or the percent of variation explained in the original dataset - and the
483 groups of taxa split by that edge form the first partition. (3) Repeat steps 1 and
484 2, constructing subsequent ILR basis elements corresponding to remaining edges
485 in the phylogeny and made orthogonal to all previous partitions by limiting the
486 comparisons to taxa within the groups of taxa un-split by previous partitions.

487 Explicitly, the first iteration of phylofactorization considers a set of candi-
488 date ILR coordinates, $\{x_e^*\}$ corresponding to the two groups of taxa split by
489 each edge, e . Then, regression is performed on each of the ILR coordinates,
490 $x_e^* \sim f(X)$ for an appropriate function, f and a set of independent variables,
491 X . The edge, e_1^+ , which maximizes the objective function is chosen as the first
492 phylogenetic factor. In this paper, our objective function is the difference be-
493 tween the null deviance of the ILR coordinate and the deviance of the generalized
494 linear model explaining that ILR coordinate as a function of the independent
495 variables. We use this objective function as a measure of the amount of variance
496 explained by regression on each edge because the total variance in a composi-
497 tional dataset is constant and equal to the sum of the variances of all ILR
498 coordinates corresponding to any sequential binary partition. Consequently, at
499 each iteration there is a fixed amount of the total variance remaining in the
500 dataset, and so at the candidate ILR coordinate which captures the greatest
501 fraction of the total variance in the dataset is the one with the greatest amount
502 of variance explained by the regression. After identifying e_1^+ , we cut the tree in
503 two sub-trees along the edge, e_1^+ .

504 For the second iteration, another set of candidate ILR coordinates is constructed
505 such that their underlying balancing elements are orthogonal to the first ILR
506 coordinate. Orthogonality is ensured by constructing ILR coordinates contrast-
507 ing the abundances of taxa along each edge, restricting the contrast to all taxa
508 within the sub-tree in which the edge is found. A new edge, e_2^+ , which maximizes
509 the objective function is chosen as the second factor, the sub-tree containing this
510 edge is cut along this edge to produce two sub-trees, and the process is repeated
511 until a desired number of factors is reached or until a stopping criterion is met.
512 More details on the algorithm, along with a discussion on objective functions,
513 is contained in the SI.

514 While one could use other methods of amalgamating abundances along edges,
515 the conceptual importance of using the ILR transform is twofold: the ILR trans-
516 form has proven asymptotic normality properties for compositional data to allow

517 the application of standard multivariate methods [11], and the ILR transform
518 serves as a measure of contrast between two groups. The log-ratio used in phylo-
519 factor is an averaged ratio of abundances of taxa on two sides of an edge (see
520 supplement for more detail), thus phylofactorization searches the tree for the
521 edge which has the most predictable difference between taxa on each side of the
522 edge, or, put differently, the edge which best differentiates taxa on each side.
523 Thus, each edge that differentiates taxa and their responses to independent
524 variables is considered a phylogenetic “factor” driving variation in the data.

525 The output of phylofactorization is a set of orthogonal, sequentially “less im-
526 portant” ILR basis elements, their predicted balances, and all other information
527 obtained from regression. After the first iteration of phylofactorization, we are
528 left with an ILR basis element corresponding to the edge which maximized our
529 objective function and split the dataset into two disjoint sub-trees, or sets of
530 OTUs that we henceforth refer to as “bins”, and we have an estimated ILR
531 balancing element, $\hat{x}_1^*(X)$, where X is our set of independent variables. Sub-
532 sequent factors will split the bins from previous steps, and after n iterations
533 one has n factors that can be mapped to the phylogeny, $n + 1$ bins for bin-
534 ning taxa based on their phylogenetic factors, n estimates of ILR balancing
535 elements, and an orthonormal ILR basis that can be used to project the data
536 onto a lower dimensional space. The sequential splitting of bins in phylofactor-
537 ization ensures sequentially independent inferences - having already identified
538 group B as hyper-abundant relative to group A in the example illustrated in
539 Fig. 4, downstream factors must analyze sub-compositions entirely within B
540 and within A .

541 Computational Tools

542 Phylofactorization was done using the R package “phylofactor” available at
543 <https://github.com/reptalex/phylofactor>. The R package contains detailed help
544 files that demo the use of the package, and the exact code used in analyses and
545 visualization in this paper are available in the supplementary materials. The
546 rooted ILR transform was performed as described in [10] where the sequential
547 binary partition was the rooted phylogeny.

548 Power Analysis of Rooted ILR and Phylofactorization

549 To compare the ability of phylofactorization and the rooted ILR to identify
550 clades of OTUs with shared associations with independent variables, we simu-
551 lated random communities of $D = 50$ OTUs and $p = 40$ samples by simulating
552 random absolute abundances, $N_{i,j}$, such that $\log N_{i,j}$ were i.i.d Gaussian ran-
553 dom variables with mean $\mu = 8$ and standard deviation $\sigma = 0.5$. The OTUs
554 were connected by a random tree (the tree remained constant across all simula-
555 tions), and then either 1 or 3 clades were randomly chosen to have associations
556 with a binary “environment” independent variable with $p = 20$ samples for each
557 of its two values to represent an equal sampling of microbial communities across
558 two environments.

559 For simulations with one significant clade, the abundances of all the OTUs
560 within that clade increased by a factor a in the second environment where $a \in$
561 $\{1.5, 3, 6\}$. For simulations with three significant clades, the three clades were
562 drawn at random and randomly assigned a fold-change from the set $\{\pi^b, 0.5^b, \exp(-b)\}$
563 in a randomly chosen environment where $b \in \{1, 2, 5\}$. For each fold-change,
564 500 replicates were run to compare the power of the rooted ILR and phylofac-
565 torization in correctly identifying the affected clades.

566 Regression of rooted ILR coordinates was performed and the coordinates were
567 ranked by the difference between their null deviance and the model deviance.
568 The ability of a rooted ILR coordinate to identify the correct 1 clade or 3 clades
569 was measured by the percent of its top 1 or 3 ILR coordinates, respectively,
570 which corresponded to the node on the tree from which the affected clade(s)
571 originated. The ability of phylofactor to identify the correct 1 clade or 3 clades
572 was measured by the percent of the factors that correctly split an affected clade
573 from the rest (e.g. the percent of factors corresponding to edges along which a
574 trait arose).

575 For the 3 clade simulations, we also compared the amount of variance explained
576 by 3 factors in phylofactorization with the amount of variance explained by the
577 top 3 ILR coordinates in the rooted ILR. The amount of variance explained
578 was measured as the difference in the null deviance and the model deviance,
579 summed across all three factors or the top 3 ILR coordinates.

580 **KS-based Stopping Function for PhyloFactor**

581 While a researcher can iterate through phylofactorization until a full basis of
582 $D - 1$ ILR coordinates is constructed, there is value in stopping the iteration
583 when all of the clades have been identified or at a conservative underestimate of
584 the true number of phylogenetic factors. We implemented a stopping function
585 based on a Kolmogorov-Smirnov (KS) test of the distribution of P-values from
586 analyses of variance of the regressions on candidate ILR coordinates. If there
587 is no phylogenetic signal, we anticipate the true distribution of P-values to be
588 uniform (albeit with some dependence among the P-values due to overlap in the
589 OTUs used in the ILR coordinates). Thus, we tested the ability of phylofactor
590 to correctly identify the number of clades if phylofactorization is stopped when
591 a KS test of the P-values produces its own P-value $P_{KS} > 0.05$.

592 We simulated 300 replicate communities with M clades for each $M \in \{1, \dots, 10\}$.
593 For simulations with M clades, $D = 50$ and $p = 40$ communities were simulated
594 as above and fold changes, c , were drawn as log-normal random variables where
595 $\log(c_k)$ were i.i.d Gaussian random variables with $\mu = 0$ and $\sigma = 3$ for $k =$
596 $1, \dots, M$. The number of clades identified by phylofactor for a given true number
597 of clades, $K_{M,r}$, was tallied for $r = 1, \dots, 300$. We calculate the mean \bar{K}_M across
598 all replicates and, for visualization purposes, interpolate the $\alpha = 0.025$ and
599 $\alpha = 0.975$ quantiles by finding the best fit of a logistic function to the cumulative
600 distribution of $\{K_{M,r}\}_{r=1}^{300}$ for each M .

601 **Analysis of Fecal/Oral microbiome data**

602 16S amplicon sequencing data from Caporaso et al. (2011) [4] were downloaded
603 from the MG-RAST database (<http://metagenomics.anl.gov/>) along with as-
604 sociated metadata. QIIME [5] was used to trim primers from these data, and
605 to cluster OTUs with the Greengenes reference database (May 2013 version;
606 <http://greengenes.lbl.gov>). Longer sequence lengths in the greengenes database
607 (~ 1400 BP) compared to the original Illumina sequences (~ 123 BP) allows more
608 informative base pairs for phylogenetic tree construction. We used the phylo-
609 genetic tree that is included with the greengenes database for all analyses. The
610 resulting OTU table was rarefied to 6000 sequences per sample.

611 10 time points were randomly drawn from each of the male tongue, female
612 tongue, male feces and female feces datasets, giving a total of $n=20$ samples at
613 each site. Taxa present in fewer than 30 of the 40 samples were discarded, and

614 phylofactorization was done by adding pseudo-counts of 0.65 to all 0 entries in
615 the dataset [1], converting counts in each sample to relative abundances, and
616 then regressing the ILR coordinates against body site. The complete R script
617 is available in the file “Data Analysis pipeline of the FT microbiome”.

618 Complete phylofactorization of this dataset was performed by stopping the al-
619 gorithm when a KS-test on the uniformity of P-values from analyses of variance
620 of regression on candidate ILR-coordinates yielded $P_{KS} > 0.05$. These results
621 were compared with a standard, multiple hypothesis-testing analysis of CLR-
622 transformed data. The summary of the taxonomic detail at the first three factors
623 is provided in the results section, and a full list of the taxa factored at each step
624 is available in the supplement and can be further explored using the R pipeline
625 provided.

626 **Analysis of Soil microbiome data**

627 The soil microbiome dataset from [28] was included to illustrate the ability
628 of phylofactor to work on bigger microbiome datasets with continuous indepen-
629 dent variables and multiple regression. Details on sample collection, sequencing,
630 meta-data measurements and OTU clustering are available in [28]. The phy-
631 logeny was constructed by aligning representative sequences using SINA [27],
632 trimming bases that represented gaps in $\geq 20\%$ of sequences, and using fasttree
633 [26].

634 The complete dataset contained 123,851 OTUs and 580 samples. Data were
635 filtered to include all OTUs with on average 2 or more sequences counted across
636 all samples, shrinking the dataset to $D=3,379$ OTUs. The data were further
637 trimmed to include only those samples with available pH, C and N meta-data,
638 reducing the sample size to $n=551$.

639 Phylofactorization was done by adding pseudo-counts of 0.65 to all 0 entries in
640 the dataset [1], converting counts in each sample to relative abundances, and
641 performing multiple regression of pH, C and N on ILR coordinates. The first
642 three factors are used for ordination-visualization. To determine the relative
643 importance of each abiotic variable in driving phylogenetic patterns of microbial
644 community composition, we used the lmg method from the R package ‘relaimpo’
645 [17] which averages the sequential sums of squares over all orderings of regressors
646 to obtain a measure of relative importance of each regressor in the multivariate
647 model.

648 Acknowledgments

649 ADW would like to acknowledge L. Ma for his feedback and help incorporating
650 this method into the statistical literature. JS was supported in part by the
651 Duke University Medical Scientist Training Program. This paper is published
652 by support from and in loving memory of D. Nemergut.

653 Declarations

654 **Competing Interests:** The authors have no competing interests in relation
655 to this work.

656 **Availability of Data and Materials:** The data were obtained from previ-
657 ous studies and are available online through the original studies. The R pack-
658 age 'phylofactor' is available at <https://github.com/reptalex/phylofactor> and all
659 other R files used in the analysis and visualization are available online.

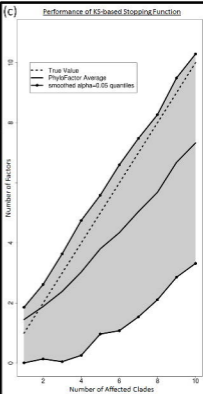
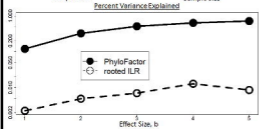
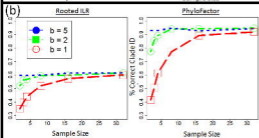
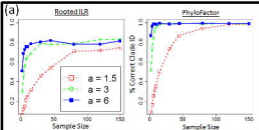
660 References

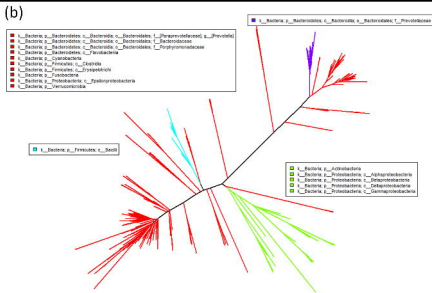
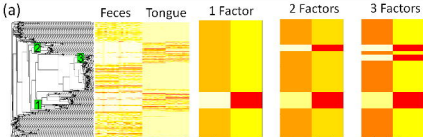
- 661 [1] John Aitchison. The statistical analysis of compositional data. 1986.
- 662 [2] Richard D Bardgett, Chris Freeman, and Nicholas J Ostle. Microbial con-
663 tributions to climate change through carbon cycle feedbacks. *The ISME*
664 *Journal*, 2(8):805–814, 2008.
- 665 [3] Roeland L Berendsen, Corne MJ Pieterse, and Peter AHM Bakker. The rhi-
666 zosphere microbiome and plant health. *Trends in plant science*, 17(8):478–
667 486, 2012.
- 668 [4] J Gregory Caporaso, Christian L Lauber, Elizabeth K Costello, Donna
669 Berg-Lyons, Antonio Gonzalez, Jesse Stombaugh, Dan Knights, Pawel
670 Gajer, Jacques Ravel, Noah Fierer, et al. Moving pictures of the human
671 microbiome. *Genome Biol*, 12(5):R50, 2011.
- 672 [5] J Gregory Caporaso, Christian L Lauber, William A Walters, Donna Berg-
673 Lyons, James Huntley, Noah Fierer, Sarah M Owens, Jason Betley, Louise
674 Fraser, Markus Bauer, et al. Ultra-high-throughput microbial community

- 675 analysis on the illumina hiseq and miseq platforms. *The ISME journal*,
676 6(8):1621–1624, 2012.
- 677 [6] Human Microbiome Project Consortium et al. Structure, function and
678 diversity of the healthy human microbiome. *Nature*, 486(7402):207–214,
679 2012.
- 680 [7] Joel Cracraft. Species concepts and speciation analysis. In *Current or-*
681 *nithology*, pages 159–187. Springer, 1983.
- 682 [8] Joel Cracraft. Species concepts in theoretical and applied biology: a system-
683 atic debate with consequences. *Species concepts and phylogenetic theory:*
684 *A debate*, pages 30–43, 2000.
- 685 [9] Tao Ding and Patrick D Schloss. Dynamics and associations of microbial
686 community types across the human body. *Nature*, 509(7500):357, 2014.
- 687 [10] Juan José Egozcue and Vera Pawlowsky-Glahn. Groups of parts and their
688 balances in compositional data analysis. *Mathematical Geology*, 37(7):795–
689 828, 2005.
- 690 [11] Juan José Egozcue, Vera Pawlowsky-Glahn, Glòria Mateu-Figueras, and
691 Carles Barcelo-Vidal. Isometric logratio transformations for compositional
692 data analysis. *Mathematical Geology*, 35(3):279–300, 2003.
- 693 [12] Paul G Falkowski, Tom Fenchel, and Edward F Delong. The microbial
694 engines that drive earth’s biogeochemical cycles. *science*, 320(5879):1034–
695 1039, 2008.
- 696 [13] Noah Fierer and Robert B Jackson. The diversity and biogeography of soil
697 bacterial communities. *Proceedings of the National Academy of Sciences of*
698 *the United States of America*, 103(3):626–631, 2006.
- 699 [14] Mariel M Finucane, Thomas J Sharpton, Timothy J Laurent, and Kather-
700 ine S Pollard. A taxonomic signature of obesity in the microbiome? getting
701 to the guts of the matter. *PloS one*, 9(1):e84689, 2014.
- 702 [15] Jonathan Friedman and Eric J Alm. Inferring correlation networks from
703 genomic survey data. *PLoS Comput Biol*, 8(9):e1002687, 2012.
- 704 [16] Keith Gregg. Engineering gut flora of ruminant livestock to reduce forage
705 toxicity: progress and problems. *Trends in biotechnology*, 13(10):418–421,
706 1995.

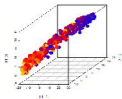
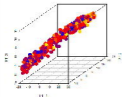
- 707 [17] Ulrike Grömping et al. Relative importance for linear regression in r: the
708 package relaimpo. *Journal of statistical software*, 17(1):1–27, 2006.
- 709 [18] B Guggenheim. Streptococci of dental plaques. *Caries research*, 2(2):147–
710 163, 1968.
- 711 [19] François Keck, Frédéric Rimet, Agnes Bouchez, and Alain Franc. phylsig-
712 nal: an r package to measure, test, and explore the phylogenetic signal.
713 *Ecology and evolution*, 6(9):2774–2780, 2016.
- 714 [20] Omry Koren, Aymé Spor, Jenny Felin, Frida Fåk, Jesse Stombaugh,
715 Valentina Tremaroli, Carl Johan Behre, Rob Knight, Björn Fagerberg,
716 Ruth E Ley, et al. Human oral, gut, and plaque microbiota in patients
717 with atherosclerosis. *Proceedings of the National Academy of Sciences*,
718 108(Supplement 1):4592–4598, 2011.
- 719 [21] Kim-Anh Le Cao, Mary-Ellen Costello, Vanessa Anne Lakis, Francois Bar-
720 tolo, Xin-Yi Chua, Remi Brazeilles, and Pascale Rondeau. mixmc: a mul-
721 tivariate statistical framework to gain insight into microbial communities.
722 *bioRxiv*, page 044206, 2016.
- 723 [22] Jørgen J Leisner, Birgit Groth Laursen, Hervé Prévost, Djamel Drider, and
724 Paw Dalgaard. Carnobacterium: positive and negative effects in the envi-
725 ronment and in foods. *FEMS microbiology reviews*, 31(5):592–613, 2007.
- 726 [23] Josep A Martín-Fernández, Carles Barceló-Vidal, and Vera Pawlowsky-
727 Glahn. Dealing with zeros and missing values in compositional data sets
728 using nonparametric imputation. *Mathematical Geology*, 35(3):253–278,
729 2003.
- 730 [24] Jennifer BH Martiny, Stuart E Jones, Jay T Lennon, and Adam C Mar-
731 tiny. Microbiomes in light of traits: A phylogenetic perspective. *Science*,
732 350(6261):aac9323, 2015.
- 733 [25] Vera Pawlowsky-Glahn and Antonella Buccianti. *Compositional data anal-
734 ysis: Theory and applications*. John Wiley & Sons, 2011.
- 735 [26] Morgan N Price, Paramvir S Dehal, and Adam P Arkin. Fasttree 2–
736 approximately maximum-likelihood trees for large alignments. *PloS one*,
737 5(3):e9490, 2010.

- 738 [27] Elmar Pruesse, Jörg Peplies, and Frank Oliver Glöckner. Sina: accu-
739 rate high-throughput multiple sequence alignment of ribosomal rna genes.
740 *Bioinformatics*, 28(14):1823–1829, 2012.
- 741 [28] Kelly S Ramirez, Jonathan W Leff, Albert Barberán, Scott Thomas Bates,
742 Jason Betley, Thomas W Crowther, Eugene F Kelly, Emily E Oldfield,
743 E Ashley Shaw, Christopher Steenbock, et al. Biogeographic patterns in
744 below-ground diversity in new york city’s central park are similar to those
745 observed globally. In *Proc. R. Soc. B*, volume 281, page 20141988. The
746 Royal Society, 2014.
- 747 [29] Liam J Revell. phytools: an r package for phylogenetic comparative biology
748 (and other things). *Methods in Ecology and Evolution*, 3(2):217–223, 2012.
- 749 [30] Mark D Robinson, Davis J McCarthy, and Gordon K Smyth. edger: a
750 bioconductor package for differential expression analysis of digital gene ex-
751 pression data. *Bioinformatics*, 26(1):139–140, 2010.
- 752 [31] Mikhail Tikhonov, Robert W Leach, and Ned S Wingreen. Interpreting 16s
753 metagenomic data without clustering to achieve sub-otu resolution. *The*
754 *ISME journal*, 9(1):68–80, 2015.

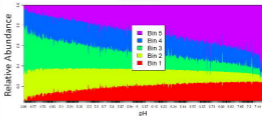
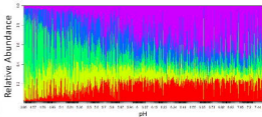


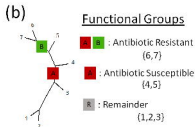
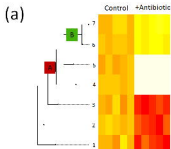


(a) ILR Ordination



(b) BPU Projection





(c) **Regress Edges**

Choose 'Best'

Split Tree

First Iteration



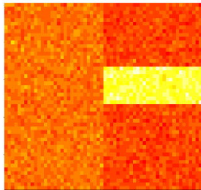
Second Iteration



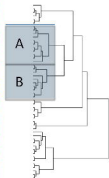
ILR Regression



Data



Phylogeny



Cor. along root path

