

1 Proteobacteria drive significant functional variability
2 in the human gut microbiome

3 Patrick H. Bradley¹, Katherine S. Pollard^{1,2*}

4 August 9, 2016

5 1. Gladstone Institutes, San Francisco, CA.

6 2. Division of Biostatistics, Institute for Human Genetics, and Institute for Computational
7 Health Sciences, University of California, San Francisco, CA.

8 * Corresponding author.

9 E-mails: patrick.bradley@gladstone.ucsf.edu, katherine.pollard@gladstone.ucsf.edu

10

Abstract

11 While human gut microbiomes vary significantly in taxonomic composition, biological
12 pathway abundance is surprisingly invariable across hosts. We hypothesized that healthy
13 microbiomes appear functionally redundant due to factors that obscure differences in gene
14 abundance across hosts. To account for these biases, we developed a powerful test of gene
15 variability, applicable to shotgun metagenomes from any environment. Our analysis of
16 healthy stool metagenomes reveals thousands of genes whose abundance differs signifi-
17 cantly between people consistently across studies, including glycolytic enzymes, lipopolysac-
18 charide biosynthetic genes, and secretion systems. Even housekeeping pathways contain a
19 mix of variable and invariable genes, though most deeply conserved genes are significantly
20 invariable. Variable genes tend to be associated with Proteobacteria, as opposed to taxa
21 used to define enterotypes or the dominant phyla Bacteroidetes and Firmicutes. These re-
22 sults establish limits on functional redundancy and predict specific genes and taxa that may
23 drive physiological differences between gut microbiomes.

24 Impact Statement

25 A statistical test for gene variability reveals extensive functional differences between healthy
26 human microbiomes.

27 Keywords

28 human gut microbiome, Proteobacteria, Bacteroidetes, Firmicutes, variance, shotgun metage-
29 nomics, statistical methods, functional redundancy, enterotypes

30 **1 Background**

31 The microbes that inhabit the human gut encode a wealth of proteins that contribute to a broad
32 range of biological functions, from modulating the human immune system [1, 2, 3] to par-
33 ticipating in metabolism [4, 5]. Shotgun metagenomics is revolutionizing our ability to iden-
34 tify protein-coding genes from these microbes and associate gene levels with disease [6], drug
35 efficacy [7] or side-effects [8], and other host traits. For instance, gut microbiota associated
36 with a traditional high-fiber agrarian diet encoded gene families involved in cellulose and xy-
37 lan hydrolysis, which were absent in age-matched controls eating a typical Western diet [9].
38 The functional capabilities of the gut microbiome go beyond statistical associations; a num-
39 ber of microbial genes have now been causally linked to host physiology. Examples include the
40 colitis-inducing cytolethal distending toxins of *Helicobacter hepaticus* [10] and the enzymes of
41 commensal bacteria that protect against these toxins by producing anti-inflammatory polysac-
42 charide A [11].

43 It is therefore surprising that healthy human gut microbiomes have been characterized as
44 functionally stable (i.e., invariable), with largely redundant gene repertoires in different hosts.
45 Several lines of evidence support this conclusion. First, biological pathway abundance tends to
46 be less variable across metagenomes than it is between isolate genomes [12], suggesting strong
47 selection for microbes that encode functions necessary for adaptation to the gut environment.
48 Second, the relative abundances of pathways are strikingly invariable compared to the relative
49 abundances of bacterial phyla in the same metagenomes [13, 12]. Thus, it appears that humans
50 harbor phylogenetically distinct gut communities that all do more or less the same things, ex-
51 cept in the context of disease or other extreme host phenotypes.

52 Functional redundancy deserves a closer look, however, because physiologically meaning-
53 ful differences in gene abundances between healthy human microbiomes could easily have
54 been missed. One primary factor may be that prior work did not look at quantitative abun-
55 dances of individual genes, but instead mainly summarized function at the level of Clusters
56 of Orthologous Groups (COG) categories and KEGG modules [13, 12, 14]. Summarizing genes
57 into groups will not have power to detect one component of a pathway or protein complex that
58 varies in abundance across hosts if other components are less variable. This masking of variable
59 genes is likely to occur because the presence and abundance of most COG categories and KEGG
60 modules will be dominated by core components (i.e., housekeeping genes) that are widely dis-
61 tributed across the tree of life and abundant in metagenomes. The only previous analyses of
62 individual genes asked whether they were universally detected across all individuals sampled
63 [12, 14]; however, universally-detected genes may still vary substantially in abundance, and
64 conversely, lower-abundance invariable genes may not be universally detected merely due to

65 sampling. This approach is also sensitive to read depth [12] and sample size [14]. Based on
66 these observations, we were motivated to quantitatively investigate functional redundancy at
67 the level of individual gene families.

68 To enable high-resolution, quantitative analysis of functional stability in the microbiome,
69 we developed a statistical test that identifies individual gene families whose abundances are
70 either significantly variable or invariable across samples. Our method incorporates solutions
71 to three major challenges to studying functional redundancy with shotgun metagenomics data.
72 The first key innovation of our approach is using a test statistic that captures residual variability
73 after accounting for overall gene abundance. This modeling choice is important because abun-
74 dant genes will be variable just by chance due to the correlation between mean and variance in
75 any sequencing experiment. Conversely, phylogenetically restricted genes will have relatively
76 low variance due to being less abundant. Furthermore, gene abundances can be sparse (i.e.,
77 zero in many samples). For all of these reasons simply ranking genes based on their variances
78 would yield many false positives and false negatives.

79 A second benefit of our modeling approach is that we can adjust for systematic differences
80 in a gene's measured level between studies to allow for quantitative integration of data from
81 multiple sources. Meta-analysis is essential for gaining sufficient power to detect variable genes
82 across the range of mean abundance levels. It also ensures robustness and generalizability of
83 discovered inter-individual differences, which occur by chance in small sets of metagenomes.

84 Finally, our method does not require predefined cases and controls, but instead enables
85 discovery of genes that drive functional differences between microbiomes without prior knowl-
86 edge of which groups of samples to compare. This is critical for the current phase of micro-
87 biome research, when many drivers of microbial community composition are unknown. Gene
88 families that contribute to survival in one particular type of healthy gut environment should
89 emerge as variable between hosts and their functions may point to drivers of community com-
90 position, mechanisms of microbe-host interactions, and biomarkers of presymptomatic disease
91 (e.g., pre-diabetes).

92 We applied our test to healthy gut metagenomes ($n = 123$) spanning three different shotgun
93 sequencing studies and found both significantly invariable (3,768) and variable (1,219) gene
94 families ($FDR < 5\%$). Many pathways, including some commonly viewed as housekeeping or
95 previously identified as invariable across gut microbiota (e.g., central carbon metabolism and
96 secretion), included significantly variable gene families. Phylogenetic distribution (PD) corre-
97 lated overall with variability in gene family abundance, and exceptions to this trend highlight
98 functions that may be involved in adaptation, such as two-component signaling and special-
99 ized secretion systems. Finally, we show that Proteobacteria, and not the major phyla Bac-
100 teroidetes and Firmicutes, are a major source for genes with the greatest variability in abun-

101 dance across hosts, suggesting a relationship between inflammation and gene-level differences
102 in gut microbial functions. This approach to discovering functions that distinguish microbial
103 communities is applicable to any body site or environment.

104 2 Results

105 2.1 A new test captures the variability of microbial gene families

106 We present a model that enables gene family abundance to be quantitatively compared across
107 metagenomes for thousands of microbial genes. In shotgun metagenomics data, different gene
108 families vary widely in average abundance (Figure 1). Gene family abundances can also vary
109 by study, both because of biological differences between populations, and for technical reasons
110 including library preparation, amplification protocol, and sequencing technology (see, e.g., Fig-
111 ure 1 G-H). To account for such effects, we fit a linear model of log abundance $D_{g,s}$ for gene g
112 in sample s as a function of the overall mean abundance μ_g and a term $\beta_{g,y}$ that quantifies the
113 offset for each study y :

$$D_{g,s} = \mu_g + \sum_{y \in Y} I_{y,s} \beta_{g,y} + \epsilon_{g,s} \quad (1)$$

114 where $I_{y,s}$ is an indicator variable that is 1 if sample s belongs to study y and 0 otherwise.

115 The residual $\epsilon_{g,s}$ quantifies how much the abundance of gene g in sample s differs from
116 the average abundance across samples in the same study as s . We denote the variance of the
117 residuals across samples by V_g^ϵ . When this statistic is small, the gene has similar abundance
118 across samples after accounting for study effects. A large value of V_g^ϵ indicates that samples
119 have very different abundances.

120 To assess the statistical significance of gene family variability, we compare the residual vari-
121 ance V_g^ϵ to a data-driven null distribution based on the negative binomial distribution (Figure
122 1—figure supplement 1, Methods). This approach is necessary because there is no straightfor-
123 ward formula for the p-value of V_g^ϵ . Our method looks for deviations from the null hypothe-
124 sis that gene families in the dataset have the same mean-variance relationship (i.e., the same
125 overdispersion). This choice of null is very important: if we were instead to simply test for
126 high variance, regardless of mean abundance, highly abundant gene families (e.g., single-copy
127 proteins in the bacterial ribosome) would be significantly variable despite being nearly uni-
128 versally present at equal abundance in each bacterial genome, because genes with high mean
129 abundance would have high variance in any sequencing experiment. Meanwhile, thousands of
130 lower-abundance gene families would appear to be significantly invariable simply by virtue of

131 having relatively low read counts.

132 We validated this approach using simulated data (see Methods, Figure 1—figure supplement 3)
133 and found that the residual variance test has high power and good control over the false posi-
134 tive rate when the overdispersion parameter k used in the null distribution was accurately esti-
135 mated. To make the test more robust to factors affecting the estimation of k (Figure 1—figure supplement 4),
136 we used simulation to control the false discovery rate empirically (Table 1). Our statistical test
137 can be applied to shotgun metagenomes to sensitively and specifically identify variable genes
138 in any environment without prior knowledge of factors that stratify relatively high versus low
139 abundance samples.

140 **2.2 Thousands of variable gene families in the gut microbiome**

141 To describe variation within healthy gut microbiota across different human populations, we
142 randomly selected 123 metagenomes of healthy individuals from the Human Microbiome Project
143 (HMP) [13], controls in a study of type II diabetes (T2D) [15], and controls in a study of glucose
144 control (GC) [16]. These span American, Chinese, and European populations, respectively (see
145 Methods). We mapped these metagenomes to KEGG Orthology families with ShotMAP [17]
146 and counted reads for 17,417 gene families. Accurately normalizing gene read counts so that
147 they were comparable across samples and studies is critical to our meta-analytical approach
148 and any quantitative evaluation of shotgun metagenomes. We therefore quantified gene family
149 abundance using log-transformed reads per kilobase of genome equivalents (log-RPKG) [18].

150 We found 2,357 gene families with more variability than expected and 5,432 with less (leav-
151 ing 9,628 non-significant) at an empirical FDR of 5% (Figure 2—figure supplement 1). Restrict-
152 ing the analysis to gene families with at least one annotated representative from a bacterial
153 or archaeal genome in KEGG, we obtained 1,219 significantly variable and 3,813 significantly
154 invariable gene families (and 2,194 non-significant). The differences in the residual variation
155 of these gene families can be visualized using a heatmap of the residual $\epsilon_{g,s}$ values (Figures
156 2—figure supplement 2, 2—figure supplement 3). The large number of genes that were less
157 variable than expected given their means supports the hypothesis of some functional redun-
158 dancy in the gut microbiome, potentially due to selection for core functions that make microbes
159 more successful in the gut environment. However, our discovery of thousands of significantly
160 variable genes across a range of abundance levels demonstrates that the gut microbiome is less
161 invariable than prior work suggested.

162 This result highlights the importance of a quantitative, gene-level evaluation of functional
163 stability. Importantly, the magnitude of the residual variance statistic V_g^c is not the sole deter-
164 minant of significance, as observed by the overlap in distributions of V_g^c between the variable,

165 invariable, and non-significant gene families. For example, both low-abundance gene families
166 with many zero values and high-abundance but invariable gene families will tend to have low
167 residual variance, but the evidence for invariability is much stronger for the second group. Our
168 test accurately discriminates between these scenarios, tending to call the second group signifi-
169 cantly invariable and not the first (Figure 2—figure supplement 1, inset).

170 **2.3 Biological pathways contain both invariable and variable components**

171 To test our hypothesis that the appearance of pathways and functional categories with similar
172 abundance across samples is driven by a subset of core components, we examined individual
173 gene variability within KEGG modules. As expected, we observed an overall signal of stability
174 at this broad level of gene groupings. Many of the pathways previously identified as invariable
175 (e.g., aminoacyl-tRNA metabolism, central carbon metabolism) indeed have more invariable
176 than variable genes. However, individual genes show a much more complex picture. Even the
177 most invariable pathways also include significantly variable genes (Figure 2). For example, the
178 highly conserved KEGG module set “aminoacyl-tRNA biosynthesis, prokaryotes” included one
179 variable gene at an empirical FDR of 5%, SepRS. SepRS is an O-phosphoseryl-tRNA synthetase,
180 which is an alternative route to biosynthesis of cysteinyl-tRNA in methanogenic archaea [19].
181 Methanogen abundance has previously been noted to be variable between individual human
182 guts: while DNA extraction for archaea may be less reliable than for bacteria, even optimized
183 methods showed large standard deviations across individuals [20]. Another gene in this cate-
184 gory was variable at a weaker level of significance (10% empirical FDR): PoxA, a variant lysyl-
185 tRNA synthetase. Recent experimental work has shown that this protein has a diverged, novel
186 functionality, lysinyllating the elongation factor EF-P [21, 22].

187 By comparison, 77% of the tested prokaryotic gene families in the KEGG module set “central
188 carbohydrate metabolism” were significantly invariable, and 5.6% (5 genes) were significantly
189 variable (Figure 3—figure supplement 1) at an empirical FDR of 5%. In this case, the variable
190 gene families highlight the complexities of microbial carbon utilization. Glucose can be metab-
191 olized by two alternative pathways: the well-known Embden-Meyerhof-Parnas (EMP) pathway
192 (i.e., classical “glycolysis”), or the Entner-Doudoroff pathway (ED). Both take glucose to pyru-
193 vate, but with differing yields of ATP and electron carriers; ED also allows growth on sugar acids
194 like gluconate [23]. Our analysis indicates that hosts differ in how much their gut microbial
195 communities use ED. While all genes in the “core module” of glycolysis dealing with 3-carbon
196 compounds were significantly invariable across individuals, we found that the ED-specific gene
197 family edd, which takes 6-phosphogluconate to 2-keto-3-deoxy-phosphogluconate (KDPG), was
198 significantly variable.

199 We also discovered significant variability in other enzymes involved in unusual sugar-phosphate
200 and tricarboxylic acid metabolism (Figure 3—figure supplement 1). Multifunctional and pri-
201 marily archaeal variants of fructose-bisphosphate aldolase (K16306, K01622) were significantly
202 variable across hosts, while the typical FBA enzyme (FbaA) was significantly invariable. Another
203 difference was seen in genes potentially contributing to ribose-phosphate generation. While
204 typical pentose-phosphate pathway genes (e.g., zwf and gnd) were invariable, the bifunctional
205 gene family Fae/Hps, thought to be involved in an alternative route to ribose-phosphate, was
206 significantly variable [24]. Finally, a subunit of fumarate reductase, frdD, was also significantly
207 variable. Fumarate reductase catalyzes the reverse reaction from the typical TCA cycle enzyme
208 succinate dehydrogenase and can be used for redox balance during anaerobic growth [25]. Con-
209 versely, the standard succinate dehydrogenase genes sdhA, sdhB and sdhC were significantly
210 invariable. These results suggest that using our test to identify variable genes within otherwise
211 invariable pathways can reveal diverged functionality as well as families that play domain or
212 clade-specific roles.

213 We found that the majority of significantly variable gene families annotated to “bacterial se-
214 cretion system” (16 out of 18) were involved in specialized secretion systems, especially the type
215 III and type VI systems (Figure 3). These secretion systems are predominantly found in Gram
216 negative bacteria and are often involved in specialized cell-to-cell interactions, between mi-
217 crobes and between pathogens or symbionts and the host. They allow the injection of effector
218 proteins, including virulence factors, directly into target cells [26, 27]. Type VI secretion systems
219 have also been shown to be determinants of antagonistic interactions between bacteria in the
220 gut microbiome [28, 29].

221 In contrast, gene families in the Sec (general secretion) and Tat (twin-arginine translocation)
222 pathways were nearly all significantly *invariable* at an empirical FDR of 5%, with only one gene
223 in each being found to be significantly variable. This contradicts previous suggestions that the
224 Sec and Tat pathways were some of the most variable in the human microbiome [13]. This
225 discrepancy is probably due to our accounting for the mean-variance relationship in shotgun
226 data; the Sec and Tat systems are abundant and phylogenetically diverse [30] and will therefore
227 have high variance just by chance compared to low-abundance genes. Our test adjusts for this
228 feature of sequencing experiments and shows that these genes are in fact less variable than
229 expected given their mean abundance.

230 Our results further demonstrate that analyzing functional variability at the level of pathways
231 can obscure gene-family-resolution trends of potential biomedical importance. The variabil-
232 ity of individual gene families involved in lipopolysaccharide (LPS) metabolism may exemplify
233 such a case (Figure 4). LPS (also known as “endotoxin”) is a macromolecular component of
234 Gram-negative bacterial outer membrane, consisting of a lipid anchor called “lipid A,” a “core

235 oligosaccharide” moiety, and a polysaccharide known as the “O-antigen” (which may be ab-
236 sent). Lipid A is sensed directly by the human innate immune system via the Toll-like receptor
237 TLR4. Furthermore, lipid A variants with different covalent modifications (e.g., differentially
238 acylated [31], phosphorylated [32], and palmitoylated [33] variants) have been shown to have
239 different immunological properties. Hexaacylated lipid A, as found in *E. coli*, stimulates TLR4
240 and induces the release of pro-inflammatory cytokines; conversely, pentaacylated lipid A vari-
241 ants, as found in *Bacteroides*, tend not to induce TLR4 signaling, and can even prevent the hex-
242 aacylated variety from inducing inflammation [34]. This inflammation may have a variety of
243 downstream effects on health. For example, elevated serum LPS levels are observed in obese
244 individuals [35, 36] and individuals with inflammatory bowel disease [35], and have been linked
245 to an increase in coronary heart disease events [37]. Conversely, a recent study advanced the
246 hypothesis that dampening of TLR4 signaling in childhood by *Bacteroides* species may actually
247 *increase* later susceptibility to autoimmune disease [34].

248 We found that all but one gene family involved in the biosynthesis of lipid A, as well as all
249 gene families involved in the biosynthesis of the core oligosaccharide components ketodeoxy-
250 octonate (Kdo) and glyceromannoheptose (GMH), were significantly invariable (16 out of 17).
251 The lone exception catalyzes the the final lipid A acylation step, adding a sixth acyl chain;
252 this gene family was significantly variable ($FDR \leq 5\%$). Furthermore, we observe several vari-
253 able gene families annotated as performing covalent modifications of LPS, including hydroxyl-
254 (LpxO), palmitoyl- (PagP), and palmitoleoylation (LpxP), as well as deacylation and dephospho-
255 rylation. Previous experimental work has shown that these modifications can lead to differen-
256 tial TLR4 activation [33, 38]. We also observe that gene families involved in O-antigen synthesis
257 and ligation to lipid A tended to be variable (5 out of 6). These results suggest that healthy in-
258 dividuals may differ in the amount of hexa- vs. pentaacylated LPS, and in the amounts of other
259 LPS chemical modifications, and thus in their baseline level of TLR4-dependent inflammation.
260 Importantly, since the majority of gene families annotated to LPS biosynthesis were invariable,
261 this result would have been missed by considering the pathway as a unit.

262 **2.4 Many invariable gene families are deeply conserved**

263 Conservation of gene families across the tree of life is one factor we might expect to affect gene
264 variability. For instance, ribosomal proteins should appear to be invariable merely because they
265 are shared by all members of a given kingdom of life. To explore the relationship between gene
266 family taxonomic distribution and variability in abundance across hosts, we constructed trees
267 of the sequences in each KEGG family using ClustalOmega and FastTree. We then calculated
268 phylogenetic distribution (PD), using tree density to correct for the overall rate of evolution [39]

269 (Figure 5a).

270 Overall, invariable gene families with below-median PD tended to be involved in carbohy-
271 drate metabolism and signaling. Specifically, these 2,046 gene families were enriched for the
272 pathways “two-component signaling” (FDR-corrected p-value $q = 1.5 \times 10^{-15}$), “starch and su-
273 crose metabolism” ($q = 1.8 \times 10^{-3}$), “amino sugar and nucleotide sugar metabolism” ($q = 0.063$),
274 “ABC transporters” ($q = 2.4 \times 10^{-5}$), and “glycosaminoglycan [GAG] degradation” ($q = 0.053$),
275 among others (Supplementary File 1). Enriched modules included a two-component system
276 involved in sporulation control ($q = 0.018$), as well as transporters for rhamnose ($q = 0.14$),
277 cellobiose ($q = 0.14$), and alpha- and beta-glucosides ($q = 0.14$ and $q = 0.19$, respectively).
278 These results are consistent with the hypothesis that one function of the gut microbiome is
279 to encode carbohydrate-utilization enzymes the host lacks [40]. Additionally, recent experi-
280 ments have also shown that the major gut commensal *Bacteroides thetaiotaomicron* contains
281 enzymes adapted to the degradation of sulfated glycans including GAGs [41, 42], and that many
282 *Bacteroides* species can in fact use the GAG chondroitin sulfate as a sole carbon source [43].

283 Out of the 298 significantly-variable gene families with above-median PD, we found no path-
284 way enrichments but three module enrichments. These included the archaeal ($q = 1.5 \times 10^{-3}$)
285 and eukaryotic ($q = 8.7 \times 10^{-9}$) ribosomes, which reflects differences in the relative abundance
286 of microbes from these domains of life across hosts (Figure 2b). The third conserved but vari-
287 able module was the type VI secretion system ($q = 0.039$). Intriguingly, specialized secre-
288 tion systems were also observed to vary within gut-microbiome-associated species in a strain-
289 specific manner, using a wholly separate set of data [44]. Finally, gene families described as
290 “hypothetical” were enriched in the high-PD but variable gene set ($p = 2.4 \times 10^{-8}$, odds ratio =
291 2.2) and depleted in the low-PD but invariable set ($p = 5.4 \times 10^{-13}$, odds ratio = 0.41).

292 Transporters were recently observed to show strain-specific variation in copy number across
293 different human gut microbiomes [44], and analyses by Turnbaugh et al. identified membrane
294 transporters as enriched in the “variable” set of functions in the microbiome [12]. However,
295 we mainly found transporters enriched amongst gene families with similar abundance across
296 hosts, despite being phylogenetically restricted (low-PD but invariable genes; Supplementary
297 File 2). Part of this difference is likely due to our stratifying by phylogenetic distribution, a step
298 previous studies did not perform.

299 **2.5 Proteobacteria are the major source of variable genes**

300 To assess which taxa contributed these variable and invariable genes, we first computed corre-
301 lations between phylum relative abundances (predicted using MetaPhlAn2 [45]) and gene fam-
302 ily abundances. This analysis revealed that the predicted abundance of Proteobacteria (and, to

303 a lesser extent, the abundance of the archaeal phylum Euryarchaeota) tended to be correlated
304 with variable gene families (Figure 6b).

305 Proteobacteria were a comparatively minor component of these metagenomes (median =
306 1%), compared to Bacteroidetes (median = 59%) and Firmicutes (median = 33%). However,
307 some hosts had up to 41% Proteobacteria. Overgrowth of Proteobacteria has been associated
308 with metabolic syndrome [46] and inflammatory bowel disease [47]. Also, Proteobacteria can
309 be selected (over Bacteroidetes and Firmicutes) by intestinal inflammation as tested by TLR5-
310 knockout mice [48], and some Proteobacteria can induce colitis in this background [49], po-
311 tentially leading to a feedback loop. Thus, the variable gene families we discovered could be
312 biomarkers for dysbiosis and inflammation in otherwise healthy hosts.

313 We also examined correlations between gene abundance and three taxonomic summary
314 statistics that have been previously linked to microbiome function: average genome size (AGS)
315 [18], the Bacteroidetes/Firmicutes ratio [12, 50], and α -diversity (Shannon index). All of these
316 statistics were *less* often correlated with variable gene families than with invariable or non-
317 significant gene families (see Supplementary File 7, Figure 6—figure supplement 1). These statis-
318 tics therefore do not explain the variability of gene families in this dataset.

319 Finally, previous research has suggested the existence of a small number of “enterotypes”
320 in the human gut microbiome, each with distinct taxonomic composition. A recent large-scale
321 study confirmed that abundances of the taxa Ruminococcaceae, Bacteroides, and Prevotella
322 explained the most taxonomic variation across individuals [51]. These enterotypes appear to
323 be linked to long-term diet, with Prevotella highest in individuals with the most carbohydrate
324 intake, and Bacteroides correlating with protein and animal fat. However, while these clades
325 contribute most to taxonomic variation, all were actually *depleted* for associations with vari-
326 able genes. In contrast, the Proteobacterial family Enterobacteriaceae was much more likely to
327 be associated with variable gene families (Figure 6—figure supplement 2). This suggests that
328 compared to previously-identified enterotype marker taxa, levels of Proteobacteria, and poten-
329 tially Euryarchaeota, better explain person-to-person variation in gut microbial gene function.
330 These less abundant phyla were missed in enterotype studies, likely because 1. enterotypes
331 were identified by methods that will tend to weight higher-abundance taxa more, and 2. en-
332 terotypes were identified from taxonomic, not functional data.

333 Because Proteobacteria are a relatively well annotated yet low abundance phylum, we ex-
334 plored whether either of these characteristics drive their association with variable genes. Im-
335 portantly, genes correlated with Actinobacteria did not tend to be variable, even though Pro-
336 teobacteria and Actinobacteria had similar levels of abundance (minimum 0%, median 1%,
337 maximum 20%). Thus, phylum prevalence and abundance do not explain the variability of Pro-
338 teobacterial genes. To investigate annotation bias, we first compared the numbers of genomes

339 in KEGG for each phylum. There are 1,111 Proteobacterial genomes compared to 575 for Firmi-
340 cutes, 276 for Actinobacteria, and only 97 for Bacteroidetes. Proteobacteria consequently had
341 the most “private” gene families not annotated in any other phylum (1,417), compared to 538
342 for Firmicutes, 342 for Euryarchaeota, 215 for Actinobacteria, and 21 for Bacteroidetes. Con-
343 sidering only these private gene families, Proteobacteria and Euryarchaeota were enriched for
344 variable genes, as before, whereas variable genes were depleted in the other three phyla (Figure
345 6—figure supplement 3). This suggests that the level of annotation does not predict the amount
346 of variable genes. In a further test, we repeated the entire statistical test on a subset of genes,
347 sampling one part phylum-specific genes drawn equally from Proteobacteria, Actinobacteria,
348 Firmicutes, and Euryarchaeota, and one part genes annotated to all four phyla (see Methods).
349 Again, Proteobacteria- and Euryarchaeota-specific genes were significantly variable more of-
350 ten than those from either Actinobacteria or Firmicutes (Figure 6—figure supplement 4). We
351 therefore conclude that phylum abundance and annotation bias do not drive the enrichment
352 of variable genes in Proteobacteria.

353 **2.6 Bacterial phyla have unique sets of variable genes**

354 The variable gene families we identified seem to include both genes whose variance is explained
355 by phylum-level variation (e.g., Proteobacteria), and genes that vary within fine-grained tax-
356 onomic classifications, such as strains within species. Also, some gene families may confer
357 adaptive advantages in the gut only within certain taxa. To detect gene families that are vari-
358 able or invariable within a phylum, we repeated the test, but using only reads that mapped best
359 to sequences from each of the four most abundant bacterial phyla (Bacteroidetes, Firmicutes,
360 Actinobacteria, and Proteobacteria). Most (77%) gene families showed phylum-specific effects.
361 Invariable gene families tended to agree, but the reverse was true for variable gene families:
362 19.4% of gene families that were invariable in one phylum were invariable in all, compared to
363 just 0.34% (8 genes) in the variable set (Figure 7A-B). This trend was robust to the FDR cutoff
364 (Figure 7—figure supplement 1). Gene families invariable in all four phyla were enriched for
365 basal cellular machinery, as expected (Supplementary File 3).

366 The relationship between phylum-specific and overall gene family abundance variability
367 differed by phylum. Proteobacteria-specific variable gene families tended to be variable overall
368 (59%), whereas the proportions of gene families that were also variable overall were much lower
369 for Bacteroidetes- (12%), Firmicutes- (29%), and Actinobacteria-specific (18%) gene families
370 (Figure 7C). This supports the hypothesis that Proteobacterial abundance is a dominant driver
371 of functional variability in the human gut microbiome. It further suggests that many overall-
372 variable gene families are not merely markers for the amount of Proteobacteria (or some other

373 phylum), but are also variable at finer taxonomic levels, such as the species or even the strain
374 level [44, 52].

375 Comparing the two dominant phyla in the gut, Bacteroidetes and Firmicutes, we further ob-
376 serve that the overall proportions of variable and invariable families were similar across path-
377 ways, with some interesting exceptions. For example, lipopolysaccharide (LPS) biosynthesis
378 had many invariable gene families in Bacteroidetes and very few in Firmicutes, which we ex-
379 pected given that LPS is primarily made by Gram-negative bacteria. Conversely, both two-
380 component signaling and the PTS system had many more invariable gene families in Firmi-
381 cutes than in Bacteroidetes (Figure 7—figure supplement 2A). However, phylum-specific vari-
382 able gene families tended not to overlap (median overlap: 0%, compared to 46% for invariable
383 gene families). This was even true for pathways where the overall proportion of variable and
384 invariable gene families is similar, such as cofactor and vitamin biosynthesis and central car-
385 bohydrate metabolism (Figure 7—figure supplement 2B). Thus, unique genes within invariable
386 pathways vary in their abundance across microbiome phyla.

387 Furthermore, the enriched biological functions of the phylum-specific variable gene fam-
388 ilies differed by phylum (Supplementary File 4). For instance, Proteobacterial-specific vari-
389 able gene families were enriched (Fisher’s test enrichment $q = 0.13$) for the biosynthesis of
390 siderophore group nonribosomal peptides, which may reflect the importance of iron scaveng-
391 ing for the establishment of both pathogens (e.g. *Yersinia*) and commensals (*E. coli*) [53]. An-
392 other phylum-specific variable function appeared to be the Type IV secretion system (T4SS)
393 within Firmicutes ($q = 0.021$). Homologs of this specialized secretion system have been shown
394 to be involved in a wide array of biochemical interactions, including the conjugative transfer of
395 plasmids (e.g. antibiotic-resistance cassettes) between bacteria [54]. We conclude that our ap-
396 proach enables the identification of substantial variation within all four major bacterial phyla in
397 the gut, much of which is not apparent when data are analyzed at broader functional resolution
398 or without stratifying by phylum.

399 **2.7 Variable genes are not biomarkers for body mass index, sex or age**

400 To explore associations of gene variability with measured host traits, we used a two-sided par-
401 tial Kendall’s τ test that controls for study effects (Methods). Body mass index, sex, and age were
402 measured in all three studies we analyzed. None of these variables correlated significantly with
403 any variable gene family abundances, even at a 25% false discovery rate. This suggests that ma-
404 jor correlates of variation in microbiota gene levels, possibly including diet and inflammation,
405 were not measured in these studies.

406 **3 Discussion**

407 This study presents a novel statistical method that provides a finer resolution estimate of “func-
408 tional redundancy” [55] in the human microbiome than was previously possible. Our test differs
409 from previous approaches to quantifying variability in microbiome function in several key ways.
410 First, we focus explicitly on the variability of gene family abundance, not differences in mean
411 abundance between predefined groups, as has been done to reveal pathways whose abundance
412 differs between body sites [56] or disease states [6]. Second, we take a finer-grained and more
413 quantitative approach to measuring variability of microbiome functions than the studies that
414 initially observed that biological pathways are relatively invariable [13, 12]. Our work identifies
415 individual gene families that break this overall trend. A third important aspect of our method is
416 that the underlying model accounts for the mean-variance relationship in count data, as well as
417 systematic biases between studies. Finally, our null distribution is estimated from the shotgun
418 data and does not require comparisons to sequenced genomes.

419 We found that basic microbial cellular machinery, such as the ribosome, tRNA-charging,
420 and primary metabolism, were universal functional components of the microbiome, both in
421 general and when each individual phylum was considered separately. This finding is consistent
422 with previous results [12], and indeed, is not surprising given the broad conservation of these
423 processes across the tree of life. In contrast, we also identified invariable gene families that
424 have narrower phylogenetic distributions. These included, for example, proteins involved in
425 two-component signaling, starch metabolism (including glucosides), and glycosaminoglycan
426 metabolism. Previous experimental work has underscored the importance of some of these
427 pathways in gut symbionts: for instance, multiple gut-associated *Bacteroides* species are ca-
428 pable of using the glycosaminoglycan chondroitin sulfate as a sole carbon source [41], and the
429 metabolism of resistant starch in general is thought to be a critical function of the omnivorous
430 mammalian microbiome [40]. These results suggest that the method we present is capable of
431 identifying protein-coding gene families that contribute to fitness of symbionts within the gut.
432 Finally, we found a number of invariable gene families whose function is not yet annotated.
433 These gene families may represent functions that are either essential or provide advantages for
434 life in the gut, and may therefore be particularly interesting targets for experimental follow-up
435 (e.g., assessing whether strains in which these gene families have been knocked out in fact have
436 slower growth rates, either in vitro or in the gut).

437 We also identified significantly variable gene families, including enzymes involved in carbon
438 metabolism, specialized secretion systems such as the T6SS, and lipopolysaccharide biosyn-
439 thetic genes. Proteobacteria, rather than Bacteroidetes or Firmicutes, emerge as a major source
440 of variable genes, including some genes whose abundance also varied within the Proteobacteria

441 (e.g., T6SS). Since Proteobacteria have been linked to inflammation and metabolic syndrome
442 [46], we speculate that baseline inflammation may be one variable influencing functions in the
443 gut microbiome. Some variable genes, including many of unknown function, had surprisingly
444 broad phylogenetic distributions.

445 Variable gene families have a variety of ecological interpretations, e.g., first-mover effects,
446 drift, host demography, and selection within particular gut environments. Computationally
447 distinguishing among these possibilities is likely to present challenges. For example, distin-
448 guishing selection from random drift will probably require longitudinal data and appropriate
449 models. Separating effects of host geography, genetics, medical history, and lifestyle will be
450 possible only when richer phenotypic data is available from a more diverse set of human pop-
451 ulations. To control for study bias and batch effects, it will be important to include multiple
452 sampling sites within each study.

453 While statistical tests focused on differences in variances are not yet common throughout
454 genomics, there is some recent precedent using this type of test to quantify the gene-level het-
455 erogeneity in single-cell RNA sequencing data [57, 58], and to identify variance effects in genetic
456 association data [59]. Like Vallejos et al. [58], we model gene counts using the negative bino-
457 mial distribution, and identify both significantly variable and invariable genes; in contrast, we
458 frame our method as a frequentist hypothesis test as opposed to a Bayesian hierarchical model.
459 Our method also accounts for study-to-study variation. Also, unlike previous approaches in
460 this domain, the method we describe does not require biological noise to be explicitly decom-
461 posed from technical noise; our method therefore does not require the use of experimentally-
462 spiked-in controls, which are not present in most experiments involving sequencing of the gut
463 microbiome. Instead, we detect differences from the average level of variability using a robust
464 nonparametric estimator, which we show through simulation leads to correct inferences under
465 reasonable assumptions.

466 A similar statistical method for detecting significant (in)variability such as the one we present
467 here could also be applied to other biomolecules measured in counts, such as metabolites, pro-
468 teins, or transcripts. Performing such analyses on human microbiota would reveal patterns in
469 the variability in the usage of particular genes, reactions, and pathways, which would expand on
470 our investigation of potential usage based on presence in the DNA of organisms in host stool.
471 Another important extension is to generalize our method for comparing hosts from different
472 pre-defined groups (disease states, countries, diets) to identify gene families that are invariable
473 in one group (e.g., healthy controls) but variable in another (e.g., patients), analogously to re-
474 cent methods for the analysis of single-cell RNA-Seq [60] and GWAS [59] data. In particular,
475 gene families whose variance differs between case and control populations could point to het-
476 erogeneity within complex diseases, interactions between the microbiome and latent variables

477 (e.g., environmental or genetic), and/or differences in selective pressure between healthy and
478 diseased guts. Investigating group differences in functional variability could thereby allow the
479 detection of different trends from the more common comparison of means.

480 **4 Materials and Methods**

481 **4.1 Data collection and processing**

482 Stool metagenomes from healthy human guts were obtained from three sources:

- 483 1. two American cohorts from the Human Microbiome Project [13], $n = 42$ samples selected;
- 484 2. a Chinese cohort from a case-control study of type II diabetes (T2D) [15], $n = 44$ samples
485 from controls with neither type II diabetes nor impaired glucose tolerance;
- 486 3. and a European cohort from a case-control study of glucose control [16], $n = 37$ samples
487 from controls with normal glucose tolerance.

488 Samples were chosen to have at least 1.5×10^7 reads and mode average quality scores ≥ 20 (esti-
489 mated via FastQC [61]). After downloading these samples from NCBI's Sequence Read Archive
490 (SRA), the FASTA-formatted files were mapped to KEGG Orthology (KO) [62] protein families
491 as previously described [17]. For consistency, each sample was rarefied to a depth of 1.5×10^7
492 reads, and additionally, as reads from HMP were particularly variable in length, they were there-
493 fore trimmed to a uniform length of 90 bp.

494 For each sample, we used ShotMAP to detect how many times a particular gene family
495 matched a read ("counts"; we added one pseudocount for reasons described below). The bit-
496 score cutoff for matching a protein family was selected based on the average read length of
497 each sample as recommended previously [17]. For every gene family in every sample, we also
498 computed the average family length (AFL), or the average length of the matched genes within
499 a family. Finally, we also computed per-sample average genome size using MicrobeCensus
500 [18] (<http://github.com/snayfach/MicrobeCensus>). These quantities were used to esti-
501 mate abundance values in units of RPKG, or reads per kilobase of genome equivalents [18].

502 These RPKG abundance values were strictly positive with a long right tail and highly corre-
503 lated with the variances (Spearman's $r = 0.99$). This strong mean-variance relationship is likely
504 simply because these abundances are derived from counts that are either Poisson or negative-
505 binomially distributed. We therefore took the natural log of the RPKG values as a variance stabi-
506 lizing transformation. Because $\log(0)$ is infinite, we added a pseudocount before normalizing
507 the counts and taking the log transform. Since there is no average family length (AFL) when

508 there are no reads for a given gene family in a given sample, we imputed it in those cases using
509 the average AFL across samples.

510 4.2 Model fitting

511 We fit a linear model to the data matrix of log-RPKG D of log-RPKG described above, with n
512 gene-families by m samples, to capture gene-specific and dataset-specific effects:

$$D_{g,s} = \mu_g + \sum_{y \in Y} I_{y,s} \beta_{g,y} + \epsilon_{g,s} \quad (2)$$

513 where $g \in [1, n]$ is a particular gene family, $s \in [1, m]$ is a particular sample, μ_g is estimated by the
514 grand or overall mean of log-RPKG $\frac{\sum_s D_{g,s}}{m}$ for a given gene family g , Y is the set of studies, $I_{y,s}$
515 is an indicator variable valued 1 if sample s is in study y and 0 otherwise, $\beta_{g,y}$ is a mean offset
516 for gene family g in study y , and the residual for a given gene family and sample are given by
517 $\epsilon_{g,s}$. For each gene family, the variance across samples of these $\epsilon_{g,s}$, which we term the “residual
518 variance” or V_g^ϵ , was our statistic of interest.

519 Overall trends in these data are explained well by this model, with an $R^2 = 0.20$. The resid-
520 uals, which are approximately symmetrically distributed around 0, represent variation in gene
521 abundance not due to study effects.

522 4.3 Modeling residual variances under the null distribution

523 Having calculated this statistic V_g^ϵ for each gene family g , we then needed to compare this statis-
524 tic to its distribution under a null hypothesis H_0 . This required us to model what the data would
525 look like if in fact there were no surprisingly variable or invariable gene families. To do this, we
526 used the negative binomial distribution to model the original count data (before adding pseu-
527 docounts and normalization to obtain RPKG).

528 The negative binomial distribution is commonly used to model count data from high through-
529 put sequencing. It can be conceptualized as a mixture of Poisson distributions with different
530 means, which themselves follow a Gamma distribution. Like the Poisson distribution, the neg-
531 ative binomial distribution has an intrinsic mean-variance relationship. However, instead of a
532 single mean-variance parameter as in the Poisson, the negative binomial can be described with
533 two, a mean parameter and a “size” parameter, which we refer to here as k such that $k = \frac{\mu^2}{\sigma^2 - \mu}$.
534 k ranges from $(0, \infty)$, with smaller values corresponding to more overdispersion (i.e., higher
535 variance given the mean) and larger values approaching, in the limit, the Poisson distribution.

536 To model the case where no gene family has unusual variance given its mean value, i.e.,
537 our null hypothesis, we assumed that the data were negative-binomially distributed with the

538 observed means $\mu_{g,y}$ for each gene g and study y , but where the amount of overdispersion
 539 was modeled with a single size parameter k_y for each study y . This has similarities to previous
 540 approaches to model RNAseq distributions [63] and to identify (in)variable genes from single-
 541 cell RNAseq data [58] (see also Discussion).

$$H_0: V_g^e = V_g^e | D_{g,s} \sim NB(\mu_{g,y}, k_y)$$

$$H_{alt}: V_g^e \neq V_g^e | D_{g,s} \sim NB(\mu_{g,y}, k_y)$$

542 To estimate this \widehat{k}_y , the overall size parameter for a given study y , we estimated the mode of per-
 543 gene-family size parameters $k_{g,y}$ within data set y , using the method-of-moments estimator
 544 for each $k_{g,y}$. We accomplished this by fitting a Gaussian kernel density estimate to the log-
 545 transformed $k_{g,y}$ values, and then finding the \widehat{k}_y value that gave the highest density. (From
 546 simulations, we found that the mode method-of-moments was more robust than the median
 547 or harmonic mean: see Figure 1—figure supplement 2.) We could then easily generate count
 548 data under this null distribution, add a pseudocount and normalize by AFL and AGS, fit the
 549 above linear model, and obtain null residual variances $V_g^{\epsilon_0}$ using exactly the same procedure
 550 described above.

551 Statistical significance was obtained by a two-tailed test:

$$p_g = \frac{\# \left(\left(\frac{V_g^{\epsilon_0} - \overline{V_g^{\epsilon_0}}}{V_g^{\epsilon_0}} \right)^2 \geq \left(\frac{V_g^e - \overline{V_g^e}}{V_g^{\epsilon_0}} \right)^2 \right) + 1}{B + 1}$$

552 Here, B refers to the number of null test statistics $V_g^{\epsilon_0}$ (in this case, $B = 750$), and the over-
 553 lined test statistics refer to their mean across the null distribution.

554 The resulting p-values were then corrected for multiple testing by converting to FDR q-
 555 values using the procedure of Storey et al. [64] as implemented in the `qvalue` package in R
 556 [65]. An alternative approach to determining significance is based on the bootstrap. While us-
 557 ing a parametric null distribution allows us to explicitly model the null hypothesis, it also breaks
 558 the structure of covariance between gene families, which may be substantial because genes are
 559 organized into operons and individual genomes within a metagenome. This structure can, op-
 560 tionally, be restored using a strategy outlined by Pollard and van der Laan [66]. Instead of using
 561 the test statistics $V_g^{\epsilon_0}$ obtained under the parametric null as is, we can use these test statistics to
 562 center and scale bootstrap test statistics $V_g^{\epsilon'}$, which we derive from applying a cluster bootstrap
 563 with replacement from the real data and then fitting the above linear model (2) to the resampled
 564 data to obtain bootstrap residual variances:

$$V_g^{\epsilon_{0'}} = \left(\left(\frac{V_g^{\epsilon'} - \overline{V}_g^{\epsilon'}}{sd(V_g^{\epsilon'})} \right) \times sd(V_g^{\epsilon_0}) \right) + \overline{V}_g^{\epsilon_0}$$

565 A similar non-parametric bootstrap approach has previously been successfully applied to test-
566 ing for differences in gene expression [67].

567 4.4 Visualization

568 As expected, when the residuals are plotted in a heatmap as in Figure 2—figure supplement 2,
569 variable gene families were generally brighter (i.e., more deviation from the mean) than in-
570 variable gene families, though not exclusively: this is because our null distribution, unlike the
571 visualization, models the expected mean-variance relationship. We visualized this information
572 by scaling each gene family by its expected standard deviation under the negative binomial null
573 (i.e., by the mean root variance $\sum_{b \in [1, B]} \sqrt{V_{g_b}^{\epsilon_0} / B}$) (Figure 2—figure supplement 3).

574 In Figure 3, for comparability with existing literature, gene families in the T6SS were named
575 by mapping to the COG IDs used in Coulthurst [27], except when multiple KOs mapped to
576 the same COG ID; in these cases, the original KO gene names were kept. Schematics of the
577 T3SS, T6SS, Tat, and Sec pathways were modeled on previous reviews [68, 69, 27] and on the
578 KEGG database [62]. The pathway diagram in Figure 4 is based on representations in the KEGG
579 database [62], MetaCyc [70], and reviews by Wang and Quinn [71] and Whitfield and Trent [72].
580 These reviews were also used to identify KEGG Orthology gene families that were involved in
581 lipopolysaccharide metabolism but not yet annotated under that term.

582 4.5 Power analysis

583 The test we present controls α as expected if the correct size parameter k is estimated from the
584 data (Figure 1—figure supplement 2d-e). Estimating this parameter accurately is known to be
585 difficult, however, particularly for highly over-dispersed data [73], and in this case we must also
586 estimate this parameter from a mixture of true positives and nulls. We found that the mode
587 of per-gene-family method-of-moments estimates was more robust to differences in the ratio
588 of variable to invariable true positives (Figure 1—figure supplement 2a-c) than the median or
589 harmonic mean (the harmonic mean mirrors the approach in Yu et al. [63]).

590 Power analysis was performed on simulated datasets comprising three simulated studies.
591 For each study, 1,000 gene families were simulated over $n \in \{60, 120, 480, 960\}$ samples. Null
592 data were drawn from a negative-binomial distribution with a randomly-selected size param-
593 eter k in common to all gene families, which was drawn from a log-normal distribution (log-
594 mean= -0.65 , sd= 0.57). Gene family means were also drawn from a log-normal (log mean=

595 2.94, $sd = 2.23$). True positives were drawn from a similar negative-binomial distribution, but
596 where the size parameter was multiplied by an effect size z (for variable gene families) or its
597 reciprocal $1/z$ (for invariable gene families). The above test was then applied to the simulated
598 data, and the percent of Type I and II errors was calculated by comparing to the known gene
599 family labels from the simulation. Using similar parameters to those estimated from our real
600 data, we saw that α decreased and power approaches 1 with increasing sample size (see Figure
601 1—figure supplement 3) and that $n = 120$ appears to be sufficient to achieve control over α .

602 However, at $n = 120$, we also noted that α appeared to be greater for variable vs. invari-
603 able gene families (Figure 1—figure supplement 4), possibly because accurately detecting addi-
604 tional overdispersion in already-overdispersed data may be intrinsically difficult. We therefore
605 performed additional simulations to determine q -value cutoffs corresponding to an empirical
606 FDR of 5%. We calculated appropriate cutoffs based on datasets with 43% true positives and
607 a variable:invariable gene family ratio ranging from 0.1 to 10, taking the median cutoff value
608 across these ratios (Supplementary File 1). Using these cutoffs, the overall dataset had 45% true
609 positives and a variable:invariable gene family ratio of 0.43.

610 **4.6 Calculating phylogenetic distribution of gene families**

611 The phylogenetic distribution (PD) of KEGG Orthology (KO) families was estimated using tree
612 density [39]. We first obtained sequences of each full-length protein annotated to a particular
613 KO, and then performed a multiple alignment of each family using ClustalOmega [74]. These
614 multiple alignments were used to generate trees via FastTree [75]. For both the alignment and
615 tree-building, we used default parameters for homologous proteins.

616 For all families represented in at least 5 different archaea and/or bacteria (6,703 families to-
617 tal), we then computed tree densities, or the sum of edge lengths divided by the mean tip height.
618 Using tree density instead of tree height as a measure of PD corrects for the rate of evolution,
619 which can otherwise cause very highly-conserved but slow-evolving families like the ribosome
620 to appear to have a low PD [39]. Empirically, this measure is very similar to the number of pro-
621 tein sequences (Figure 5—figure supplement 1), but is not as sensitive to high or variable rates
622 of within-species duplication: for example, families such as transposons, which exhibit high
623 rates of duplication as well as copy-number variation between species, have a larger number of
624 sequences than even very well-conserved proteins such as RNA polymerase, but have similar
625 or even lower tree densities, indicating that they are not truly more broadly conserved.

626 Many protein families (8,931 families) did not have enough observations in order to reliably
627 calculate tree density, with almost all of these being annotated in only a single bacterium/archaeum.
628 For these, we predicted their PD by extrapolation. To predict PD, we used a linear model that

629 predicted tree density based on the total number of annotations (including annotations in eu-
630 karyotes). In five-fold cross-validation, this model actually had a relatively small mean absolute
631 percentage error (MAPE) of 13.1%. We also considered a model that took into account the taxo-
632 nomic level (e.g., phylum) of the last common ancestor of all organisms in which a given protein
633 family was annotated, but this model performed essentially identically (MAPE of 13.0%). Pre-
634 dicted tree densities are given in Supplementary File 6. The PD of gene families varied from 1.2
635 (an iron-chelate-transporting ATPase only annotated in *H. pylori*) to 434.9 (the rpoE family of
636 RNA polymerase sigma factors).

637 **4.7 Gene family enrichment**

638 We were interested in whether particular pathways were enriched in several of the gene family
639 sets identified in this work. For subsets of genes (such as those with specifically low PD), a 2-
640 tailed Fisher's exact test (i.e., hypergeometric test) was used instead to look for cases in which
641 the overlap between a given gene set and a KEGG module or pathway was significantly larger
642 or smaller than expected. The background set was taken to be the intersection of the set of
643 gene families observed in the data with the set of gene families that had pathway- or module-
644 level annotations. p -values were converted to q -values as above. Finally, enrichments were
645 enumerated by selecting all modules or pathways below $q \leq 0.25$ that had positive odds-ratios
646 (i.e., enriched instead of depleted).

647 **4.8 Associations with clinical and taxonomic variables**

648 We were interested in using a non-parametric approach to detect association of residual RPKG
649 with clinical and taxonomic variables (e.g., the inferred abundance of a particular phylum or
650 other clade via MetaPhlan2). To take into account potential study effects in clinical and taxo-
651 nomic variables without using a parametric modeling framework, we used partial Kendall's τ
652 correlation as implemented in the ppcor package for R [76], coding the study effects as binary
653 nuisance variables. Kendall's τ was used over Spearman's ρ because of better handling of ties
654 (an issue with taxonomic variables especially, since many, particularly at the finer-grained lev-
655 els, were often zero). The null distribution was obtained by permuting the clinical/taxonomic
656 variables within each study 250 times, and then re-assessing the partial τ . Finally, p -values were
657 calculated by taking the fraction of null partial correlations equally or more extreme (i.e., distant
658 from zero) than the real partial correlations.

659 Taxonomic relative abundances were predicted from the shotgun data using MetaPhlan2
660 with the `--very-sensitive` flag [45].

661 Two approaches were used to test for annotation bias. First (Figure 6—figure supplement 3),

662 gene families private to a phylum (i.e., those annotated in only a single bacterial/archaeal phy-
663 lum) were identified from the KEGG database. We then tested whether these private gene fam-
664 ilies were enriched or depleted for significantly variable gene families (5% FDR) using Fisher’s
665 exact test. Second (Figure 6—figure supplement 4), we performed a test in which we sampled
666 215 private gene families from each of Proteobacteria, Firmicutes, Actinobacteria, and Eur-
667 yarchaeota, totaling 860, plus 860 gene families annotated in all four phyla. (Since Bacteroidetes
668 only had 21 private genes, that phylum was dropped from this analysis.) Enrichment/depletion
669 for variable gene families within each phylum was performed as above.

670 **4.9 Phylum-specific tests**

671 We created taxonomically-restricted data sets in which the abundance of each gene family was
672 computed using only metagenomic reads aligning best to sequences from each of the four
673 most abundant bacterial phyla (Bacteroidetes, Firmicutes, Actinobacteria, and Proteobacte-
674 ria). Phylum-specific data were obtained from the overall data as follows. First, the NCBI taxon-
675 omy was parsed to obtain species annotated below each of the four major bacterial phyla (Bac-
676 teroidetes, Firmicutes, Actinobacteria, and Proteobacteria); these species were then matched
677 with KEGG species identifiers. Next, the original RAPSearch2 [77] results were filtered, so that
678 the only reads remaining were those for which their “best hit” in the KEGG database originally
679 came from the genome of a species belonging to the specific phylum in question (e.g., *E. coli* for
680 Proteobacteria). Finally, when performing the test, normalization for average genome size was
681 accomplished by normalizing gene family counts by the median abundance of a set of 29 bac-
682 terial single-copy gene families [78], which had been filtered in the same phylum-specific way
683 as all other gene families; this approach is similar to the MUSiCC method for average genome
684 size correction [79]. This also controls for overall changes in phylum abundance. Finally, we
685 estimated the average level of overdispersion \widehat{k}_y for individual studies based on the full dataset
686 (not phylum-restricted), since the expectation that < 50% of gene families were differentially
687 variable might not hold for each individual phylum. We used the same q -value cutoffs as in
688 the overall test to set an estimated empirical FDR (Table 1). Otherwise, tests were performed as
689 above.

690 **4.10 Codebase**

691 The scripts used to conduct the test and related analyses are available at the following URL:

692 <http://www.bitbucket.org/pbradz/variance-analyze>

693 Counts of reads mapped to KEGG Orthology (KO) groups and average family lengths for all
694 of the samples used in this study can be obtained at FigShare:

- 695 • <https://figshare.com/s/fcf1abf369155588ae41> (overall)
- 696 • <https://figshare.com/s/90d44cffdfb1d214ef83> (phylum-specific)

697 **5 Author contributions**

698 PHB performed the experiments and analyses. PHB and KSP developed the test, designed the
699 experiments, wrote the paper, and read and approved the final manuscript.

700 **6 Declarations**

701 **6.1 Acknowledgements**

702 The authors would like to thank Stephen Nayfach for downloading and organizing metage-
703 nomic data and metadata, and for providing and checking code for metagenome annotation,
704 Dongying Wu for suggesting the tree density metric to measure phylogenetic distribution, Aram
705 Avila-Herrera for help with phenotype-to-abundance associations, and Clifford Anderson-Bergman,
706 other members of the Pollard group, and Peter Turnbaugh for helpful discussions.

707 **6.2 Information about HMP clinical data**

708 Clinical covariates for HMP were obtained from dbGaP accession #phs000228.v3.p1. Funding
709 support for the development of NIH Human Microbiome Project - Core Microbiome Sampling
710 Protocol A (HMP-A) was provided by the NIH Roadmap for Medical Research. Clinical data
711 for HMP-A were jointly produced by the Baylor College of Medicine and the Washington Uni-
712 versity School of Medicine. Sequencing data for HMP-A were produced by the Baylor College
713 of Medicine Human Genome Sequencing Center, The Broad Institute, the Genome Center at
714 Washington University, and the J. Craig Venter Institute. These data were submitted by the
715 EMMES Corporation, which serves as the clinical data collection site for the HMP. Authors read
716 and agreed to abide by the Genomic Data User Code of Conduct.

717 7 Figures

718

719 **Figure 1: The residual variance statistic captures variation in gene families after account-**
720 **ing for between-study variation.** Left panels (“original abundances”) show filled circles repre-
721 senting log-RPKG abundances for gene families from the KEGG Orthology (KO), with per-study
722 means shown in solid horizontal lines and the distance from these means shown as dashed
723 vertical lines. Right-hand panels (“residuals”) show the same gene families after fitting a linear
724 model that accounts for these per-study means, with an accompanying density plot showing
725 the distribution of these residuals. V_g^c values in bold underneath density plots are the calcu-
726 lated variances of these residuals. These gene families are sets of orthologs corresponding to
727 the genes A) *tatA*, B) *devR*, c) *waaW*, d) *thrC*, E) *gspA*, F) *tssB*, G) *dctS*, and H) *ecnB*. Panels A-B
728 show two invariable gene families with relatively high (A) and low (B) average abundance; sim-
729 ilarly, panels C-D show two variable gene families with relatively low (C) and high (D) relative
730 abundances. Panels E-F show two gene families involved in secretion with similar abundances,
731 but low (E) vs. high (F) variability. Finally, panels G-H show that both invariable (G) and variable
732 (H) gene families can have substantial study-specific effects. (All gene families displayed were
733 significantly (in)variable using the test we present, $FDR \leq 5\%$.)

734 **Figure 2: Most pathways include a mixture of both variable and invariable gene families.** A)
735 Stacked bar plots show the fraction of invariable (blue), non-significant (gray), and variable
736 (red) gene families annotated to KEGG Orthology pathway sets (rows), at different false discov-
737 ery rate (FDR) cutoffs (color intensity). Only gene families with at least one annotated bacte-
738 rial or archaeal homolog are counted. B) Fraction of strongly invariable, non-significant, and
739 strongly variable gene families within the ribosomes of different kingdoms. Row labels with
740 only one kingdom indicate gene families unique to that kingdom, while rows with multiple
741 kingdoms (e.g. “Eukaryotes/archaea”) indicate gene families shared between these two king-
742 doms. As expected, the bacterial ribosome was completely invariable.

743 **Figure 3: Variable and invariable gene families involved in bacterial secretion separate by**
744 **gene function.** A) Schematic diagram showing the type III (T3SS), type VI (T6SS), Sec, and
745 Tat secretion system gene families measured in this dataset. Gene families are color-coded by
746 whether they were variable (red), invariable (blue), or neither (gray), with strength of color cor-
747 responding to the FDR cutoff (color intensity). Insets show a summary of how many gene fam-
748 ilies in KEGG modules corresponding to a particular secretion system are variable or invariable
749 and at what level of significance. B) Heatmaps showing scaled residual log-RPKG for gene fam-
750 ilies (rows) involved in bacterial secretion. Variable (red) and invariable (blue) gene families are
751 clustered separately, as are samples within a particular study (columns). log-RPKG values are
752 scaled by the expected variance from the negative-binomial null distribution. Genes in specific
753 secretion systems are annotated with colored squares (T6SS: red-orange; T3SS: orange; Tat: yel-
754 low; Sec: teal).

755 **Figure 4: Central Kdo and lipid A biosynthesis is invariable, but many genes involved in**
756 **covalent modifications to LPS are variable.** A) Pathway schematic showing a selection of mea-
757 sured gene families involved in lipopolysaccharide metabolism. Gene families are color-coded
758 by whether they were variable (red) or invariable (blue), with strength of color correspond-
759 ing to the FDR cutoff (color intensity). Central Kdo and lipid A metabolism is highlighted
760 in light grey. Abbreviated metabolites are GlcNAc (N-acetylglucosamine), Kdo (ketodeoxyoc-
761 tonate), ribose-5-phosphate (R5P), sedoheptulose-7-phosphate (S7P), and glyceromannohep-
762 tose (GMH). Aminoarabinose refers to 4-amino-4-deoxy-L-arabinose. B) Heatmaps showing
763 scaled residual *log*-RPKG for gene families (rows) involved in lipopolysaccharide metabolism,
764 as in Figure 3.

765 **Figure 5: Phylogenetic distribution (PD) of gene families partially explains gene family vari-**
766 **ability.** Scatter plot shows \log_{10} PD (x-axis) vs. \log_{10} residual variance statistic (y-axis). Red
767 points are significantly variable while blue points are significantly invariable. Gene families in
768 specific functional groups are also highlighted in different colors, specifically the bacterial ribo-
769 some (green), the type VI secretion system (or “T6SS”; orange), the KinABCDE-Spo0FA sporu-
770 lation control two-component signaling system (yellow), and hypothetical genes (tan squares).
771 Gene families that were significantly invariable (ribosome and sporulation control) or signifi-
772 cantly variable (hypothetical genes and the T6SS) at an estimated 5% FDR are outlined in black.
773 The bacterial ribosome, as expected, had very high PD and is strongly invariable. The Type VI
774 secretion system genes, in contrast, were conserved but variable, while some genes involved
775 in the Kin-Spo sporulation control two-component signaling pathway have low PD but were
776 invariable. Only gene families with at least one annotated bacterial or archaeal homolog are
777 shown.

778 **Figure 6: Variable gene families correlate with the predicted abundance of Proteobacte-**
779 **ria.** Bar plots give the fraction of gene families in each category (significantly invariable, non-
780 significant, and significantly variable, 5% FDR) that were significantly correlated to predicted
781 relative abundances of phyla, as assessed by MetaPhlan2, using partial Kendall’s τ to account
782 for study effects and a permutation test to assess significance. Asterisks give the level of signif-
783 icance by chi-squared test of non-random association between gene family category and the
784 number of significant associations. (***: $p \leq 10^{-8}$ by chi-squared test after Bonferroni correc-
785 tion; **: $p \leq 10^{-4}$.)

786 **Figure 7: Phylum-specific tests reveal hidden variability in the most prevalent bacterial**
787 **phyla.** A-B) Venn diagrams showing the number of significantly variable (A) and invariable (B)
788 gene families across Proteobacteria, Bacteroidetes, and Firmicutes, FDR \leq 5%. C) Bars indicate
789 the fraction of phylum-specific variable gene families that were also variable overall (yellow,
790 “both tests”) or that were specific to a particular phylum (red, “phylum-specific test only”).

8 Additional Files

792 Figure 1—figure supplement 1: **Schematic shows overview of data processing and method.** A)
793 Data is processed by taking reads from multiple datasets (represented by letters here) with a
794 certain number of samples (represented by S_A , S_B , etc.). These reads will eventually map to
795 multiple gene families G . MicrobeCensus [18] is used to estimate average genome size, while
796 Shotmap [17] is used to map reads, yielding both matrices of counts (right hand side) and ma-
797 trices of average lengths of the best-hit proteins (“average family length” or AFL). AFL and AGS
798 estimates are used to normalize counts. B) We calculate our statistic and assign p-values as fol-
799 lows. First, we normalize counts from Shotmap using AFL and AGS, log-transform the resulting
800 reads per kilobase of genome (RPKG), then apply a simple linear model to fit dataset- and gene-
801 family-specific effects. The resulting residuals (“residual log RPKG”) form a matrix of G genes by
802 $S_A+S_B+S_C$ samples. We take the variance across all samples for each gene to obtain a $1 \times G$ vector
803 of residual variances. To get a null distribution, we can either use data generated from a negative
804 binomial fit, or, optionally, from a negative binomial fit integrated with (shaded section) boot-
805 strap resampling. For the negative binomial fit, from the count matrices, we estimate the mean
806 of each gene in each dataset, as well as dataset-specific overdispersion parameters k . We then
807 use these to make simulated count datasets (“ $\times B$ ” indicating that this card is replicated once
808 for each of B simulations), which we process as in the case of the real data, yielding simulated
809 log-RPKG matrices and simulated residual variances for each gene family. For the resampling (if
810 applicable), we sample with replacement from each count dataset, yielding resampled counts.
811 We process these in the same way to obtain resampled residual variances. Finally, if using the
812 resampled data, we center and scale the resampled residual variances using per-gene-family
813 means and standard deviations from the simulated residual variances; otherwise, we simply
814 take the values from applying the test to the negative binomial simulations. These form the
815 background distribution (bottom panel, solid curve) for each gene in G (“ $\times G$ ” indicating that
816 this card is replicated once for each of G genes). The actual observed residual variance (dashed
817 line) is then compared to this distribution to obtain p-values (gray shaded area).

818 Figure 1—figure supplement 2: **Size parameter estimator choice affects accuracy of estima-**
819 **tion.** For each mock dataset y , simulated null data is generated from a negative binomial distri-
820 bution, fixing the size parameter k_y but allowing the mean $\mu_{g,y}$ to vary for each of 1,000 genes;
821 simulated true-positive gene families are drawn from a negative-binomial distribution with size
822 equal to zk_y or k_y/z , where z is the effect size. A-C) The choice of estimator affects the accu-
823 racy of size estimates. The mode method-of-moments estimator (C, y-axis) more accurately
824 estimates the true size specified in the simulation (x-axis) than the harmonic mean (A, y-axis)
825 or median (B, y-axis), and is more tolerant to differences in the ratio of true-positive variable
826 and invariable gene families (colors). D-E) When the size parameter is known, α (D) and power
827 (E) are well controlled, with α approximately equal to 0.05 at $p \leq 0.05$ and power approaching
828 1. Here, each simulation comprises three mock studies with different size parameters, mirror-
829 ing our actual data. Bar heights are means from 4 simulations and error bars are ± 2 SD. The
830 proportion of variable:invariable gene families was 0.5 and 44% of genes were true positives.

831 Figure 1—figure supplement 3: **Size parameter estimation affects power and α , with the**
832 **mode method-of-moments giving the best control.** α (A) was minimized and power (B) was
833 maximized when the mode method-of-moments estimator was used to get estimates of the
834 study-specific dispersion parameters \widehat{k}_y . Bars are from 4 simulations. The proportion of vari-
835 able:invariable gene families was 0.4 and 43% of genes were true positives.

836 Figure 1—figure supplement 4: **The mode estimator is robust to changes in the proportion of**
837 **true positives and the ratio of variable to invariable gene families.** α (A-C) and power (D-F)
838 as a function of the proportion of true positives (x-axis) and the ratio of variable to invariable
839 true positives (y-axis) for $n = 120$. $\alpha = 0.05$ and power = 1 are shown in color-bars to the left of
840 each heatmap for reference. α and power are calculated overall (left), for variable gene families
841 (center), and for invariable gene families (right). In general, α was better controlled for the in-
842 variable gene families than for the variable gene families; we therefore used different empirical
843 cutoffs for each set of genes.

844 Figure 2—figure supplement 1: **We identify significantly variable and invariable gene fami-**
845 **lies.** Density plots of distributions of residual variance (V_G) statistics for significantly invariable
846 (blue dashed line), non-significant (black solid line), and significantly variable (red dashed line)
847 gene families. The distributions had the expected trend (e.g., significantly variable gene families
848 tend to have higher residual variance) but also overlap, indicating the importance of the calcu-
849 lated null distribution. The inset shows the proportion of zero values for the non-significant
850 (black) and significantly invariable (blue) gene families with V_G falling in the lowest range (ver-
851 tical dashed lines), indicating that the test differentiates between gene families that only appear
852 invariable because they have few observations, and gene families that are consistently abun-
853 dant yet invariable.

854 Figure 2—figure supplement 2: **Heatmap showing significantly variable and invariable gene**
855 **families (unscaled).** Heatmap showing residual \log -RPKG abundances (i.e., after normalizing
856 for between-study effects and gene-specific abundances) of significantly invariable (blue) and
857 significantly variable (red) gene families. Variable and invariable gene families are clustered
858 separately, while samples are clustered within each dataset.

859 Figure 2—figure supplement 3: **Heatmap showing significantly variable and invariable gene**
860 **families (scaled).** As with 2—figure supplement 2, but residual \log -RPKG abundances scaled
861 by their expected variance under the negative binomial null model (see Methods).

862 **Figure 3—figure supplement 1: Carbon metabolism contains variable and invariable gene**
863 **families.** A) Pathway schematic showing a selection of measured gene families in-
864 volved in central carbohydrate metabolism. Gene families are color-coded by whether
865 they were variable (red) or invariable (blue), with strength of color corresponding to the
866 FDR cutoff (color intensity). Genes involved in the Entner-Doudoroff pathway (*edd*),
867 pentose metabolism (*fae-hps*), hexose metabolism (*K01622*, *K16306*), and tricarboxylic
868 acid cycle intermediate metabolism (*frdCD*) were found to be variable across healthy
869 hosts. Abbreviated metabolites are glucose-6-phosphate (G6P), fructose-6-phosphate (F6P),
870 fructose-1,6-bisphosphate (FBP), glyceraldehyde-3-phosphate (GAP), dihydroxyacetone phos-
871 phate (DHAP), 6-phosphogluconolactone (6PGL), 6-phosphogluconate (6PG), 2-keto-3-deoxy-
872 phosphogluconate (KDPG), ribulose-5-phosphate (R5P), ribose-5-phosphate (R5P), pyru-
873 vate (*pyr*), hexulose-6-phosphate (*Hu6P*), formaldehyde (HCHO), 2-amino-3,7-dideoxy-D-
874 threo-hept-6-ulosonate (*ADTH*), and tetrahydromethanopterin (*H₄MPT*). B) Heatmaps show-
875 ing scaled residual *log*-RPKG for gene families (rows) involved in central carbohydrate
876 metabolism. Variable (red) and invariable (blue) gene families are clustered separately, as are
877 samples within a particular study (columns). *log*-RPKG values are scaled by the expected vari-
878 ance from the negative-binomial null distribution.

879 **Figure 5—figure supplement 1: Number of leaves is correlated with tree density, but tree den-**
880 **sity corrects for the overall rate of evolution.** The number of leaves (i.e., individual sequences)
881 is plotted vs. tree density on a log-log scatter plot, with each circle representing one gene fam-
882 ily. Two outliers with lower density than expected are plotted in colors: a putative transposase
883 (green) and a *Staphylococcus* leukotoxin (red). Both families have large numbers of sequences
884 from the same organism.

885 **Figure 6—figure supplement 1: Variable gene families are less-often correlated to measured**
886 **host characteristics.** Bar plots give the fraction of gene families with at least one bacterial or
887 archaeal representative in each category (significantly invariable, non-significant, and signifi-
888 cantly variable) that were significantly correlated to various sample characteristics, using par-
889 tial Kendall's τ to account for study effects and a permutation test to assess significance. These
890 sample characteristics are average genome size (AGS), the ratio of Bacteroidetes to Firmicutes
891 (B/F ratio), and a measure of α -diversity (Shannon index). (***: $p \leq 10^{-8}$ by chi-squared test
892 after Bonferroni correction; **: $p \leq 10^{-4}$.)

893 Figure 6—figure supplement 2: **Variable gene families are less often correlated to enterotype-**
894 **associated taxa, and more often correlated to the Proteobacterial clade *Enterobacteriaceae*.**
895 Bar plots give the fraction of gene families with at least one bacterial or archaeal representa-
896 tive in each category (significantly invariable, non-significant, and significantly variable) that
897 were significantly correlated to the predicted abundance of specific bacterial clades (the genera
898 *Bacteroides* and *Prevotella*, and the families *Ruminococcaceae* and *Enterobacteriaceae*). Sig-
899 nificance was assessed as in Figure 6—figure supplement 1. (***: $p \leq 10^{-8}$ by chi-squared test
900 after Bonferroni correction; **: $p \leq 10^{-4}$.)

901 Figure 6—figure supplement 3: **Genes only annotated in Proteobacteria or Euryarchaeota,**
902 **but not Actinobacteria or Firmicutes, are more likely to be variable.** Bar plots give the fraction
903 of gene families with at least one bacterial or archaeal representative in each category (signif-
904 icantly invariable, non-significant, and significantly variable) that were annotated *only* in the
905 phylum listed (x-axis). Significance was assessed as in Figure 6—figure supplement 1, using a
906 Holm correction for significance. p-values are color-coded by whether a phylum was enriched
907 (red), depleted (blue), or neither (gray) for variable gene families (Holm-corrected $p \leq 0.1$).

908 Figure 6—figure supplement 4: **Genes only annotated in Proteobacteria or Euryarchaeota,**
909 **but not Actinobacteria or Firmicutes, are more likely to be variable in a test that uniformly**
910 **samples from phylum-specific genes.** Bar plots are as per Figure 6—figure supplement 3, but
911 test results come from a test that sampled equal parts phylum-specific genes and genes anno-
912 tated in all four listed phyla, with phylum-specific genes themselves uniformly sampled across
913 phyla. Significance was assessed as in Figure 6—figure supplement 3. p-values are color-coded
914 by whether a phylum was enriched (red), depleted (blue), or neither (gray) for variable gene
915 families (Holm-corrected $p \leq 0.1$).

916 Figure 7—figure supplement 1: **Phyla show similar trends of overlap at a generous FDR cutoff.**
917 A-B) Venn diagrams showing the number of significantly variable (A) and invariable (B) gene
918 families across Proteobacteria, Bacteroidetes, and Firmicutes, $FDR \leq 25\%$. Compare to Figure
919 7A-B.

920 Figure 7—figure supplement 2: **Comparison between Bacteroidetes- and Firmicutes-specific**
921 **variable and invariable genes.** A) Bars indicate the fraction of phylum-specific variable gene
922 families that were also variable overall (red, “both tests”) or that were specific to a particular
923 phylum (yellow, “phylum-specific test only”). A) For the Bacteroidetes- (left) and Firmicutes-
924 (right) specific tests, the proportion of invariable (blue), non-significant (gray), and variable
925 (red) gene families, at an estimated 5% FDR (using cutoffs from overall test). Pathways with at
926 least 5 total gene families across both phyla are shown. B) Rectangular Venn diagrams showing
927 the proportion of *Bacteroides*-specific (left), shared (center, bright), and Firmicutes-specific
928 (right) invariable (blue) and variable (red) gene families for each of the pathways enumerated
929 in B.

930

Table 1: q -value cutoffs to reach a given empirical FDR, estimated from simulation.

empirical FDR	q value cutoff, variable	q value cutoff, invariable
5%	0.0238	0.108
10%	0.0669	0.180
25%	0.181	0.294

931

932 Supplementary file 1: Module and pathway enrichments for variable and invariable gene sets
933 (Fisher's exact test $q \leq 0.25$).

934 Supplementary file 2: Module and pathway enrichments for variable/high-PD and
935 invariable/low-PD gene sets (Fisher's exact test $q \leq 0.25$).

936 Supplementary file 3: Module and pathway enrichments for gene families with invariable abun-
937 dances in every phylum-specific test (Fisher's exact test, $q \leq 0.25$).

938 Supplementary file 4: Module and pathway enrichments for gene families variable in each
939 phylum-specific test (Fisher's exact test, $q \leq 0.25$).

940 Supplementary file 5: SRA IDs and characteristics (read length, average genome size from Mi-
941 crobeCensus) for samples used in this study.

Supplementary file 6: Predicted tree densities.

942 Supplementary file 7: Supplementary note on correlation of variable and invariable gene fami-
943 lies with taxonomic summary statistics

944 Figure 1—Source data 1: Matrix of read counts (after rarefaction) for every gene family in each
945 sample included in the present study.

946 Figure 1—Source data 2: Matrix of average family lengths for every gene family in each sample
947 included in the present study.

948 Figure 1—Source data 3: Log-RPKG abundances for every gene family mapped in the present
949 study.

950 Figure 1—Source data 4: Residual log-RPKG abundances (i.e., after fitting the linear model) for
951 every gene family mapped in the present study.

952 Figure 2—Source data 1: Counts of invariable, non-significant, and variable gene families per
953 pathway.

954 Figure 2—Source data 2: Counts of invariable, non-significant, and variable gene families for
955 ribosomes in each domain of life.

956 Figure 3—Source data 1: Residual log-RPKG scaled by the expected variance under the null
957 model (see Methods).

958 Figure 5—Source data 1: \log_{10} phylogenetic distribution (PD), \log_{10} residual variance statis-
959 tics (residvar), significance at 5% FDR (invariable coded as “dn”, variable coded as “up”, non-
960 significant coded as “ns”), presence in at least one bacterial/archaeal genome in KEGG, and
961 annotations for all measured gene families.

962 Figure 6—Source data 1: Counts of significant associations of invariable, non-significant, and
963 variable gene families with phylum-level abundances.

964 Figure 6—Source data 2: Counts of significant associations of invariable, non-significant, and
965 variable gene families with taxonomic summary statistics.

Figure 7—Source data 1: q -values for gene families in the overall test.

Figure 7—Source data 2: q -values for gene families in phylum-specific tests.

966 Figure 7—Source data 3: JSON-formatted lists of significantly (in)variable or non-significant
967 gene families at 5% (“strong”), 10% (“med”), and 25% FDR (“weak”); overall test.

968 Figure 7—Source data 4: JSON-formatted lists of significantly (in)variable or non-significant
969 gene families at 5% (“strong”), 10% (“med”), and 25% FDR (“weak”); phylum-specific tests.

970 **9 References**

971 **References**

- 972 [1] Slack E, Hapfelmeier S, Stecher B, Velykoredko Y, Stoel M, Lawson MAE, Geuking MB, Beut-
973 ler B, Tedder TF, Hardt WD, et al. “Innate and Adaptive Immunity Cooperate Flexibly to
974 Maintain Host-Microbiota Mutualism.” *Science (New York, NY)*, **325** (2009)(5940):617–20.
975 ISSN 1095-9203. doi:10.1126/science.1172747.
- 976 [2] Atarashi K, Tanoue T, Shima T, Imaoka A, Kuwahara T, Momose Y, Cheng G, Yamasaki S,
977 Saito T, Ohba Y, et al. “Induction of Colonic Regulatory T Cells by Indigenous Clostrid-
978 ium Species.” *Science (New York, NY)*, **331** (2011)(6015):337–41. ISSN 1095-9203. doi:
979 10.1126/science.1198469.
- 980 [3] Hapfelmeier S, Lawson MAE, Slack E, Kirundi JK, Stoel M, Heikenwalder M, Cahenzli J,
981 Velykoredko Y, Balmer ML, Endt K, et al. “Reversible Microbial Colonization of Germ-
982 Free Mice Reveals the Dynamics of IgA Immune Responses.” *Science (New York, NY)*, **328**
983 (2010)(5986):1705–9. ISSN 1095-9203. doi:10.1126/science.1188454.
- 984 [4] Sonnenburg JL, Xu J, Leip DD, Chen CH, Westover BP, Weatherford J, Buhler JD, and Gor-
985 don JL. “Glycan Foraging in Vivo by an Intestine-Adapted Bacterial Symbiont.” *Science*
986 (*New York, NY*), **307** (2005)(5717):1955–9. ISSN 1095-9203. doi:10.1126/science.1109051.
- 987 [5] Wikoff WR, Anfora AT, Liu J, Schultz PG, Lesley SA, Peters EC, and Siuzdak G.
988 “Metabolomics Analysis Reveals Large Effects of Gut Microflora on Mammalian Blood
989 Metabolites.” *Proceedings of the National Academy of Sciences of the United States of Amer-*
990 *ica*, **106** (2009)(10):3698–703. ISSN 1091-6490. doi:10.1073/pnas.0812874106.
- 991 [6] Wang J and Jia H. “Metagenome-Wide Association Studies: Fine-Mining the Micro-
992 biome.” *Nature Reviews Microbiology*, **14** (2016)(8):508–22. ISSN 1740-1534. doi:
993 10.1038/nrmicro.2016.83.
- 994 [7] Haiser HJ, Gootenberg DB, Chatman K, Sirasani G, Balskus EP, and Turnbaugh PJ. “Pre-
995 dicting and Manipulating Cardiac Drug Inactivation by the Human Gut Bacterium *Eg-*
996 *gerthella lenta*.” *Science (New York, NY)*, **341** (2013)(6143):295–8. ISSN 1095-9203. doi:
997 10.1126/science.1235872.
- 998 [8] Wallace BD, Wang H, Lane KT, Scott JE, Orans J, Koo JS, Venkatesh M, Jobin C, Yeh LA, Mani
999 S, et al. “Alleviating Cancer Drug Toxicity by Inhibiting a Bacterial Enzyme.” *Science (New*
1000 *York, NY)*, **330** (2010)(6005):831–5. ISSN 1095-9203. doi:10.1126/science.1191175.

- 1001 [9] De Filippo C, Cavalieri D, Di Paola M, Ramazzotti M, Poullet JB, Massart S, Collini S, Pierac-
1002 cini G, and Lionetti P. “Impact of Diet in Shaping Gut Microbiota Revealed by a Compara-
1003 tive Study in Children from Europe and Rural Africa.” *Proceedings of the National Academy*
1004 *of Sciences of the United States of America*, **107** (2010)(33):14691–6. ISSN 1091-6490. doi:
1005 10.1073/pnas.1005963107.
- 1006 [10] Young VB, Knox KA, and Schauer DB. “Cytolethal Distending Toxin Sequence and Activ-
1007 ity in the Enterohepatic Pathogen *Helicobacter hepaticus*.” *Infection and Immunity*, **68**
1008 (2000)(1):184–91. ISSN 0019-9567.
- 1009 [11] Mazmanian SK, Round JL, and Kasper DL. “A Microbial Symbiosis Factor Prevents In-
1010 testinal Inflammatory Disease.” *Nature*, **453** (2008)(7195):620–5. ISSN 1476-4687. doi:
1011 10.1038/nature07008.
- 1012 [12] Turnbaugh PJ, Hamady M, Yatsunenko T, Cantarel BL, Duncan A, Ley RE, Sogin ML, Jones
1013 WJ, Roe BA, Affourtit JP, et al. “A Core Gut Microbiome in Obese and Lean Twins.” *Nature*,
1014 **457** (2009)(7228):480–4. ISSN 1476-4687. doi:10.1038/nature07540.
- 1015 [13] Human Microbiome Project Consortium. “Structure, Function and Diversity of the
1016 Healthy Human Microbiome.” *Nature*, **486** (2012)(7402):207–14. ISSN 1476-4687. doi:
1017 10.1038/nature11234.
- 1018 [14] Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N, Levenez
1019 F, Yamada T, et al. “a Human Gut Microbial Gene Catalogue Established by Metagenomic
1020 Sequencing.” *Nature*, **464** (2010)(7285):59–65. ISSN 1476-4687. doi:10.1038/nature08821.
- 1021 [15] Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F, Liang S, Zhang W, Guan Y, Shen D, et al. “a
1022 Metagenome-Wide Association Study of Gut Microbiota in Type 2 Diabetes.” *Nature*, **490**
1023 (2012)(7418):55–60. ISSN 1476-4687. doi:10.1038/nature11450.
- 1024 [16] Karlsson FH, Tremaroli V, Nookaew I, Bergström G, Behre CJ, Fagerberg B, Nielsen J,
1025 and Bäckhed F. “Gut Metagenome in European Women with Normal, Impaired and
1026 Diabetic Glucose Control.” *Nature*, **498** (2013)(7452):99–103. ISSN 1476-4687. doi:
1027 10.1038/nature12198.
- 1028 [17] Nayfach S, Bradley PH, Wyman SK, Laurent TJ, Williams A, Eisen JA, Pollard KS, and Sharp-
1029 ton TJ. “Automated and Accurate Estimation of Gene Family Abundance from Shotgun
1030 Metagenomes.” *PLoS Computational Biology*, **11** (2015)(11):e1004573. ISSN 1553-7358.
1031 doi:10.1371/journal.pcbi.1004573.

- 1032 [18] Nayfach S and Pollard KS. “Average Genome Size Estimation Improves Comparative
1033 Metagenomics and Sheds Light on the Functional Ecology of the Human Microbiome.”
1034 *Genome Biology*, **16** (2015):51. ISSN 1474-760X. doi:10.1186/s13059-015-0611-7.
- 1035 [19] Sauerwald A, Zhu W, Major TA, Roy H, Palioura S, Jahn D, Whitman WB, Yates JR, Ibba M,
1036 and Söll D. “RNA-Dependent Cysteine Biosynthesis in Archaea.” *Science (New York, NY)*,
1037 **307** (2005)(5717):1969–72. ISSN 1095-9203. doi:10.1126/science.1108329.
- 1038 [20] Dridi B, Henry M, El Khéchine A, Raoult D, and Drancourt M. “High Prevalence of
1039 *Methanobrevibacter smithii* and *Methanosphaera stadtmanae* Detected in the Human Gut
1040 Using an Improved DNA Detection Protocol.” *PLoS ONE*, **4** (2009)(9):e7063. ISSN 1932-
1041 6203. doi:10.1371/journal.pone.0007063.
- 1042 [21] Yanagisawa T, Sumida T, Ishii R, Takemoto C, and Yokoyama S. “a Paralog of Lysyl-TRNA
1043 Synthetase Aminoacylates a Conserved Lysine Residue in Translation Elongation Factor
1044 P.” *Nature Structural & Molecular Biology*, **17** (2010)(9):1136–43. ISSN 1545-9985. doi:
1045 10.1038/nsmb.1889.
- 1046 [22] Roy H, Zou SB, Bullwinkle TJ, Wolfe BS, Gilreath MS, Forsyth CJ, Navarre WW, and Ibba
1047 M. “the tRNA Synthetase Paralog PoxA Modifies Elongation Factor-P with (R)- β -Lysine.”
1048 *Nature chemical biology*, **7** (2011)(10):667–9. ISSN 1552-4469. doi:10.1038/nchembio.632.
- 1049 [23] Peekhaus N and Conway T. “What’s for Dinner?: Entner-Doudoroff Metabolism in Es-
1050 cherichia Coli.” *Journal of Bacteriology*, **180** (1998)(14):3495–502. ISSN 0021-9193.
- 1051 [24] Goenrich M, Thauer RK, Yurimoto H, and Kato N. “Formaldehyde Activating Enzyme (Fae)
1052 and Hexulose-6-Phosphate Synthase (Hps) in *Methanosarcina barkeri*: A Possible Func-
1053 tion in Ribose-5-Phosphate Biosynthesis.” *Archives of microbiology*, **184** (2005)(1):41–8.
1054 ISSN 0302-8933. doi:10.1007/s00203-005-0008-1.
- 1055 [25] Spencer ME and Guest JR. “Isolation and Properties of Fumarate Reductase Mutants of
1056 *Escherichia coli*.” *Journal of Bacteriology*, **114** (1973)(2):563–70. ISSN 0021-9193.
- 1057 [26] Coburn B, Sekirov I, and Finlay BB. “Type III Secretion Systems and Disease.” *Clinical
1058 Microbiology Reviews*, **20** (2007)(4):535–49. ISSN 0893-8512. doi:10.1128/CMR.00013-07.
- 1059 [27] Coulthurst SJ. “The Type VI Secretion System — a Widespread and Versatile Cell Tar-
1060 geting System”. *Research in Microbiology*, **164** (2013)(6):640–654. ISSN 09232508. doi:
1061 10.1016/j.resmic.2013.03.017.

- 1062 [28] Chatzidaki-Livanis M, Geva-Zatorsky N, and Comstock LE. “*Bacteroides fragilis* Type VI
1063 Secretion Systems Use Novel Effector and Immunity Proteins to Antagonize Human Gut
1064 Bacteroidales Species.” *Proceedings of the National Academy of Sciences of the United States
1065 of America*, **113** (2016)(13):3627–32. ISSN 1091-6490. doi:10.1073/pnas.1522510113.
- 1066 [29] Wexler AG, Bao Y, Whitney JC, Bobay LM, Xavier JB, Schofield WB, Barry NA, Russell
1067 AB, Tran BQ, Goo YA, et al. “Human Symbionts Inject and Neutralize Antibacterial
1068 Toxins to Persist in the Gut”. *Proceedings of the National Academy of Sciences*, **113**
1069 (2016)(13):201525637. ISSN 0027-8424. doi:10.1073/pnas.1525637113.
- 1070 [30] Cao TB and Saier MH. “The General Protein Secretory Pathway: Phylogenetic Analyses
1071 Leading to Evolutionary Conclusions.” *Biochimica et Biophysica Acta*, **1609** (2003)(1):115–
1072 25. ISSN 0006-3002.
- 1073 [31] Schromm AB, Brandenburg K, Loppnow H, Moran AP, Koch MHJ, Rietschel ET, and Sey-
1074 del U. “Biological Activities of Lipopolysaccharides Are Determined by the Shape of
1075 Their Lipid A Portion”. *European Journal of Biochemistry*, **267** (2000)(7):2008–2013. ISSN
1076 00142956. doi:10.1046/j.1432-1327.2000.01204.x.
- 1077 [32] Coats SR, Berezow AB, To TT, Jain S, Bainbridge BW, Banani KP, and Darveau RP. “The
1078 Lipid A Phosphate Position Determines Differential Host Toll-like Receptor 4 Responses to
1079 Phylogenetically Related Symbiotic and Pathogenic Bacteria”. *Infection and Immunity*, **79**
1080 (2011)(1):203–210. ISSN 0019-9567. doi:10.1128/IAI.00937-10.
- 1081 [33] Geurtsen J, Steeghs L, Hamstra HJ, Ten Hove J, de Haan A, Kuipers B, Tommassen J, and
1082 van der Ley P. “Expression of the Lipopolysaccharide-Modifying Enzymes PagP and PagL
1083 Modulates the Endotoxic Activity of *Bordetella pertussis*.” *Infection and Immunity*, **74**
1084 (2006)(10):5574–85. ISSN 0019-9567. doi:10.1128/IAI.00834-06.
- 1085 [34] Vatanen T, Kostic A, D’Hennezel E, Siljander H, Franzosa E, Yassour M, Kolde R, Vlamakis
1086 H, Arthur T, Hämmäläinen AM, et al. “Variation in Microbiome LPS Immunogenicity Con-
1087 tributes to Autoimmunity in Humans”. *Cell*, **165** (2016)(4):842–853. ISSN 00928674. doi:
1088 10.1016/j.cell.2016.04.007.
- 1089 [35] Gardiner KR, Halliday MI, Barclay GR, Milne L, Brown D, Stephens S, Maxwell RJ, and Row-
1090 lands BJ. “Significance of Systemic Endotoxaemia in Inflammatory Bowel Disease.” *Gut*,
1091 **36** (1995)(6):897–901. ISSN 0017-5749.

- 1092 [36] Boutagy NE, McMillan RP, Frisard MI, and Hulver MW. “Metabolic Endotoxemia with Obe-
1093 sity: Is It Real and Is It Relevant?” *Biochimie*, **124** (2016):11–20. ISSN 03009084. doi:
1094 10.1016/j.biochi.2015.06.020.
- 1095 [37] Kallio KAE, Hätönen KA, Lehto M, Salomaa V, Männistö S, and Pussinen PJ. “Endotox-
1096 emia, Nutrition, and Cardiometabolic Disorders.” *Acta Diabetologica*, **52** (2015)(2):395–
1097 404. ISSN 1432-5233. doi:10.1007/s00592-014-0662-3.
- 1098 [38] Needham BD, Carroll SM, Giles DK, Georgiou G, Whiteley M, and Trent MS. “Modulat-
1099 ing the Innate Immune Response by Combinatorial Engineering of Endotoxin”. *Proceed-*
1100 *ings of the National Academy of Sciences*, **110** (2013)(4):1464–1469. ISSN 0027-8424. doi:
1101 10.1073/pnas.1218080110.
- 1102 [39] Wu D. “Personal communication”.
- 1103 [40] Hooper LV, Midtvedt T, and Gordon JI. “How Host-Microbial Interactions Shape the Nutri-
1104 ent Environment of the Mammalian Intestine.” *Annual Review of Nutrition*, **22** (2002):283–
1105 307. ISSN 0199-9885.
- 1106 [41] Benjdia A, Martens EC, Gordon JI, and Berteau O. “Sulfatases and a Radical S-Adenosyl-L-
1107 Methionine (AdoMet) Enzyme Are Key for Mucosal Foraging and Fitness of the Prominent
1108 Human Gut Symbiont, *Bacteroides thetaiotaomicron*.” *The Journal of Biological Chemistry*,
1109 **286** (2011)(29):25973–82. ISSN 1083-351X. doi:10.1074/jbc.M111.228841.
- 1110 [42] Ulmer JE, Vilén EM, Namburi RB, Benjdia A, Beneteau J, Malleron A, Bonnaffé D,
1111 Driguez PA, Descroix K, Lassalle G, et al. “Characterization of Glycosaminoglycan (GAG)
1112 Sulfatases from the Human Gut Symbiont *Bacteroides Thetaiotaomicron* Reveals the
1113 First GAG-Specific Bacterial Endosulfatase.” *The Journal of Biological Chemistry*, **289**
1114 (2014)(35):24289–303. ISSN 1083-351X. doi:10.1074/jbc.M114.573303.
- 1115 [43] Raghavan V and Groisman EA. “Species-Specific Dynamic Responses of Gut Bacteria to
1116 a Mammalian Glycan.” *Journal of Bacteriology*, **197** (2015)(9):1538–48. ISSN 1098-5530.
1117 doi:10.1128/JB.00010-15.
- 1118 [44] Greenblum S, Carr R, and Borenstein E. “Extensive Strain-Level Copy-Number Variation
1119 Across Human Gut Microbiome Species.” *Cell*, **160** (2015)(4):583–94. ISSN 1097-4172. doi:
1120 10.1016/j.cell.2014.12.038.
- 1121 [45] Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower
1122 C, and Segata N. “MetaPhlAn2 for Enhanced Metagenomic Taxonomic Profiling”. *Nature*
1123 *Methods*, **12** (2015)(10):902–903. ISSN 1548-7091. doi:10.1038/nmeth.3589.

- 1124 [46] Shin NR, Whon TW, and Bae JW. “Proteobacteria: Microbial Signature of Dysbiosis in
1125 Gut Microbiota.” *Trends in biotechnology*, **33** (2015)(9):496–503. ISSN 1879-3096. doi:
1126 10.1016/j.tibtech.2015.06.011.
- 1127 [47] Mukhopadhyaya I, Hansen R, El-Omar EM, and Hold GL. “IBD-What Role Do Proteobacteria
1128 Play?” *Nature Reviews Gastroenterology & hepatology*, **9** (2012)(4):219–30. ISSN 1759-5053.
1129 doi:10.1038/nrgastro.2012.14.
- 1130 [48] Garrett WS, Gallini CA, Yatsunencko T, Michaud M, DuBois A, Delaney ML, Punit S, Karls-
1131 son M, Bry L, Glickman JN, et al. “Enterobacteriaceae Act in Concert with the Gut Micro-
1132 biota to Induce Spontaneous and Maternally Transmitted Colitis.” *Cell Host & Microbe*, **8**
1133 (2010)(3):292–300. ISSN 1934-6069. doi:10.1016/j.chom.2010.08.004.
- 1134 [49] Carvalho FA, Koren O, Goodrich JK, Johansson MEV, Nalbantoglu I, Aitken JD, Su Y, Chas-
1135 saing B, Walters WA, González A, et al. “Transient Inability to Manage Proteobacteria
1136 Promotes Chronic Gut Inflammation in TLR5-Deficient Mice.” *Cell Host & Microbe*, **12**
1137 (2012)(2):139–52. ISSN 1934-6069. doi:10.1016/j.chom.2012.07.004.
- 1138 [50] Ley RE, Bäckhed F, Turnbaugh P, Lozupone CA, Knight RD, and Gordon JI. “Obesity Alters
1139 Gut Microbial Ecology.” *Proceedings of the National Academy of Sciences of the United*
1140 *States of America*, **102** (2005)(31):11070–5. ISSN 0027-8424. doi:10.1073/pnas.0504978102.
- 1141 [51] Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K, Kurilshikov A, Bonder
1142 MJ, Valles-Colomer M, Vandeputte D, et al. “Population-Level Analysis of Gut Micro-
1143 biome Variation.” *Science (New York, NY)*, **352** (2016)(6285):560–4. ISSN 1095-9203. doi:
1144 10.1126/science.aad3503.
- 1145 [52] Nayfach S and Pollard KS. “Population Genetic Analyses of Metagenomes Reveal Extensive
1146 Strain-Level Variation in Prevalent Human-Associated Bacteria”. Tech. rep., 2015. doi:
1147 10.1101/031757.
- 1148 [53] Ma L, Terwilliger A, and Maresso AW. “Iron and Zinc Exploitation During Bacterial Patho-
1149 genesis.” *Metallomics: Integrated Biometal Science*, **7** (2015)(12):1541–54. ISSN 1756-591X.
1150 doi:10.1039/c5mt00170f.
- 1151 [54] Wallden K, Rivera-Calzada A, and Waksman G. “Type IV Secretion Systems: Versatility and
1152 Diversity in Function.” *Cellular Microbiology*, **12** (2010)(9):1203–12. ISSN 1462-5822. doi:
1153 10.1111/j.1462-5822.2010.01499.x.

- 1154 [55] Lozupone CA, Stombaugh JI, Gordon JI, Jansson JK, and Knight R. “Diversity, Stability and
1155 Resilience of the Human Gut Microbiota.” *Nature*, **489** (2012)(7415):220–30. ISSN 1476-
1156 4687. doi:10.1038/nature11550.
- 1157 [56] Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, and Huttenhower C.
1158 “Metagenomic Biomarker Discovery and Explanation.” *Genome Biology*, **12** (2011)(6):R60.
1159 ISSN 1465-6914. doi:10.1186/gb-2011-12-6-r60.
- 1160 [57] Brennecke P, Anders S, Kim JK, Kołodziejczyk AA, Zhang X, Proserpio V, Baying B,
1161 Benes V, Teichmann SA, Marioni JC, et al. “Accounting for Technical Noise in Single-Cell
1162 RNA-Seq Experiments.” *Nature Methods*, **10** (2013)(11):1093–5. ISSN 1548-7105. doi:
1163 10.1038/nmeth.2645.
- 1164 [58] Vallejos CA, Marioni JC, and Richardson S. “BASiCS: Bayesian Analysis of Single-Cell Se-
1165 quencing Data”. *PLoS Computational Biology*, **11** (2015)(6):e1004333. ISSN 1553-7358.
1166 doi:10.1371/journal.pcbi.1004333.
- 1167 [59] Dumitrascu B, Darnell G, Ayroles J, and Engelhardt BE. “A Bayesian Test to Identify Vari-
1168 ance Effects”. Tech. rep.
- 1169 [60] Vallejos CA, Richardson S, and Marioni JC. “Beyond Comparisons of Means: Understand-
1170 ing Changes in Gene Expression at the Single-Cell Level.” *Genome Biology*, **17** (2016)(1):70.
1171 ISSN 1474-760X. doi:10.1186/s13059-016-0930-3.
- 1172 [61] Andrews S. “FastQC: A quality control tool for high throughput sequence data.”, 2010. doi:
1173 citeulike-article-id:11583827.
- 1174 [62] Kanehisa M, Goto S, Kawashima S, Okuno Y, and Hattori M. “the KEGG Resource for De-
1175 cipherring the Genome.” *Nucleic Acids Research*, **32** (2004)(Database issue):D277–80. ISSN
1176 1362-4962. doi:10.1093/nar/gkh063.
- 1177 [63] Yu D, Huber W, and Vitek O. “Shrinkage Estimation of Dispersion in Negative Binomial
1178 Models for RNA-Seq Experiments with Small Sample Size.” *Bioinformatics (Oxford, Eng-
1179 land)*, **29** (2013)(10):1275–82. ISSN 1367-4811. doi:10.1093/bioinformatics/btt143.
- 1180 [64] Storey JD and Tibshirani R. “Statistical Significance for Genomewide Studies.” *Proceedings
1181 of the National Academy of Sciences of the United States of America*, **100** (2003)(16):9440–5.
1182 ISSN 0027-8424. doi:10.1073/pnas.1530509100.
- 1183 [65] Storey JD, Bass AJ, Dabney A, and Robinson D. “qvalue: Q-Value Estimation for False Dis-
1184 covery Rate Control”, 2015.

- 1185 [66] Pollard K and van der Laan M. “Resampling-Based Multiple Testing: Asymptotic Control
1186 of Type I Error and Applications to Gene Expression Data”, 2003.
- 1187 [67] Pollard KS and van der Laan MJ. “Choice of a Null Distribution in Resampling-Based Mul-
1188 tiple Testing”. *Journal of Statistical Planning and Inference*, **125** (2004)(1-2):85–100. ISSN
1189 03783758. doi:10.1016/j.jspi.2003.07.019.
- 1190 [68] Collinson I, Corey RA, Allen WJ, Krogh A, Larsson B, von Heijne G, Sonnhammer E, Song C,
1191 Kumar A, Saleh M, et al. “Channel Crossing: How Are Proteins Shipped Across the Bacterial
1192 Plasma Membrane?” *Philosophical Transactions of the Royal Society of London Series B, Bi-*
1193 *ological sciences*, **370** (2015)(1679):567–580. ISSN 1471-2970. doi:10.1098/rstb.2015.0025.
- 1194 [69] Portaliou AG, Tsolis KC, Loos MS, Zorzini V, and Economou A. “Type III Secretion: Build-
1195 ing and Operating a Remarkable Nanomachine”. *Trends in Biochemical Sciences*, **41**
1196 (2016)(2):175–189. ISSN 09680004. doi:10.1016/j.tibs.2015.09.005.
- 1197 [70] Caspi R, Foerster H, Fulcher CA, Kaipa P, Krummenacker M, Latendresse M, Paley S, Rhee
1198 SY, Shearer AG, Tissier C, et al. “The MetaCyc Database of Metabolic Pathways and En-
1199 zymes and the BioCyc Collection of Pathway/Genome Databases.” *Nucleic Acids Res*, **36**
1200 (2008)(Database issue):D623—D631. doi:10.1093/nar/gkm900.
- 1201 [71] Wang X and Quinn PJ. “Lipopolysaccharide: Biosynthetic Pathway and Structure Mod-
1202 ification”. *Progress in Lipid Research*, **49** (2010)(2):97–107. ISSN 01637827. doi:
1203 10.1016/j.plipres.2009.06.002.
- 1204 [72] Whitfield C and Trent MS. “Biosynthesis and Export of Bacterial Lipopolysaccharides*”.
1205 <http://dxdoiorg/101146/annurev-biochem-060713-035600>, (2014).
- 1206 [73] Lloyd-Smith JO. “Maximum Likelihood Estimation of the Negative Binomial Dispersion
1207 Parameter for Highly Overdispersed Data, with Applications to Infectious Diseases.” *PLoS*
1208 *ONE*, **2** (2007)(2):e180. ISSN 1932-6203. doi:10.1371/journal.pone.0000180.
- 1209 [74] Sievers F and Higgins DG. “Clustal Omega, Accurate Alignment of Very Large Numbers of
1210 Sequences”. In “Methods in Molecular Biology (Clifton, N.J.)”, vol. 1079, pages 105–116.
1211 2014.
- 1212 [75] Price MN, Dehal PS, and Arkin AP. “FastTree 2—Approximately Maximum-Likelihood
1213 Trees for Large Alignments.” *PLoS ONE*, **5** (2010)(3):e9490. ISSN 1932-6203. doi:
1214 10.1371/journal.pone.0009490.
- 1215 [76] Kim S. “ppcor: Partial and Semi-Partial (Part) Correlation”, 2015.

- 1216 [77] Zhao Y, Tang H, and Ye Y. "RAPSearch2: A Fast and Memory-Efficient Protein Similarity
1217 Search Tool for Next-Generation Sequencing Data." *Bioinformatics (Oxford, England)*, **28**
1218 (2012)(1):125–6. ISSN 1367-4811. doi:10.1093/bioinformatics/btr595.
- 1219 [78] Wu D, Jospin G, and Eisen JA. "Systematic Identification of Gene Families for Use As
1220 "Markers" for Phylogenetic and Phylogeny-Driven Ecological Studies of Bacteria and Ar-
1221 chaea and Their Major Subgroups." *PLoS ONE*, **8** (2013)(10):e77033. ISSN 1932-6203. doi:
1222 10.1371/journal.pone.0077033.
- 1223 [79] Manor O and Borenstein E. "MUSiCC: A Marker Genes Based Framework for Metagenomic
1224 Normalization and Accurate Profiling of Gene Abundances in the Microbiome". *Genome*
1225 *Biology*, **16** (2015)(1):53. ISSN 1465-6906. doi:10.1186/s13059-015-0610-8.

Figure 1

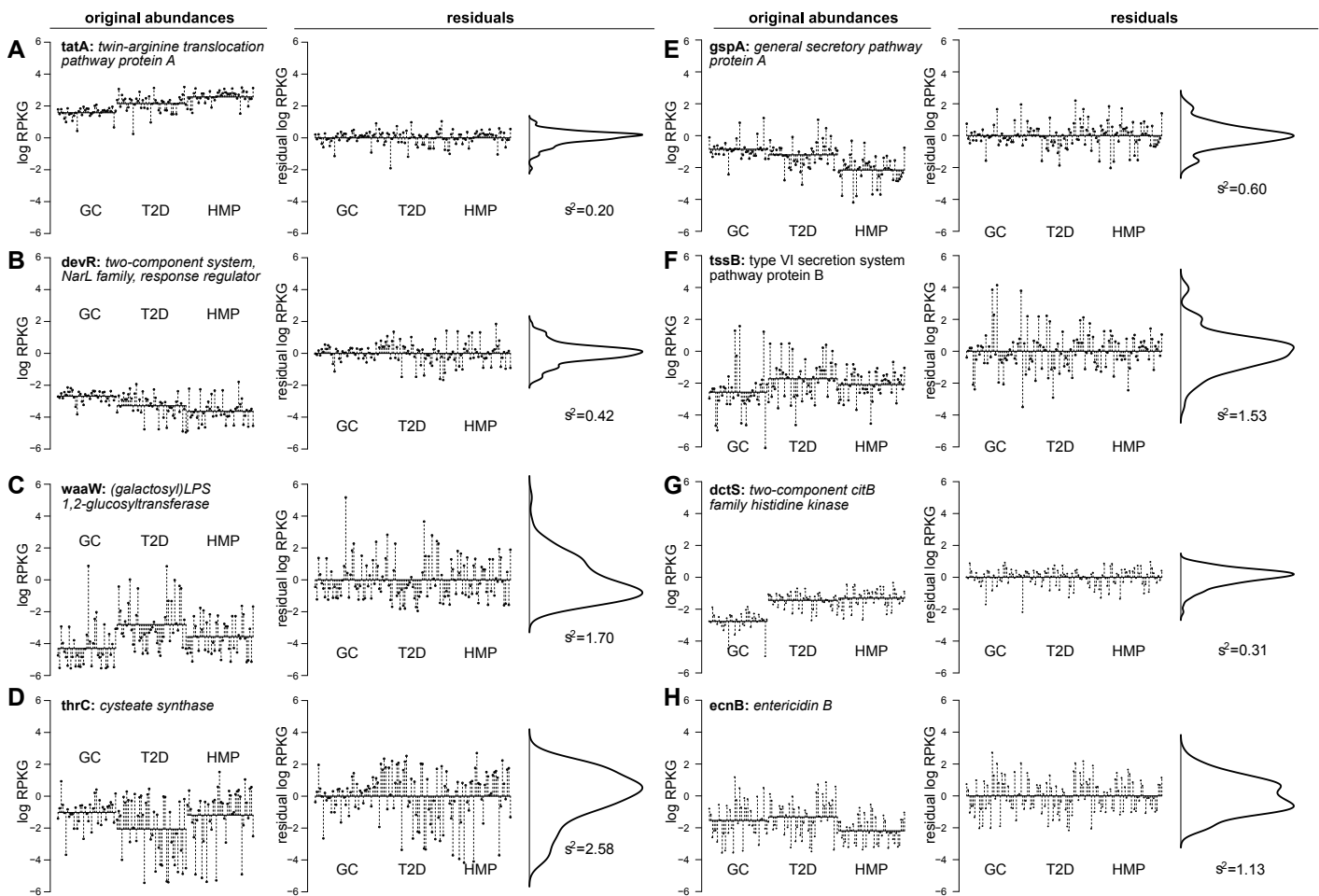


Figure 2

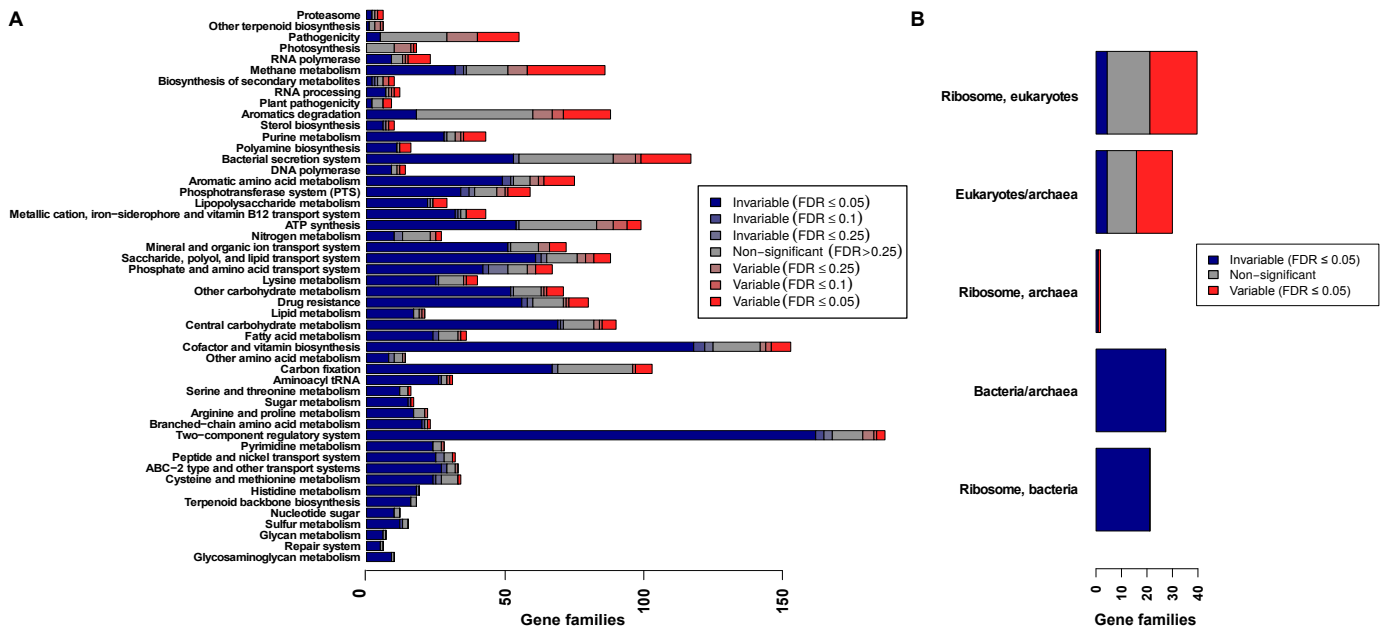


Figure 3

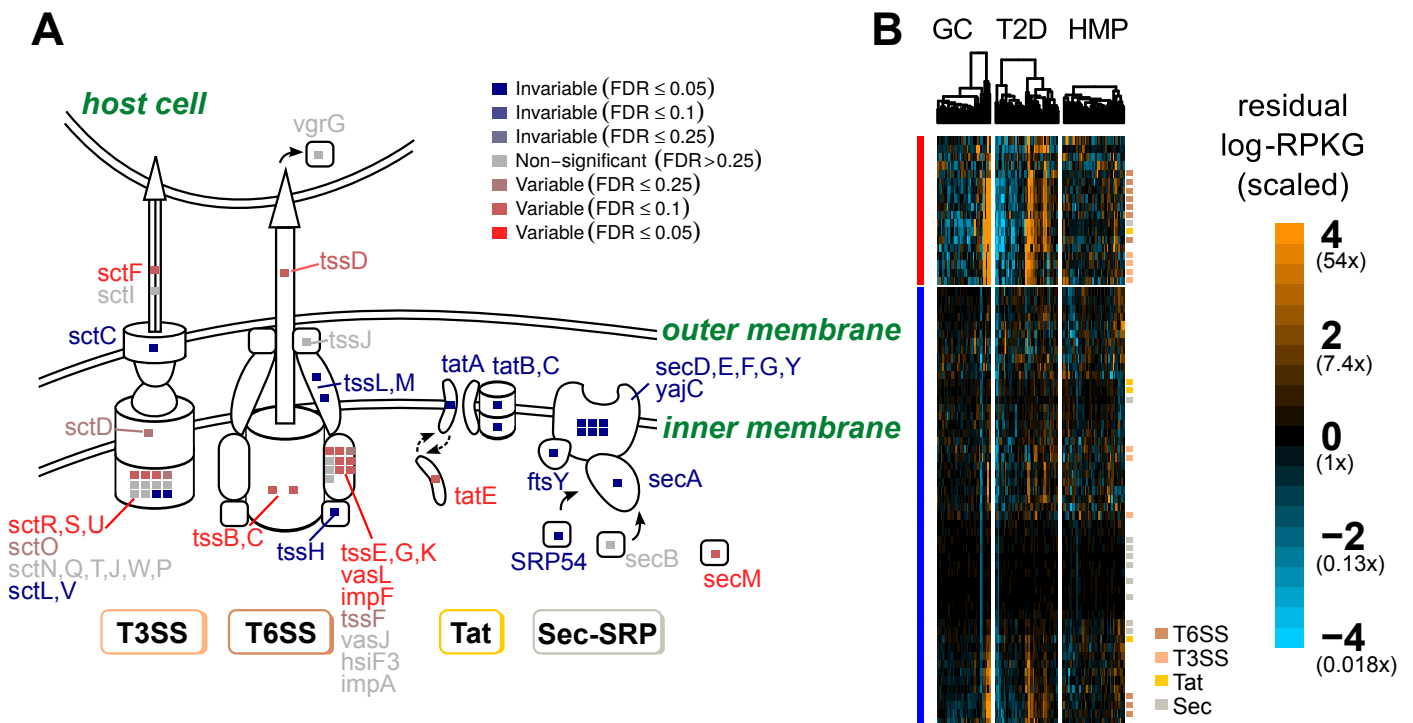


Figure 4

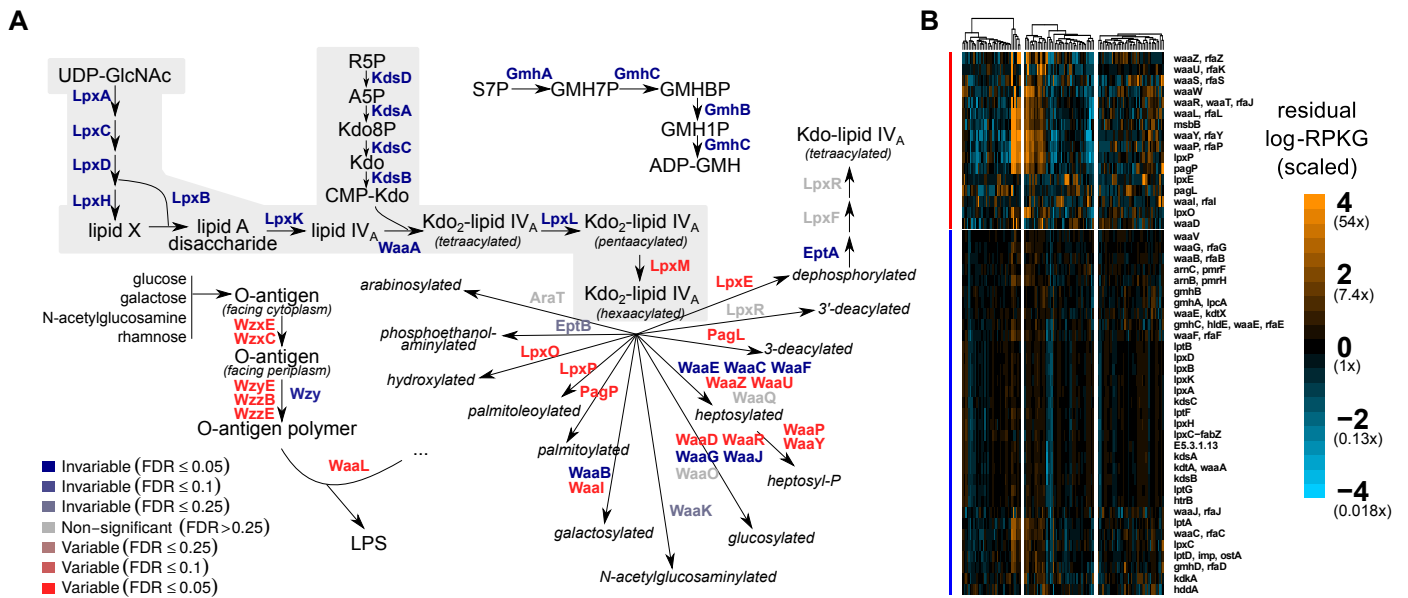


Figure 5

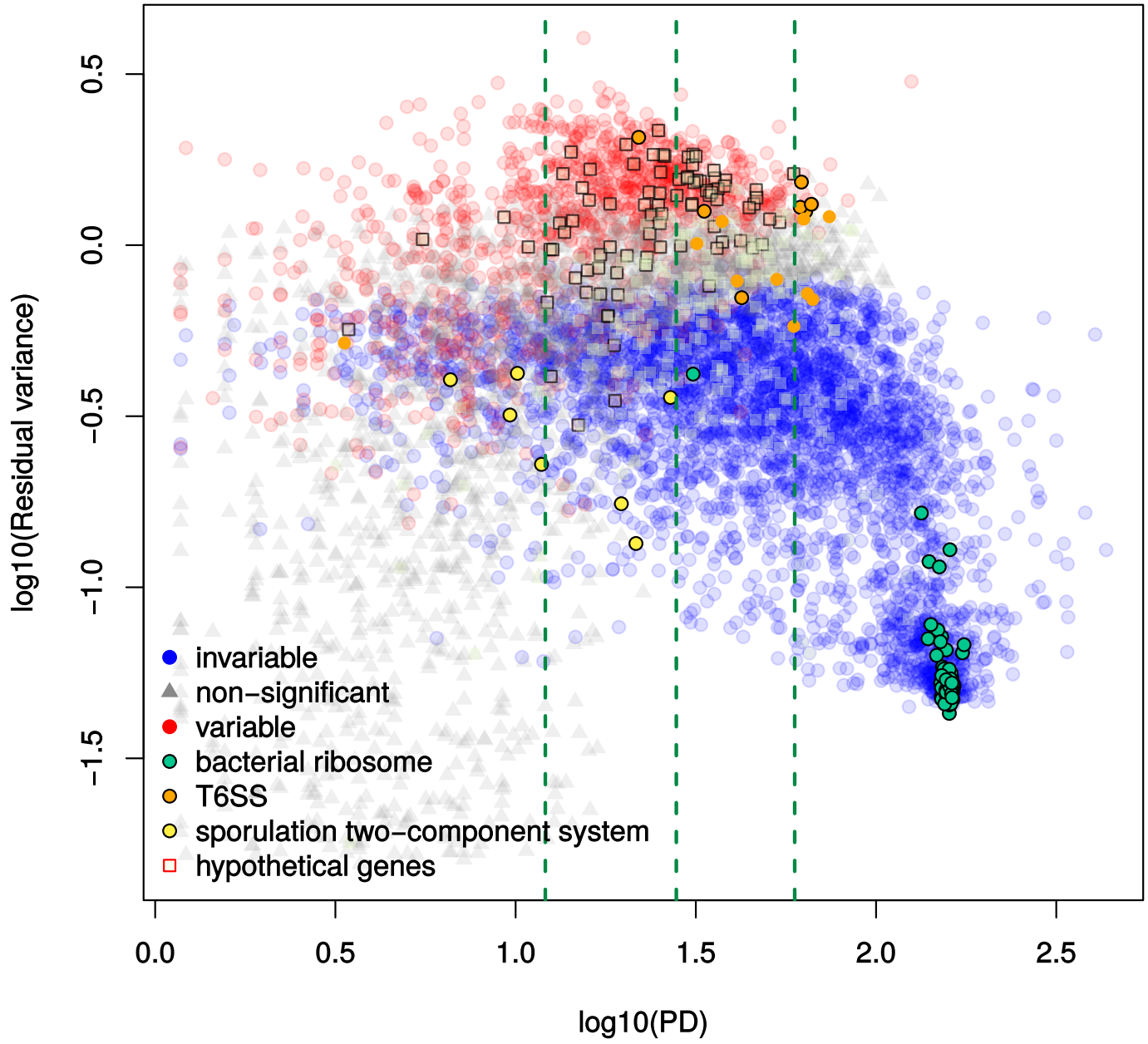


Figure 6

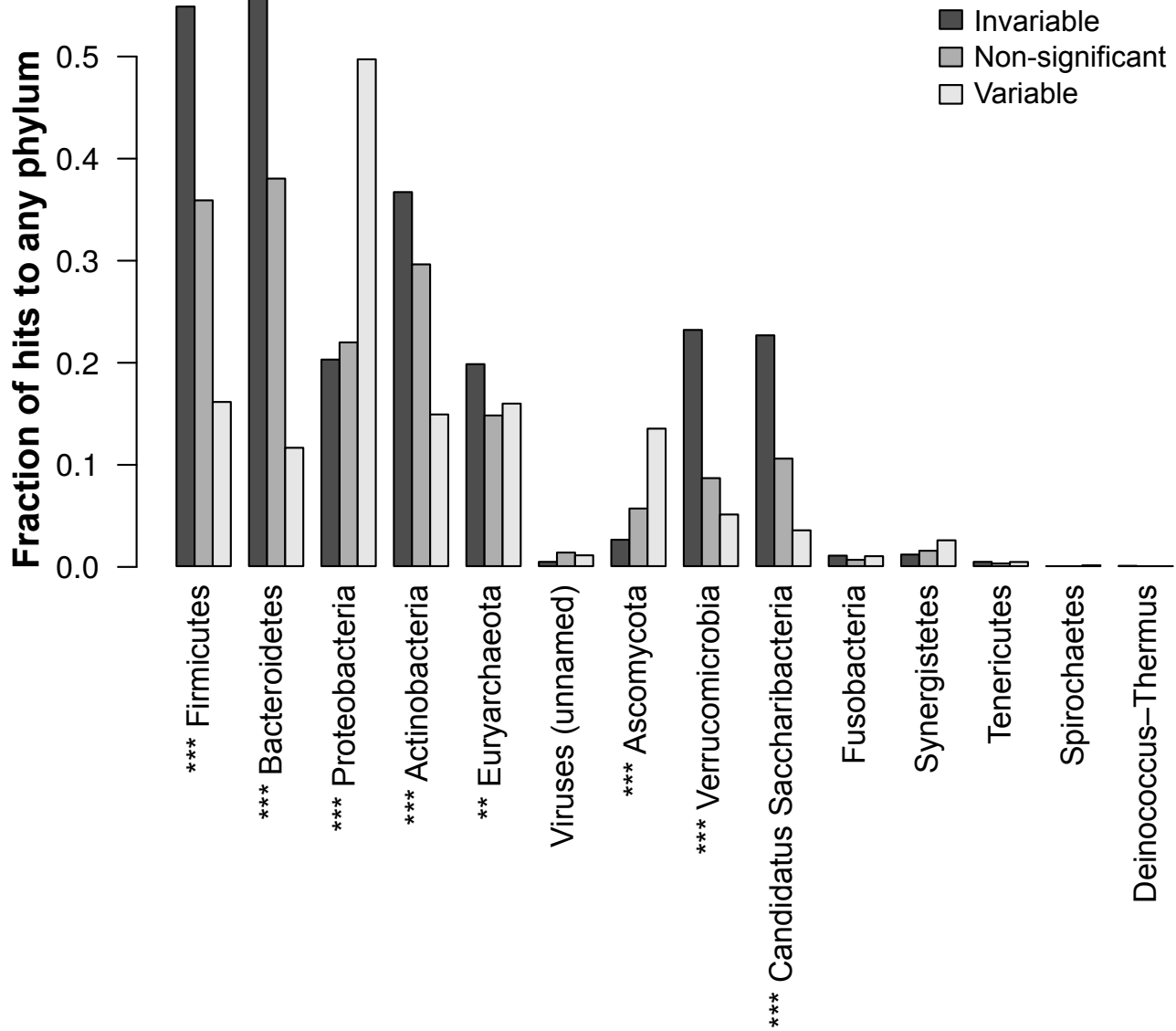


Figure 7

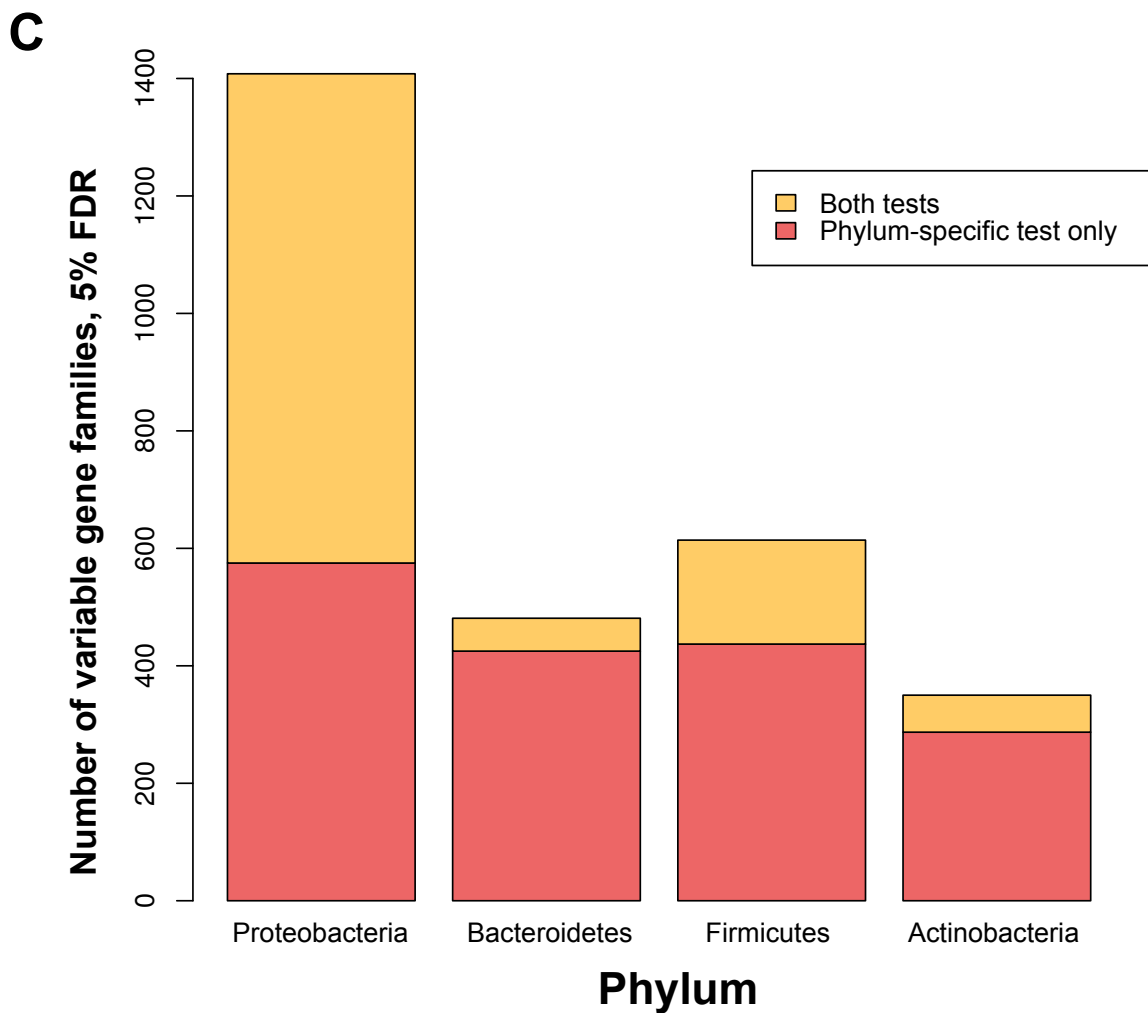
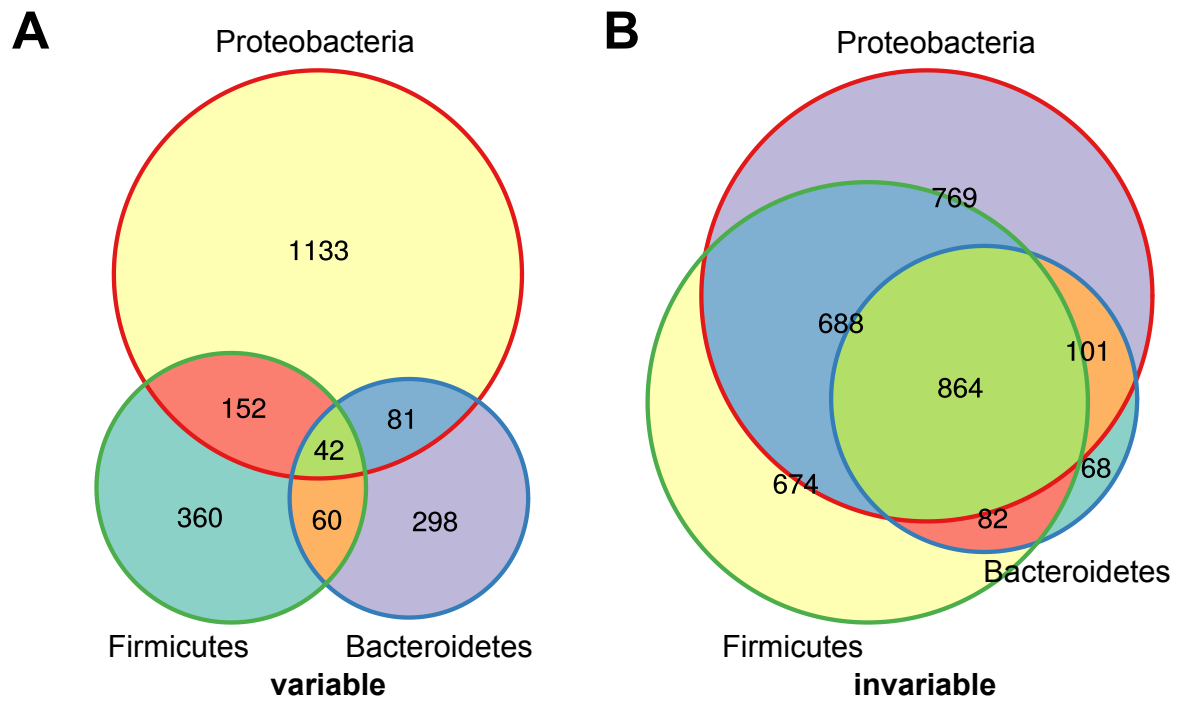
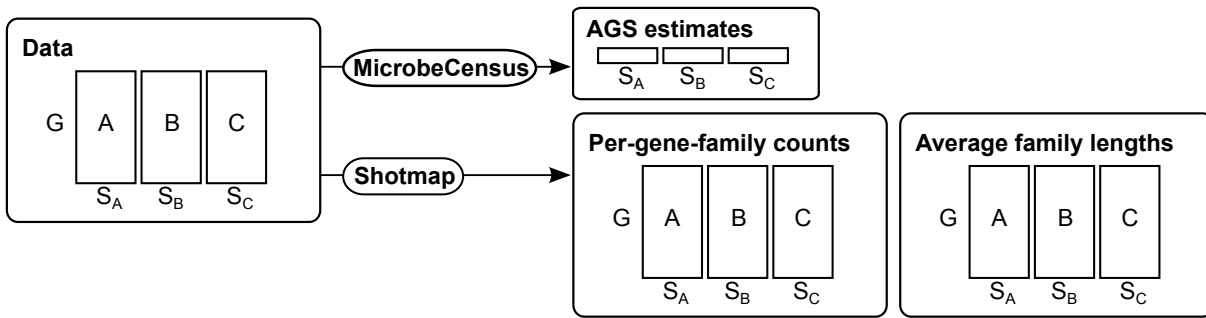


Figure 1—figure supplement 1

A



B

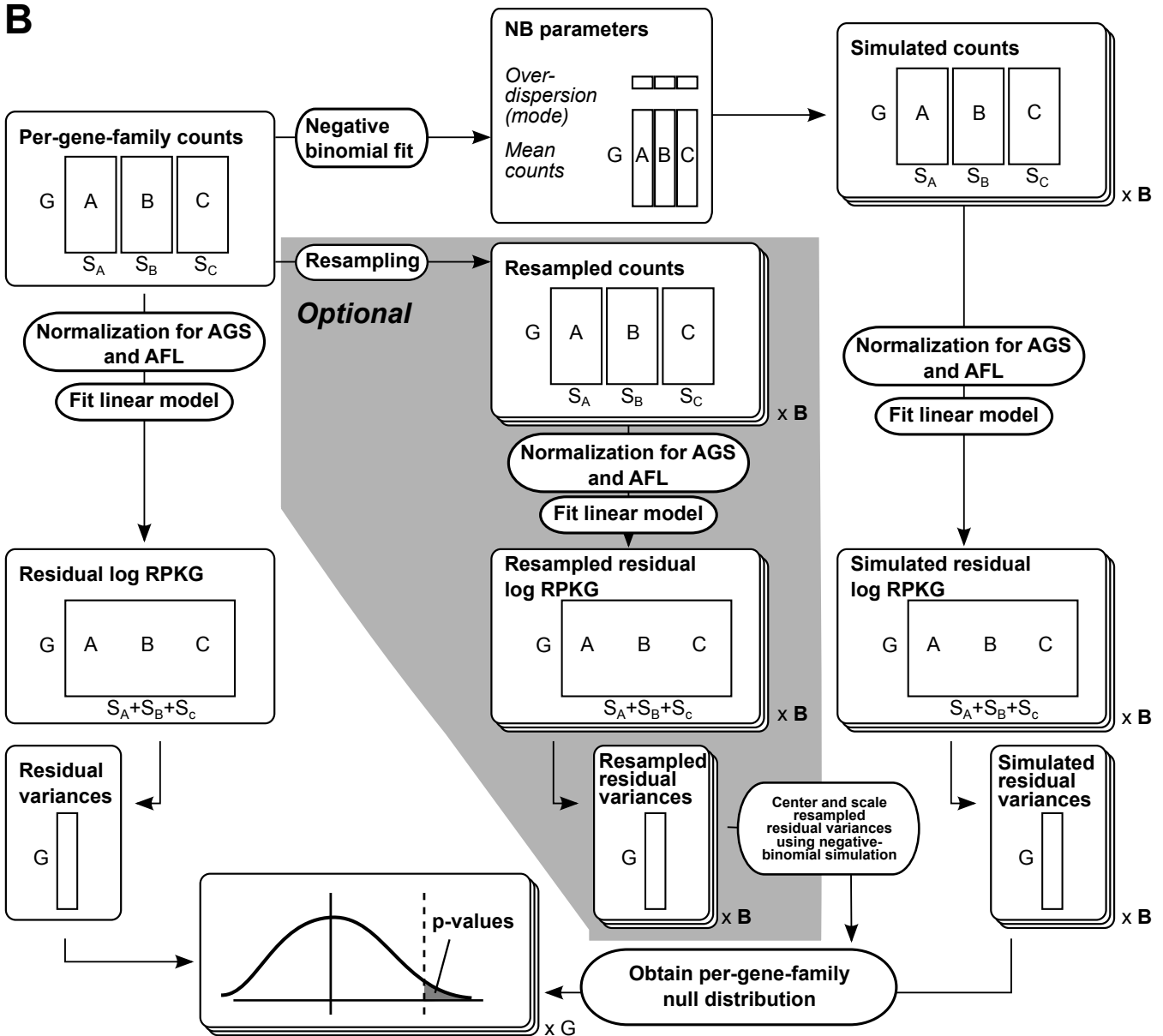


Figure 1—figure supplement 2

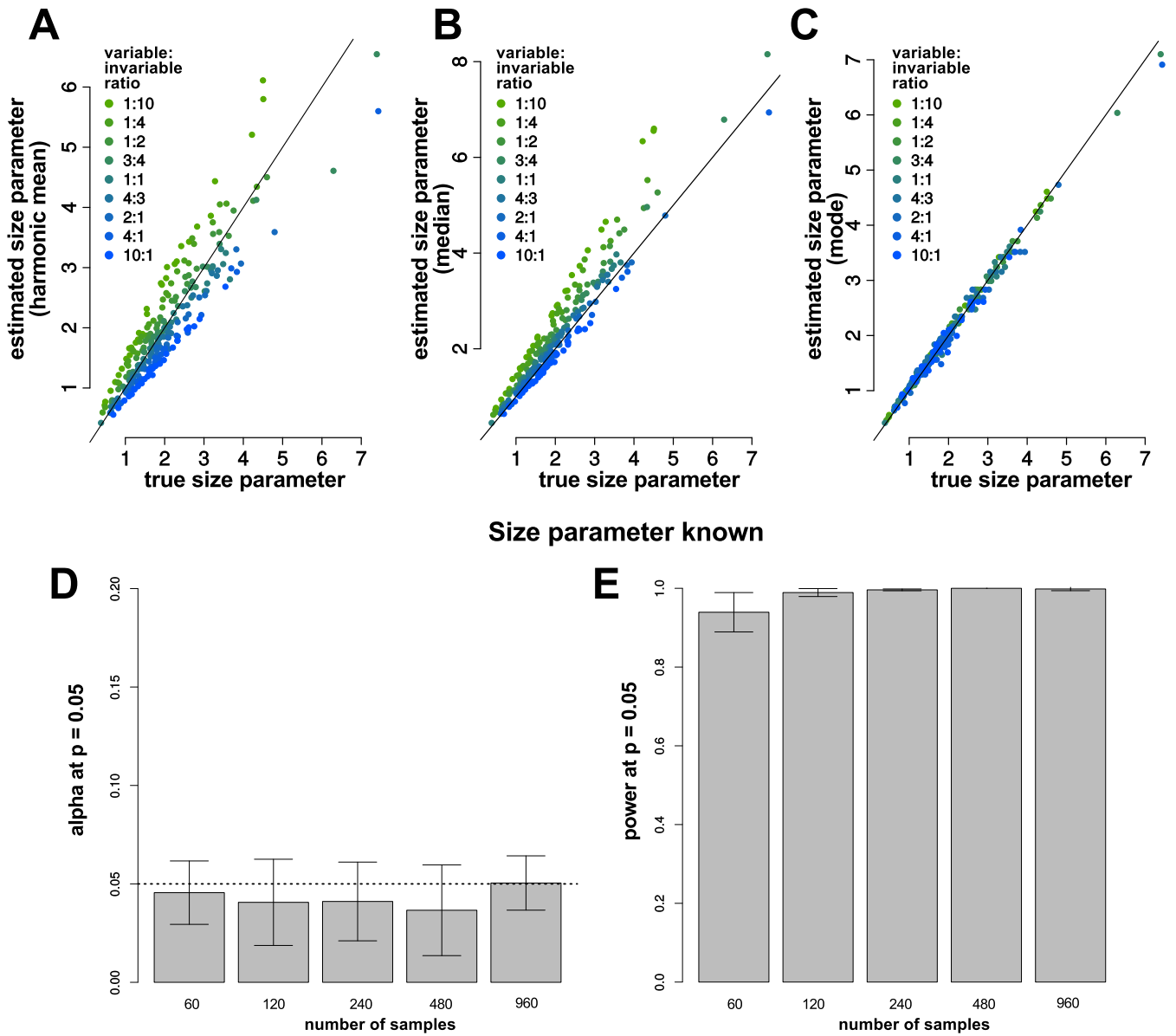


Figure 1—figure supplement 3

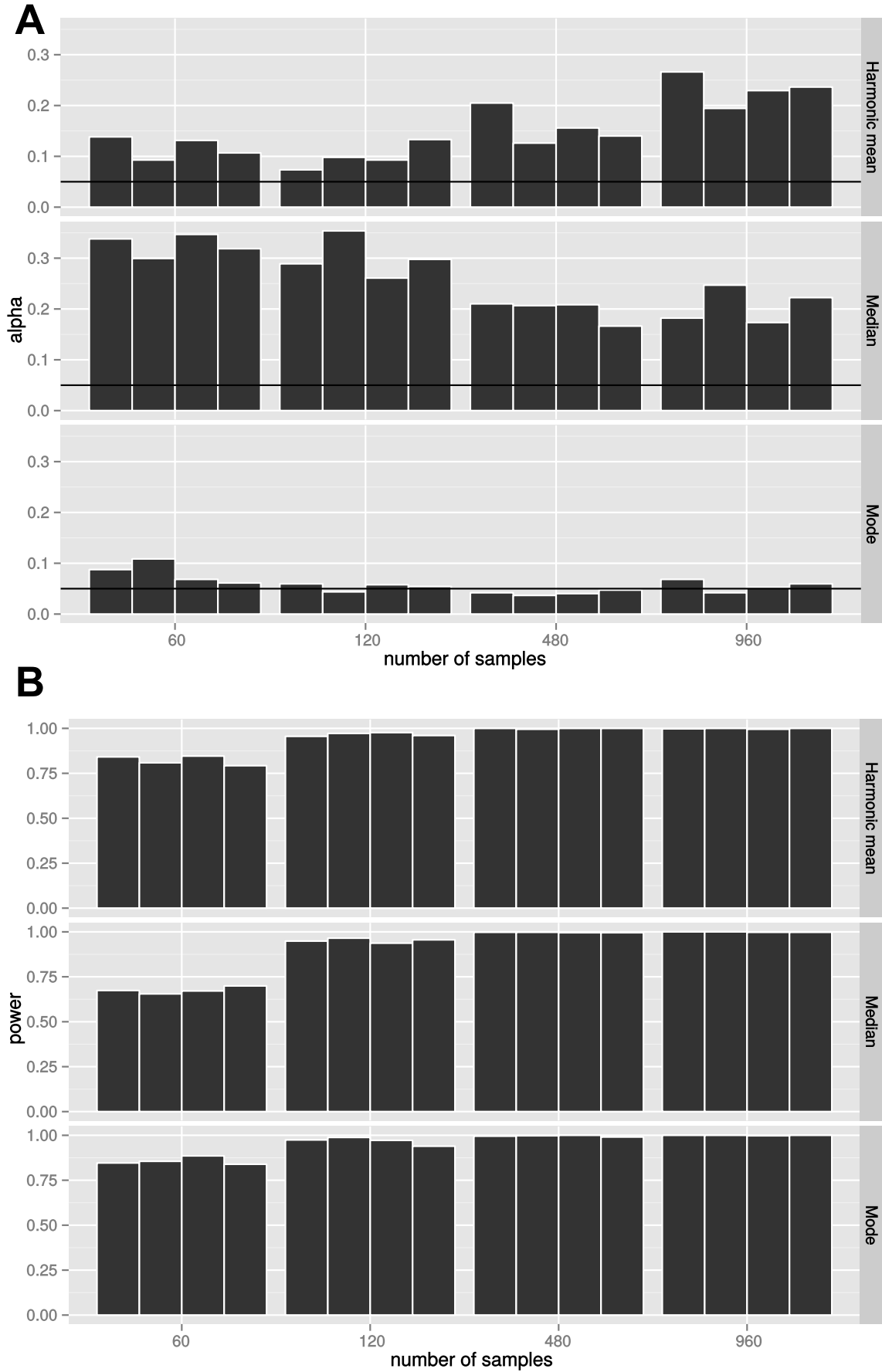


Figure 1—figure supplement 4

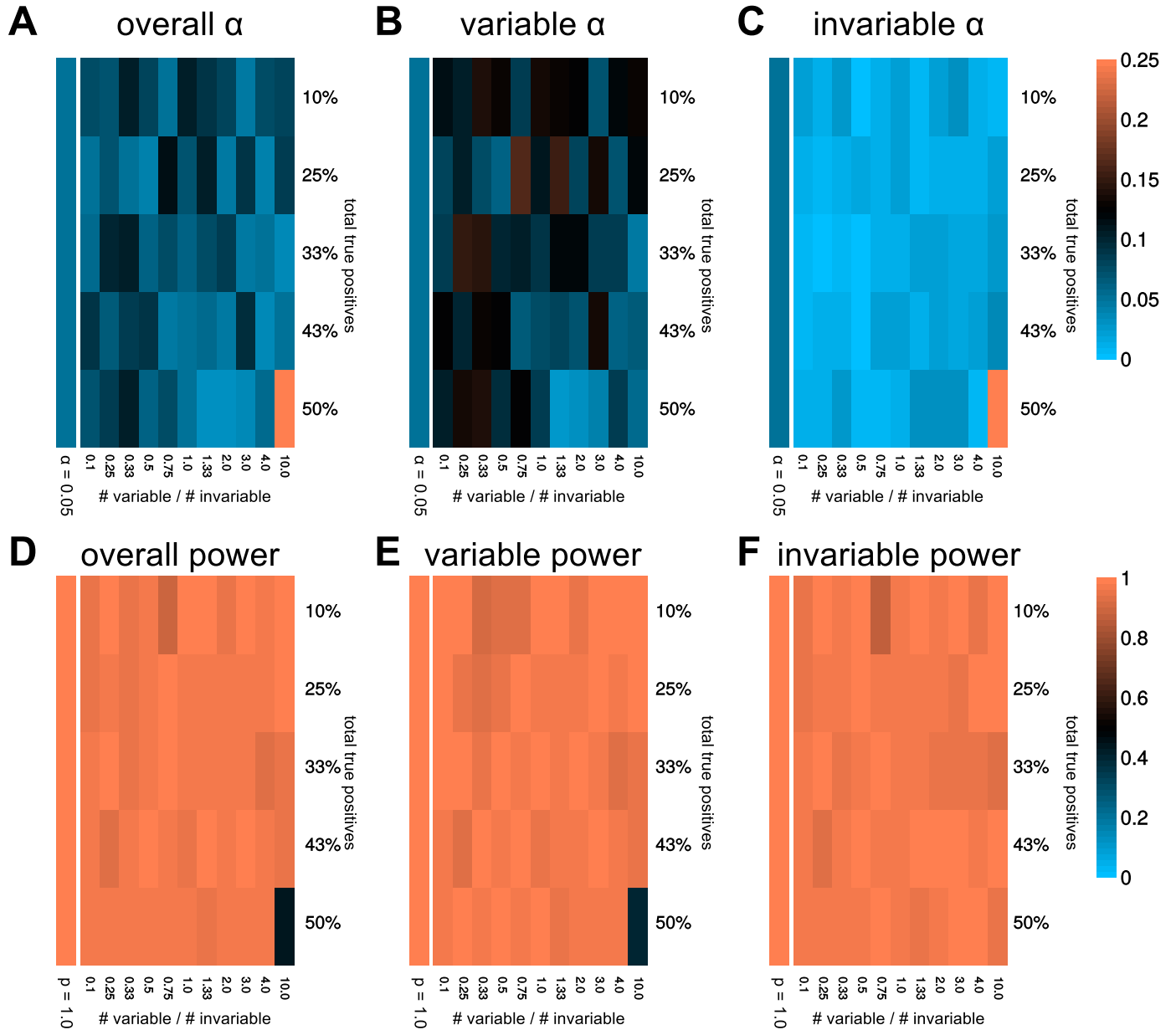


Figure 2—figure supplement 1

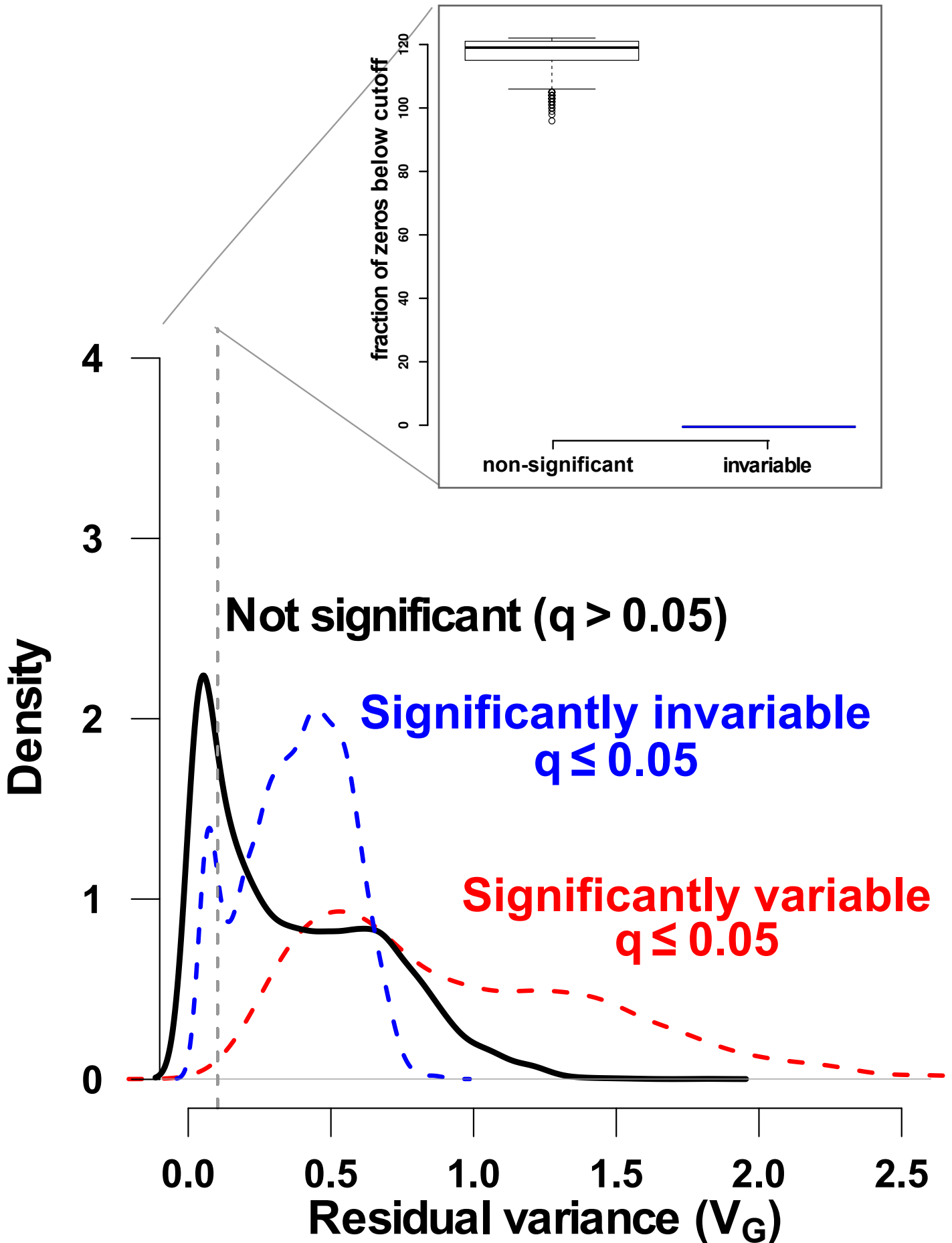


Figure 2—figure supplement 2

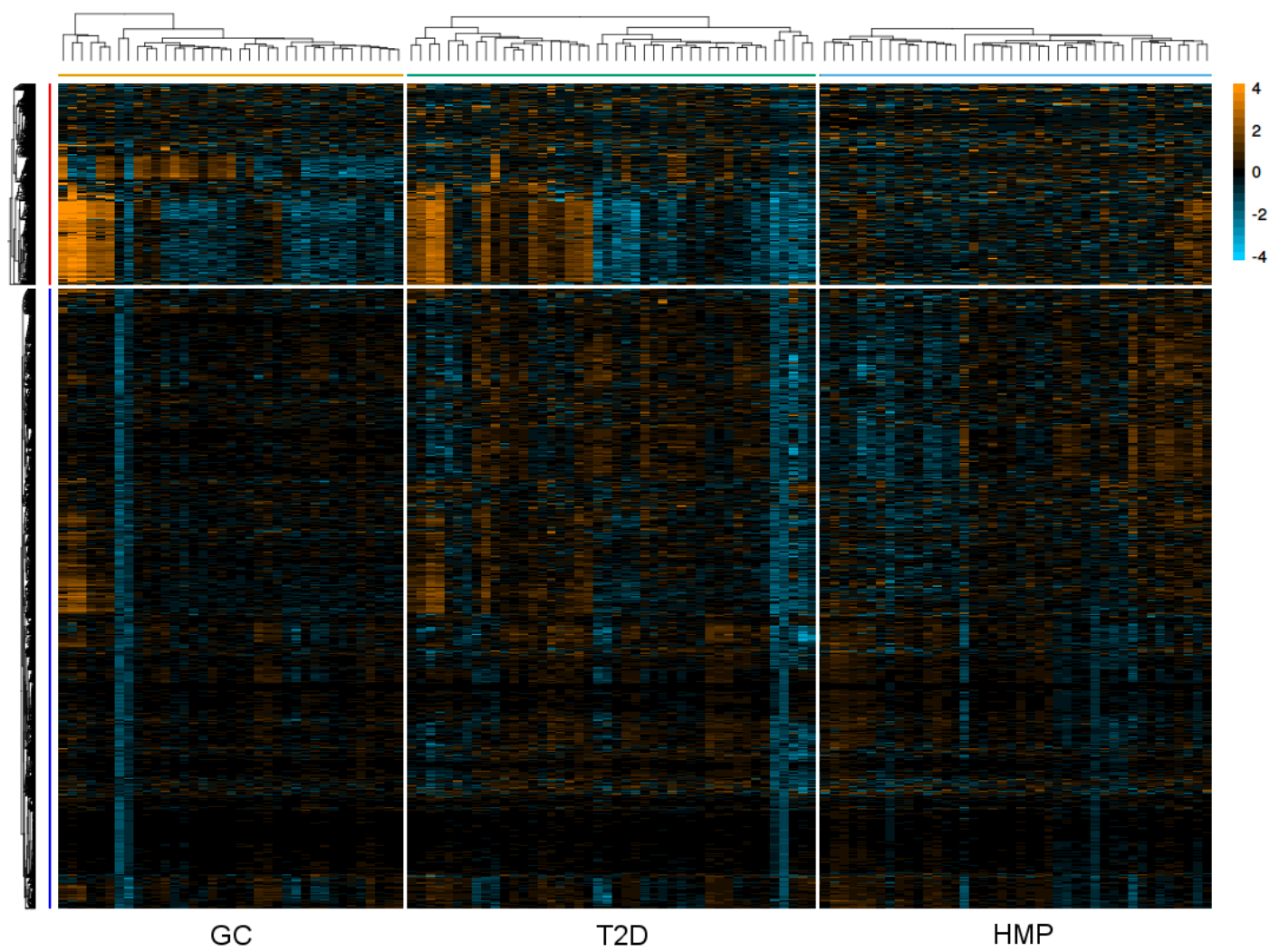


Figure 2—figure supplement 3

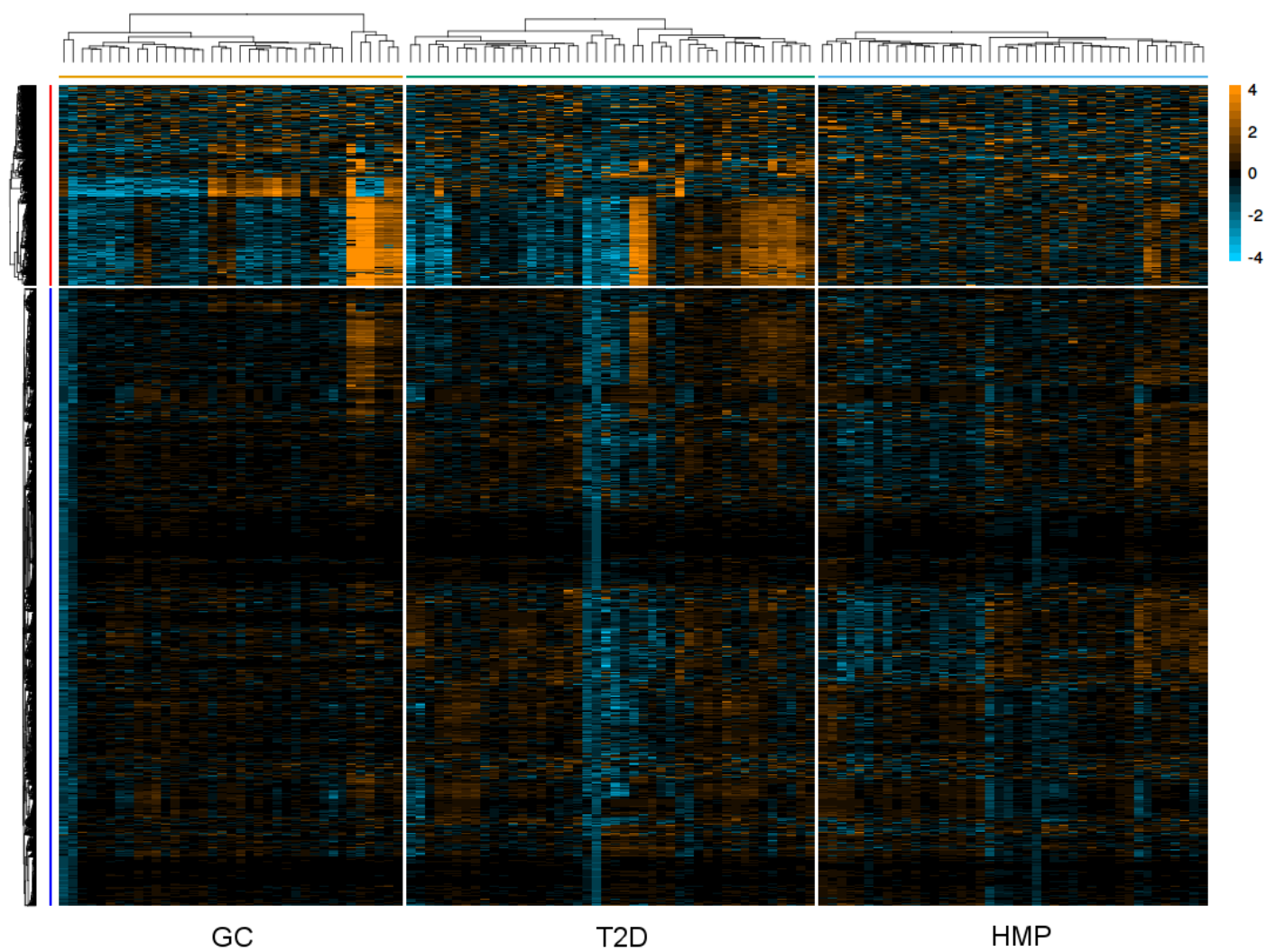


Figure 3—figure supplement 1

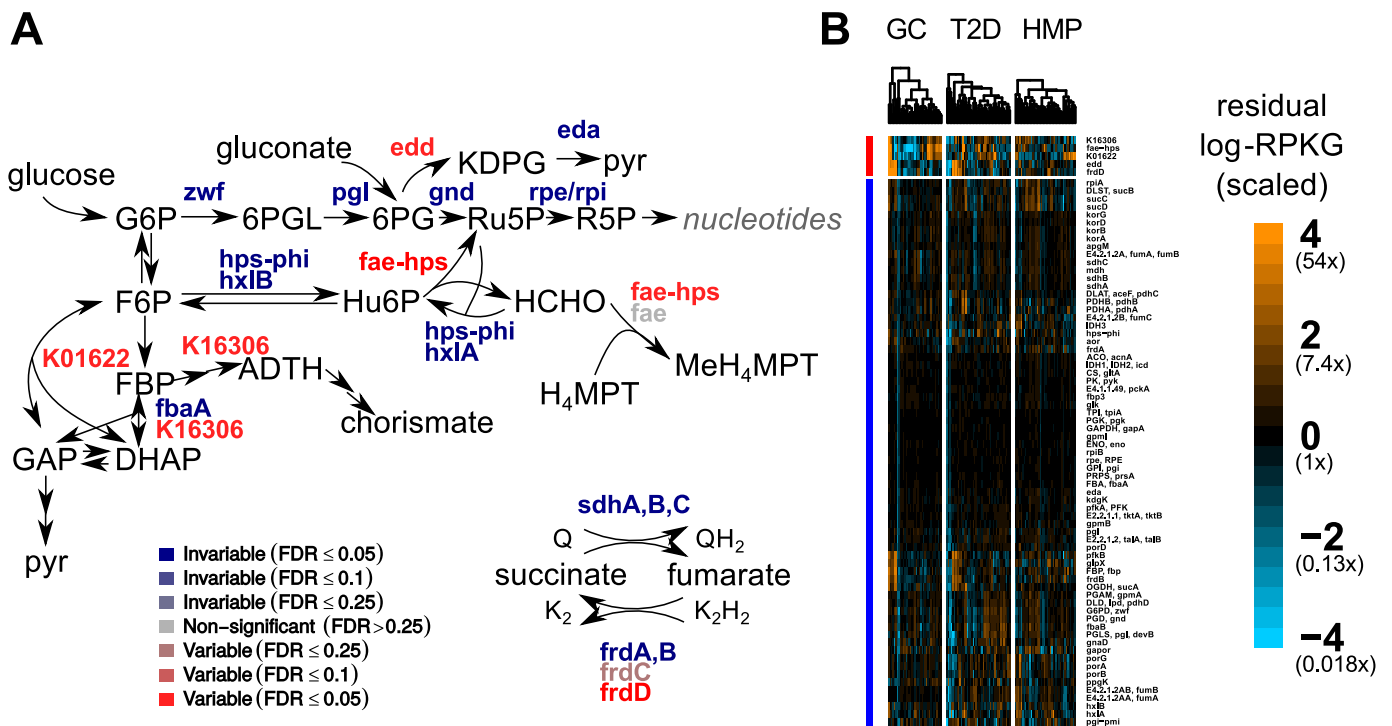


Figure 5—figure supplement 1

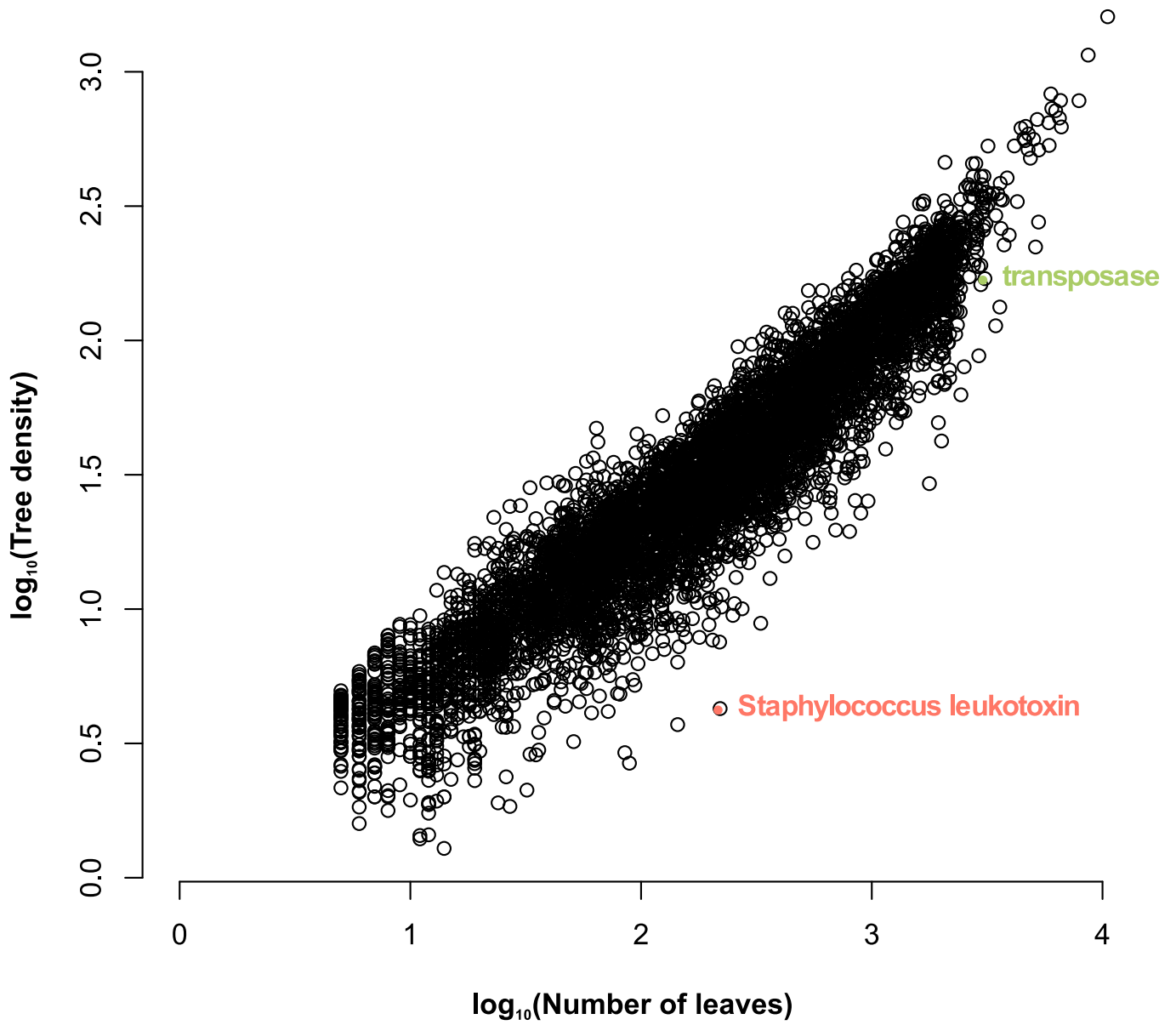


Figure 6—figure supplement 1

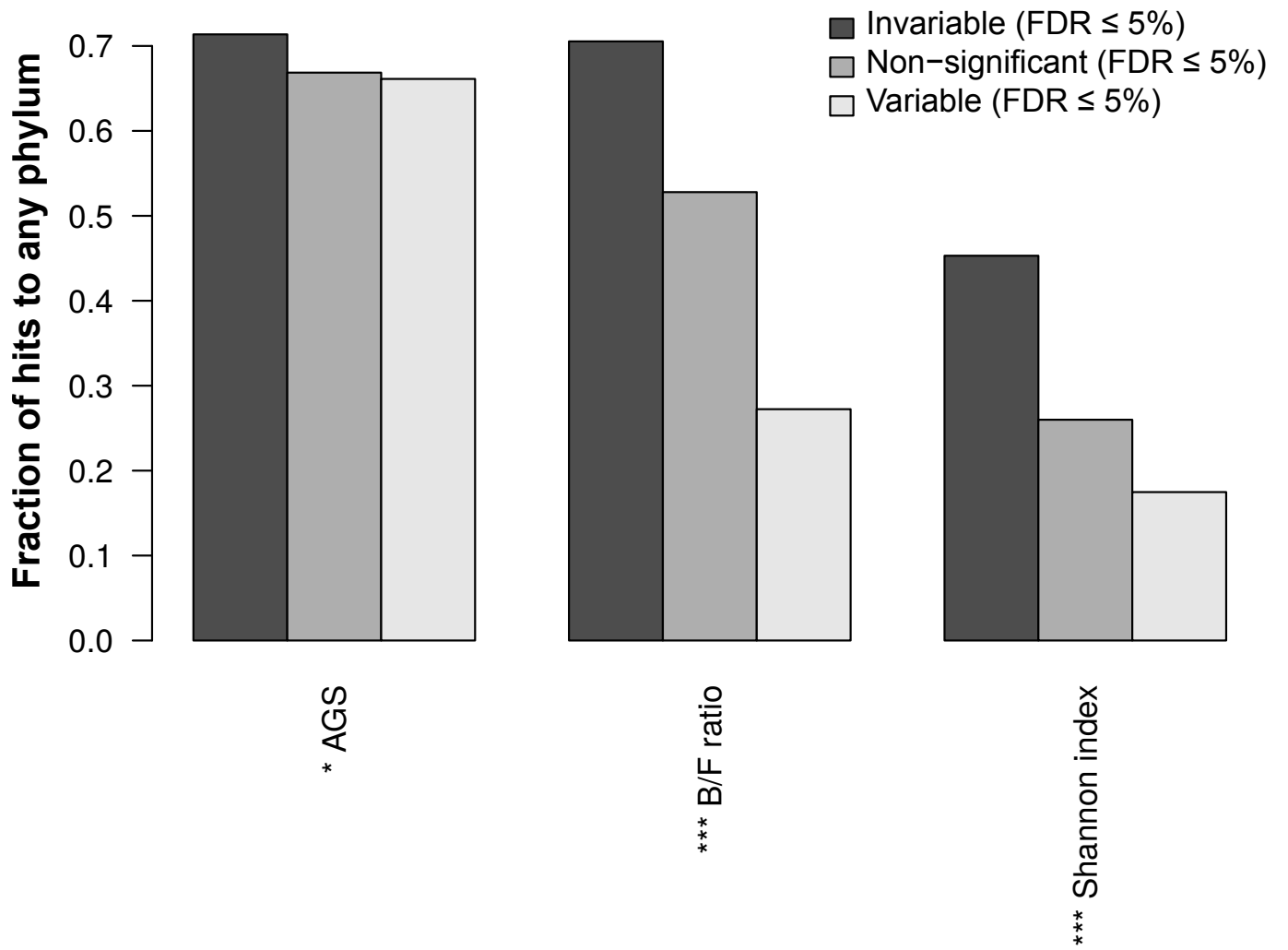


Figure 6—figure supplement 2

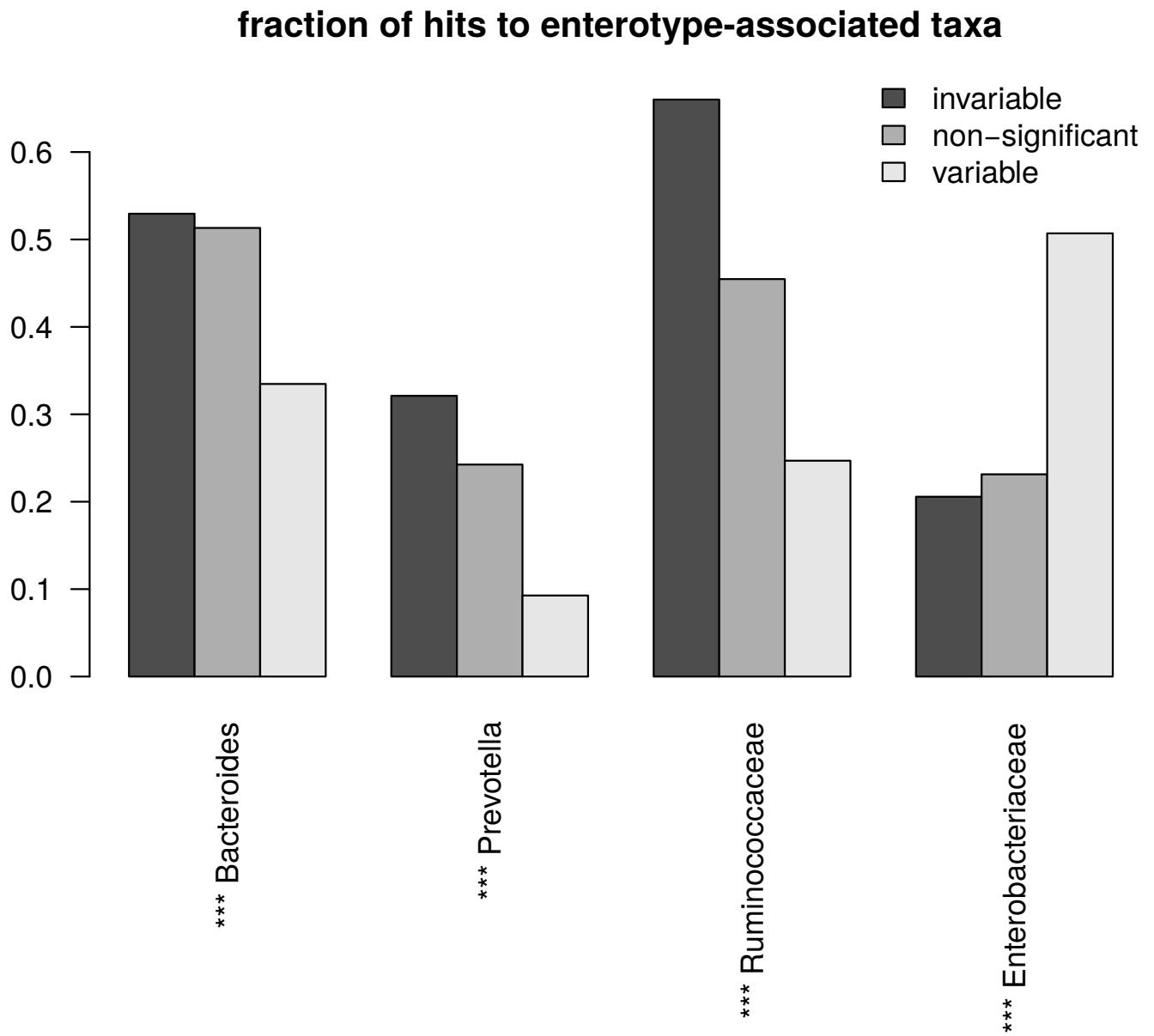


Figure 6—figure supplement 3

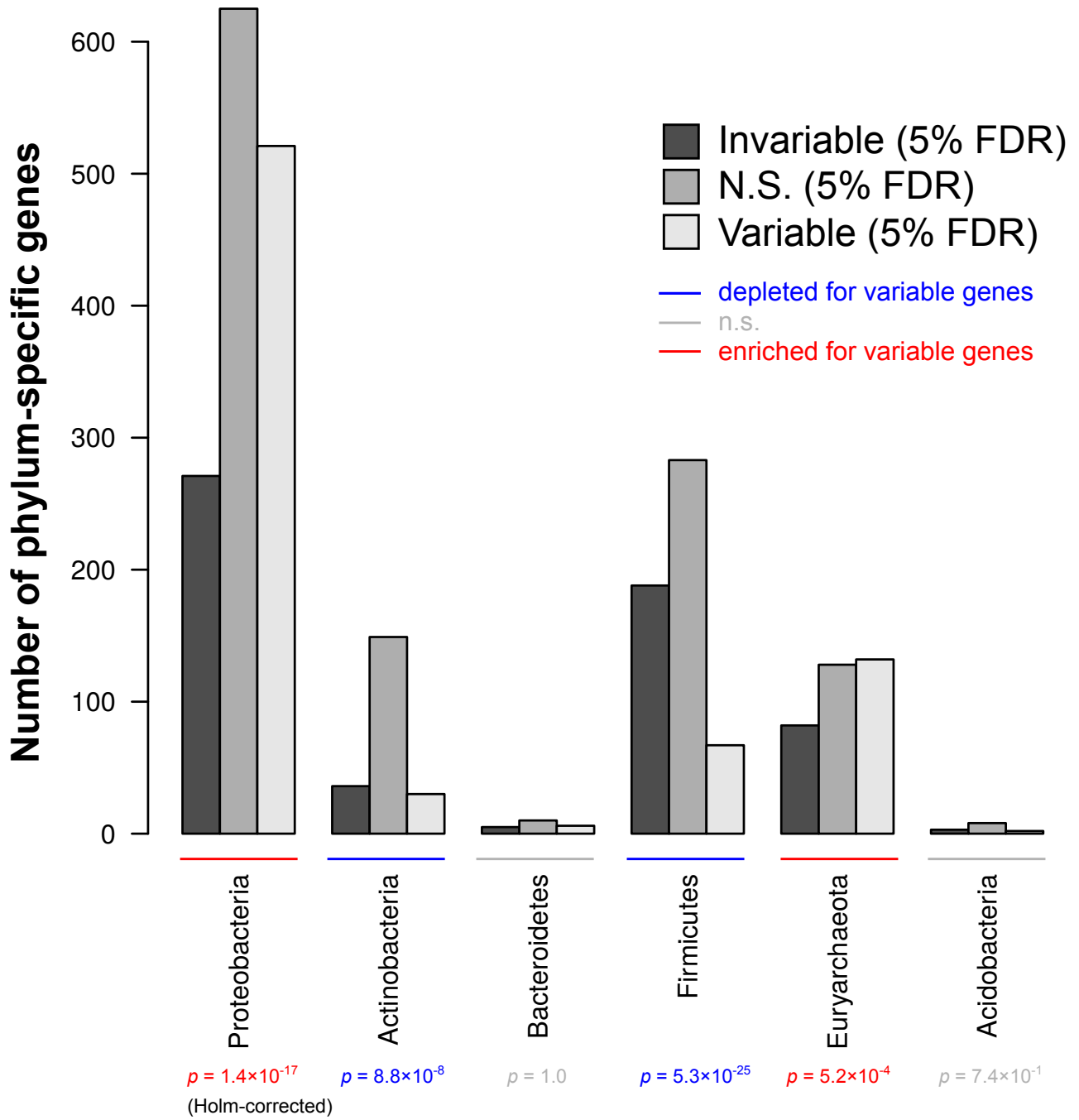


Figure 6—figure supplement 4

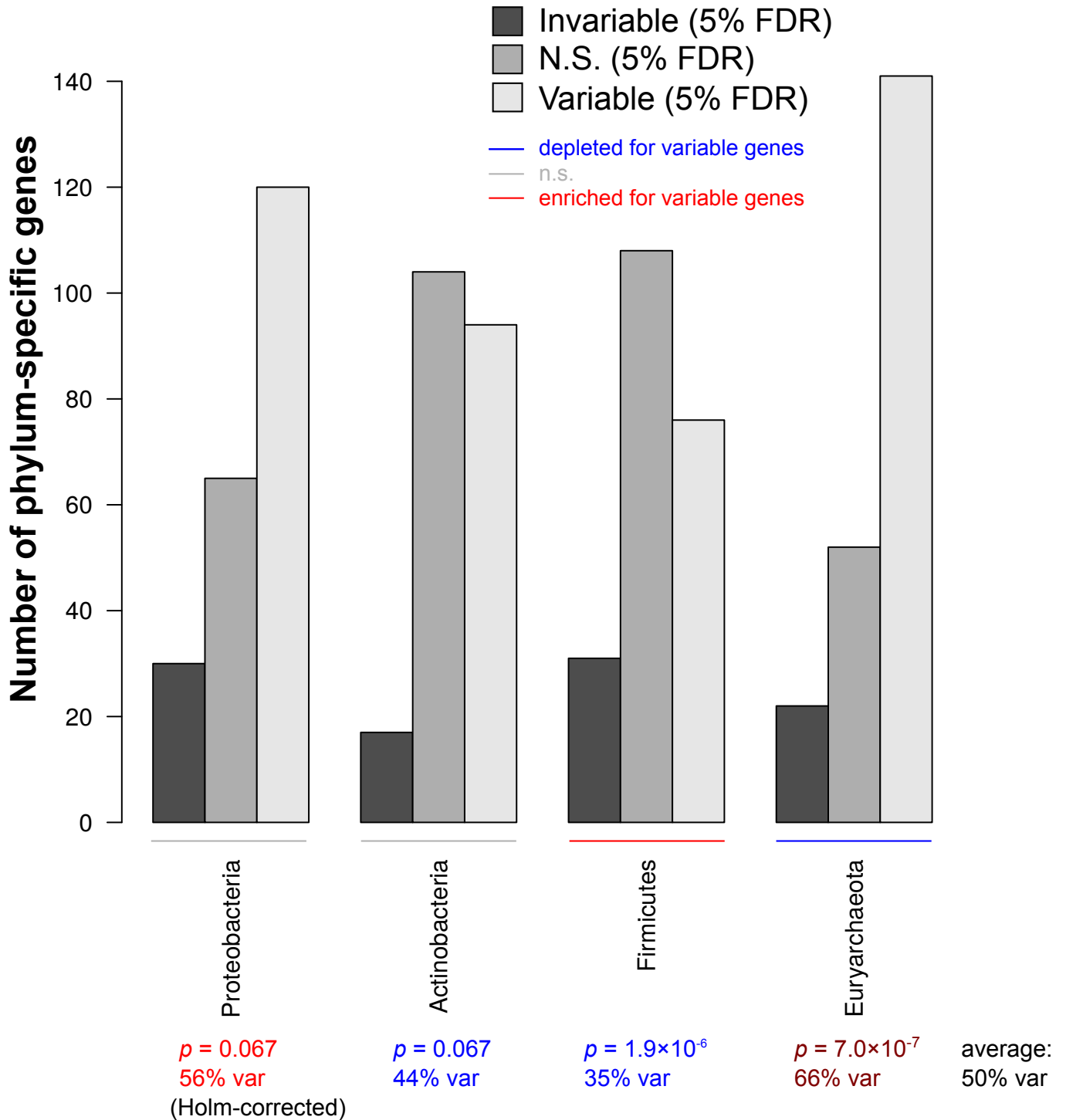


Figure 7—figure supplement 1

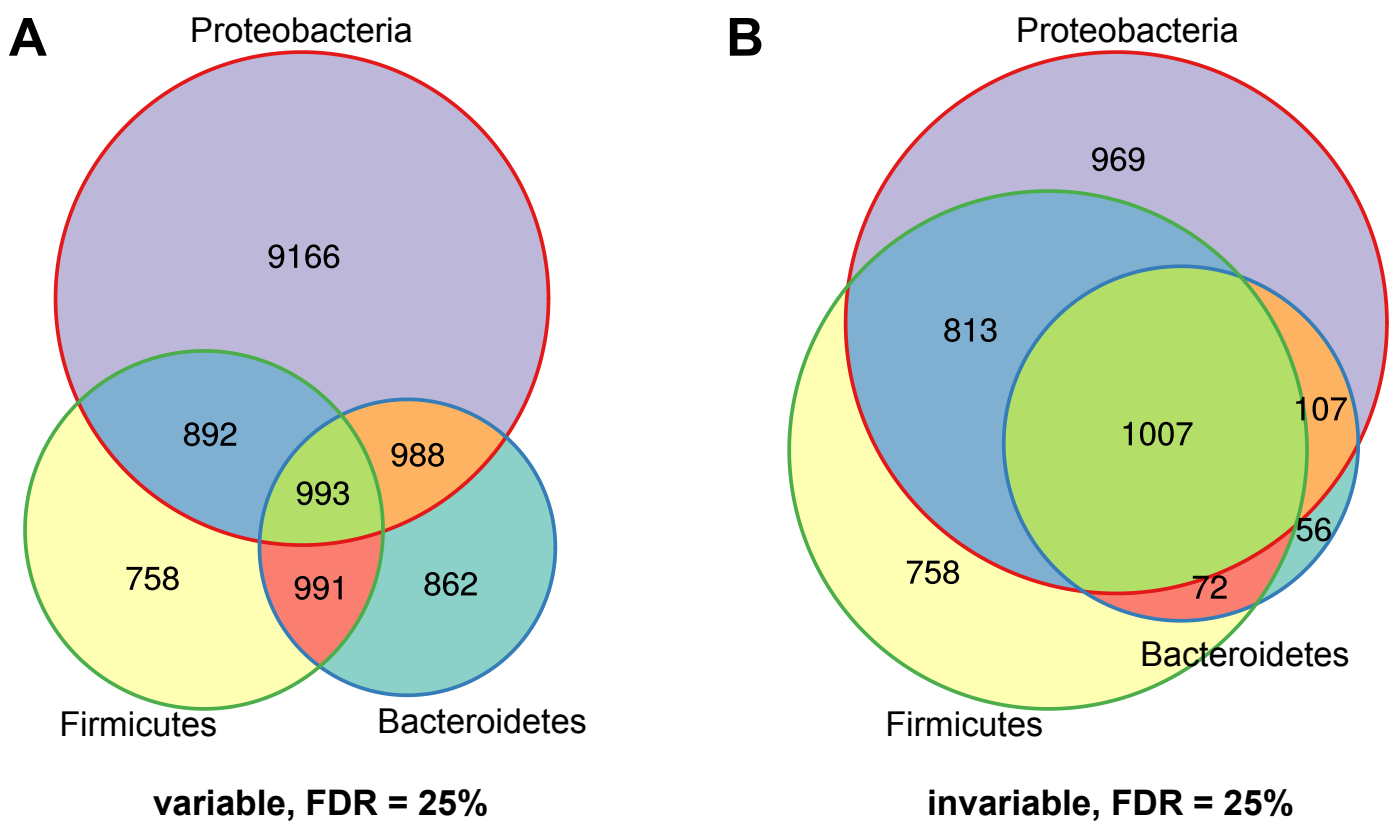
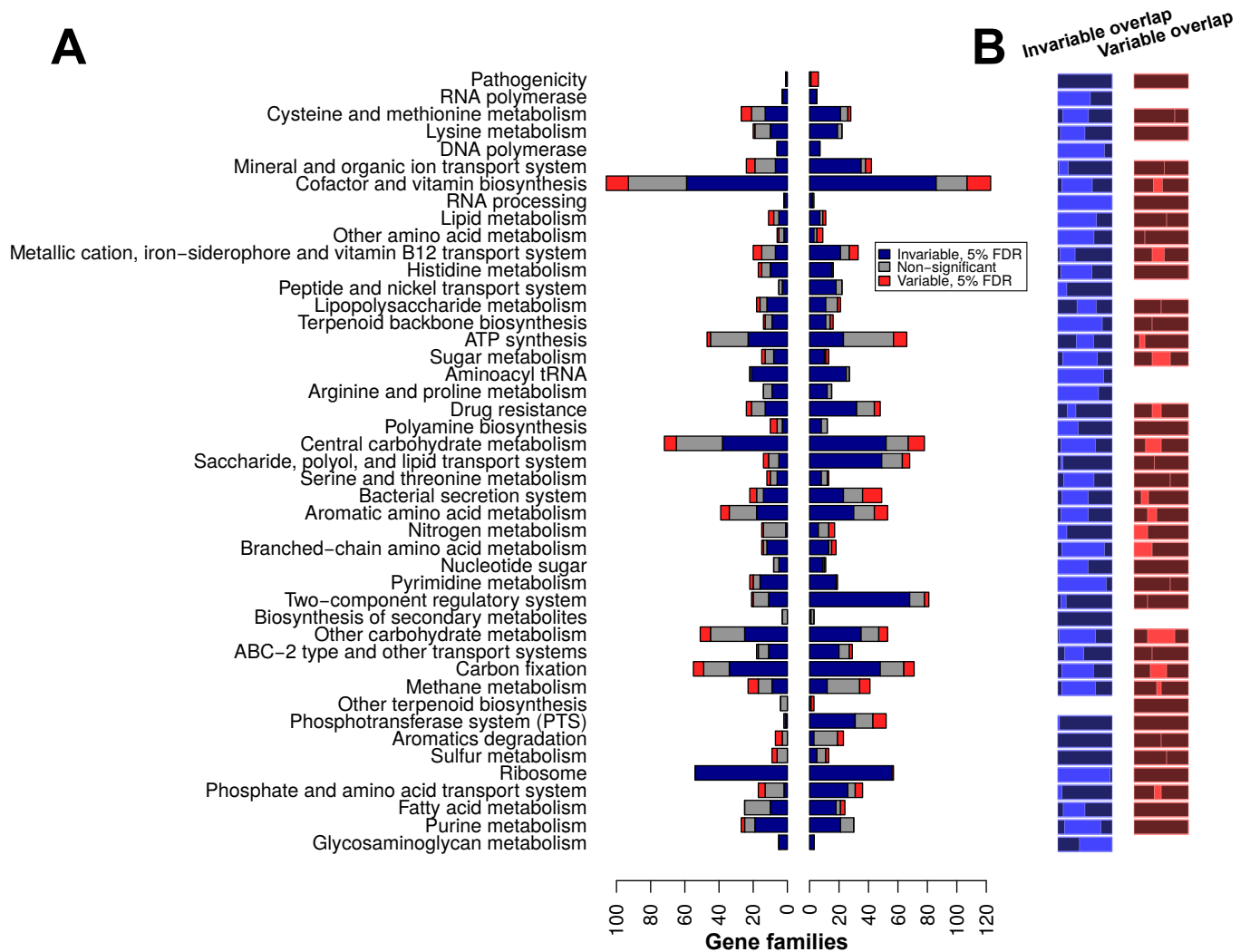


Figure 7—figure supplement 2



Supplemental Information

August 9, 2016

1 Correlation of variable and invariable gene families with taxonomic summary statistics

It has previously been suggested [1] that the genome size of gut microbiota reflects a trade-off between specialization (in which metabolic pathways for the production of reliably present nutrients may be lost over time, potentially resulting in auxotrophy) and generalization, or the ability to survive and grow in different metabolic conditions (which may require more biosynthetic genes). AGS itself has also been linked to health outcomes; for instance, individuals with Crohn's disease tend to have gut microbiota with larger genome size [2]. However, variable gene families were no more likely to be associated with AGS. Only 66% of variable gene families (with at least one bacterial or archaeal representative) had abundances that were significantly correlated with average genome size ($q \leq 0.05$), compared to 71% of invariable gene families and 66% of non-significant families at the same threshold. Thus, genome size correlates generally with gene abundance but does not predict variability of genes in healthy hosts.

The most dominant phylum-level trend across healthy human gut microbiomes is the trade-off between the two dominant phyla, Bacteroidetes and Firmicutes. The ratio of these two phyla (B/F ratio) has been linked to obesity in some studies [3, 4]; however, a later meta-analysis [5] revealed no consistent correlation across studies. Here, we found that variable genes were actually substantially *less* likely to be correlated to the B/F ratio (27%, $q \leq 0.05$) than either invariable (71%) or non-significantly-associated (55%) genes. These results parallel what we observe when we correlate gene family abundances with the α -diversity of observed bacterial species. We estimated α -diversity using the Shannon index, which is low when the distribution of species abundance is highly skewed, and high when there are many species of even abundance. Only 17% of significantly variable genes correlate significantly to the Shannon diversity ($q \leq 0.05$), versus 45% of significantly invariable and 26% of non-significant genes. We therefore conclude that bacterial and archaeal gene families identified as variable in this study are less likely to be associated with average genome size, B/F ratio, or α -diversity.

When examining the PD-stratified gene families, we noticed that the variable/high-PD gene set was also enriched for gene families described as “hypothetical” in the KEGG Orthology database; hypothetical gene families were also observed in the invariable/low-PD set, but they were statistically depleted (see main text). We were interested in whether these conserved-yet-variable hypothetical gene families could be acting as markers for minor phyla. Indeed, out of 81 genes in this group, 44 were significantly associated with Proteobacterial abundance ($q \leq 0.05$ by the above Kendall's partial τ test) and 13 were associated with Actinobacteria at the same threshold. However, 5 and 7 each were associated with Firmicutes and Bacteroidetes, indicating that even the major phyla of the human gut vary with respect to

certain as-yet-uncharacterized functions.

References

- [1] Nayfach S and Pollard KS. “Average Genome Size Estimation Improves Comparative Metagenomics and Sheds Light on the Functional Ecology of the Human Microbiome.” *Genome Biology*, **16** (2015):51. ISSN 1474-760X. doi:10.1186/s13059-015-0611-7.
- [2] Nayfach S, Bradley PH, Wyman SK, Laurent TJ, Williams A, Eisen JA, Pollard KS, and Sharpton TJ. “Automated and Accurate Estimation of Gene Family Abundance from Shotgun Metagenomes.” *PLoS Computational Biology*, **11** (2015)(11):e1004573. ISSN 1553-7358. doi:10.1371/journal.pcbi.1004573.
- [3] Turnbaugh PJ, Hamady M, Yatsunencko T, Cantarel BL, Duncan A, Ley RE, Sogin ML, Jones WJ, Roe BA, Affourtit JP, et al. “A Core Gut Microbiome in Obese and Lean Twins.” *Nature*, **457** (2009)(7228):480–4. ISSN 1476-4687. doi:10.1038/nature07540.
- [4] Ley RE, Bäckhed F, Turnbaugh P, Lozupone CA, Knight RD, and Gordon JI. “Obesity Alters Gut Microbial Ecology.” *Proceedings of the National Academy of Sciences of the United States of America*, **102** (2005)(31):11070–5. ISSN 0027-8424. doi:10.1073/pnas.0504978102.
- [5] Finucane MM, Sharpton TJ, Laurent TJ, and Pollard KS. “a Taxonomic Signature of Obesity in the Microbiome? Getting to the Guts of the Matter”. *PLoS ONE*, **9** (2014)(1):e84689. ISSN 1932-6203. doi:10.1371/journal.pone.0084689.