

1 OGU enable effective, phylogeny-aware analysis of even shallow  
2 metagenome community structures

3 Qiyun Zhu <sup>a,b,#</sup>, Shi Huang <sup>b,c</sup>, Antonio Gonzalez <sup>b</sup>, Imran McGrath <sup>b,d</sup>, Daniel McDonald <sup>b</sup>, Niina  
4 Haiminen <sup>e</sup>, George Armstrong <sup>b,c,f</sup>, Yoshiki Vázquez-Baeza <sup>c</sup>, Julian Yu <sup>a</sup>, Justin Kuczynski <sup>g</sup>, Gregory  
5 D. Sepich-Poore <sup>h</sup>, Austin D. Swafford <sup>c</sup>, Promi Das <sup>b,i</sup>, Justin P. Shaffer <sup>b</sup>, Franck Lejzerowicz <sup>b,c</sup>, Pedro  
6 Belda-Ferre <sup>b</sup>, Aki S. Havulinna <sup>j,k</sup>, Guillaume Méric <sup>l,m</sup>, Teemu Niiranen <sup>j,n,o</sup>, Leo Lahti <sup>p</sup>, Veikko  
7 Salomaa <sup>j</sup>, Ho-Cheol Kim <sup>q</sup>, Mohit Jain <sup>r,s</sup>, Michael Inouye <sup>l,t</sup>, Jack A. Gilbert <sup>b,c,i</sup>, Rob Knight <sup>b,h,u,#</sup>

8  
9 <sup>a</sup> School of Life Sciences, Arizona State University, Tempe, Arizona, USA

10 <sup>b</sup> Department of Pediatrics, School of Medicine, University of California, San Diego, California, USA

11 <sup>c</sup> Center for Microbiome Innovation, Jacobs School of Engineering, University of California San Diego,  
12 La Jolla, California, USA

13 <sup>d</sup> Division of Biological Sciences, University of California San Diego, La Jolla, California, USA

14 <sup>e</sup> IBM T. J. Watson Research Center, Yorktown Heights, New York, USA

15 <sup>f</sup> Bioinformatics and Systems Biology Program, University of California, San Diego, California, USA

16 <sup>g</sup> Google, Mountain View, CA, USA

17 <sup>h</sup> Department of Bioengineering, University of California, San Diego, La Jolla, CA, USA

18 <sup>i</sup> Scripps Institution of Oceanography, University of California San Diego, La Jolla, California, USA

19 <sup>j</sup> Department of Public Health and Welfare, Finnish Institute for Health and Welfare, Helsinki, Finland

20 <sup>k</sup> Institute for Molecular Medicine Finland (FIMM), HiLIFE, University of Helsinki, Helsinki, Finland

21 <sup>l</sup> Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne,  
22 Victoria, Australia

23 <sup>m</sup>Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria,  
24 Australia

25 <sup>n</sup> Department of Internal Medicine, University of Turku, Turku, Finland

26 <sup>o</sup> Division of Medicine, Turku University Hospital, Finland

27 <sup>p</sup> Department of Computing, University of Turku, Turku, Finland

28 <sup>q</sup> IBM Almaden Research Center, San Jose, California, USA

29 <sup>r</sup> Department of Medicine, University of California, San Diego, California, USA

30 <sup>s</sup> Department of Pharmacology, University of California, San Diego, California, USA

31 <sup>t</sup> Department of Public Health and Primary Care, Cambridge University, Cambridge, UK

32 <sup>u</sup> Department of Computer Science and Engineering, University of California, San Diego, California,  
33 USA

34

35 Qiyun Zhu and Shi Huang contributed equally to this work. Author order was determined on the basis of  
36 project seniority.

37 # Correspondence: Qiyun Zhu ([qiyun.zhu@asu.edu](mailto:qiyun.zhu@asu.edu)), Rob Knight ([robknight@eng.ucsd.edu](mailto:robknight@eng.ucsd.edu))

38

39 **Running title:** OGU's for diversity analysis of metagenomic data.

40 **Word count for the abstract:** 227.

41 **Word count for the text:** 3,307.

42

## 43 Abstract

44 We introduce Operational Genomic Unit (OGU), a metagenome analysis strategy that directly exploits  
45 sequence alignment hits to individual reference genomes as the minimum unit for assessing the diversity  
46 of microbial communities and their relevance to environmental factors. This approach is independent  
47 from taxonomic classification, granting the possibility of maximal resolution of community  
48 composition, and organizes features into an accurate hierarchy using a phylogenomic tree. The outputs  
49 are suitable for contemporary analytical protocols for community ecology, differential abundance and  
50 supervised learning while supporting phylogenetic methods, such as UniFrac and phylofactorization,  
51 that are seldomly applied to shotgun metagenomics despite being prevalent in 16S rRNA gene amplicon  
52 studies. As demonstrated in one synthetic and two real-world case studies, the OGU method produces  
53 biologically meaningful patterns from microbiome datasets. Such patterns further remain detectable at  
54 very low metagenomic sequencing depths. Compared with taxonomic unit-based analyses implemented  
55 in currently adopted metagenomics tools, and the analysis of 16S rRNA gene amplicon sequence  
56 variants, this method shows superiority in informing biologically relevant insights, including stronger  
57 correlation with body environment and host sex on the Human Microbiome Project dataset, and more  
58 accurate prediction of human age by the gut microbiomes in the Finnish population. We provide Woltka,  
59 a bioinformatics tool to implement this method, with full integration with the QIIME 2 package and the  
60 Qiita web platform, to facilitate OGU adoption in future metagenomics studies.

61

## 62 Importance

63 Shotgun metagenomics is a powerful, yet computationally challenging, technique compared to 16S  
64 rRNA gene amplicon sequencing for decoding the composition and structure of microbial communities.  
65 However, current analyses of metagenomic data are primarily based on taxonomic classification, which  
66 is limited in feature resolution compared to 16S rRNA amplicon sequence variant analysis. To solve  
67 these challenges, we introduce Operational Genomic Units (OGUs), which are the individual reference  
68 genomes derived from sequence alignment results, without further assigning them taxonomy. The OGU  
69 method advances current read-based metagenomics in two dimensions: (i) providing maximal resolution  
70 of community composition while (ii) permitting use of phylogeny-aware tools. Our analysis of real-  
71 world datasets shows several advantages over currently adopted metagenomic analysis methods and the  
72 finest-grained 16S rRNA analysis methods in predicting biological traits. We thus propose the adoption  
73 of OGU as standard practice in metagenomic studies.

74

75 **Keywords:** Operational Genomic Unit, taxonomy independent, reference phylogeny, UniFrac,

76 supervised learning, metagenomics

77

## 78 Introduction

79 The rapidly developing field of shotgun metagenomics has inherited many analytical tools from the  
80 more mature field of 16S rRNA gene amplicon studies. For example, diversity analyses provided in  
81 platforms such as QIIME 2 (1) can be used for metagenomic analyses. To date, the typical  
82 metagenomics workflow starts with taxonomic profiling, which estimates the taxonomic composition of  
83 microbial communities by matching sequencing data against a reference database (2). The resulting  
84 matches are compiled into an unstructured feature table, with values usually in the form of relative  
85 abundances of taxonomic units at a fixed rank (e.g. genus or species level), followed by relevant  
86 statistical analyses.

87 In contrast, the current standard for 16S rRNA analysis involves more advanced feature extraction,  
88 including construction of amplicon sequence variants (ASVs), which have replaced operational  
89 taxonomic units (OTUs) to deliver the finest-possible resolution from amplicon data (3). Phylogeny-  
90 aware algorithms such as UniFrac (4) have been widely-adopted to model community diversity while  
91 considering how features interrelate owing to the accessibility of reference phylogenies (5, 6), and the  
92 availability of *de novo* and *a priori* phylogenetic inference methods (7). This wisdom should be adopted  
93 as well to metagenomics. Thanks to the advances in efficient sequence alignment algorithms, and the  
94 expansions of reference genome databases (8, 9) and phylogenomic trees (10, 11), it is now possible and  
95 increasingly preferable to develop a fine-resolution, structured data analysis strategy in shotgun  
96 metagenomics.

97 Therefore, we propose an alternative method for constructing metagenomic feature tables, in which  
98 features are no longer taxonomic units, but individual reference genomes from a database, and the  
99 feature counts are the number of sequences aligned to these genomes. We refer to such features as  
100 Operational Genomic Units (OGUs). This term, in an echo of OTU but replacing “taxonomic” with

101 “genomic”, highlights the nature of the genome-based, taxonomy-free analysis. Meanwhile,  
102 “operational” indicates that this method does not rely on the direct observation of member genomes of  
103 the community, but uses pre-defined reference genomes as a proxy to model the community  
104 composition. However, like ASVs, OGU are exact and do not rely on similarity thresholds as OTUs do.

105 An OGU table represents the finest-grained resolution of observed genomes in a microbial community  
106 relative to the reference database. As such it can be used to quantify the community structure and  
107 relationships in correlation with biological traits. It can also work well with cost-efficient “shallow”  
108 shotgun metagenomics (12), where limited sequencing depth (even below the previously recommended  
109 lower threshold of 500,000 sequences per sample) is adequate for assessing community structure. It  
110 further empowers tree-based analyses, such as UniFrac and phylofactorization, which is enhanced by  
111 using the “Web of Life” (WoL) reference phylogenetic tree that we recently developed to describe  
112 accurate evolutionary relationships among genomes (10).

113 We have implemented the method for generating OGU tables in the open-source bioinformatics tool,  
114 Woltka (<https://github.com/qiyunzhu/woltka>). This program serves as a versatile interface connecting  
115 choices of upstream sequence aligners (such as Bowtie2 and BLAST) and downstream microbiome  
116 analysis pipelines (such as QIIME 2). In addition to the standalone program, the package ships with a  
117 QIIME 2 (1) plugin to facilitate adoption and integration into existing protocols. We have also made this  
118 method available through the Qiita web analysis platform (13) as part of the standard operating  
119 procedure for shotgun metagenomic data analysis, thereby enabling massive reprocessing and  
120 subsequent meta-analysis of metagenome datasets with OGUs. Thus far, we have applied the OGU  
121 method to re-analyze all public and private metagenomic datasets hosted on Qiita, totaling 143 studies  
122 and 57,063 samples, as of Mar 3, 2021.

123 Our team and collaborators have applied prototypes of the OGU method in multiple microbiome and  
124 multiomics studies and have obtained biologically relevant results (e.g., (14–16)). In this article, we  
125 systematically introduce the principles and practices of the OGU method, demonstrate its efficacy in one  
126 synthetic and two real-world microbiome datasets, and compare it with state-of-the-art metagenome  
127 analysis approaches and the alternative data type (16S rRNA gene amplicons). Given our findings, we  
128 propose the adoption of OGUs as a good practice in metagenomic analyses.

## 129 Results

### 130 **OGUs maximize resolution of community structures**

131 The rationale and benefits of the OGU method are demonstrated with a synthetic case study illustrated in  
132 [Fig. 1](#), with the underlying feature tables provided in [Table S1](#). In this simple case, three metagenomes  
133 with 12 sequences each were aligned to 10 reference genomes, which were hierarchically organized by  
134 taxonomy (left) or by phylogeny (right) ([Fig. 1A](#)). Beta diversity was calculated on feature tables at  
135 different levels: either on taxonomic units at the rank of genus or species, or directly on reference  
136 genomes (i.e., OGUs) without the need for giving them taxonomic labels.

137 As demonstrated ([Fig. 1B](#)), the genus-level analysis, which had the lowest resolution (three genera),  
138 yielded spurious proximity between samples B and C, as relative to sample A, largely determined by the  
139 differential abundance of genera G1 and G2. The species-level analysis with moderately higher  
140 resolution (five species) was able to bring A closer to B and C, mainly contributed by the identical  
141 frequencies of species S1, which could not be revealed at the genus level. The OGU-level analysis,  
142 having the highest resolution (10 features), revealed the separation between B and C due to distinct  
143 OGU composition, despite similar species counts (e.g., O5 and O7 have different counts within S3), and  
144 the proximity between A and B due to shared OGUs (O6 and O9). Additional structure was revealed by



145 using the UniFrac metric, which considers the hierarchical relationships among features, hence further  
146 joining samples (here A and B) sharing longer branches in the phylogenetic tree (even by different  
147 OGUs, such as O1 and O2) and separating those sharing shorter ones. Taxonomy may serve as a  
148 replacement of phylogeny, but it has a lower resolution than phylogeny (e.g., O1 and O2 are  
149 evolutionarily closer to each other relative to O3 but taxonomy cannot reveal this), and sometimes does  
150 not reflect the true evolutionary relationships among organisms (e.g., O4 and O5 are here placed in  
151 different genera), which can impact the accurate modeling of community structures.

152 In summary, this example illustrates the need for increasing resolution in order to better understand the  
153 diversity of microbial communities. This “resolution” has two dimensions of meaning: first, the quantity  
154 of features representing individual microbiomes; second, the granularity and accuracy of the hierarchy—  
155 if any—that defines the relationships among individual features.

### 156 **OGUs accurately represent body environment and host sex associated microbiome patterns**

157 We demonstrated the typical use of the OGU method on the classic Human Microbiome Project (HMP)  
158 shotgun metagenomic dataset (17), which contains 210 metagenomes sampled from seven body sites of  
159 male and female human subjects. We subsampled each metagenome to one million paired-end reads—a  
160 sampling depth close to the recommended lower threshold (500k reads) for “shallow” shotgun  
161 sequencing (12). The sequences were aligned to the WoL reference genome database (totaling 10,575  
162 bacterial and archaeal genomes) and the alignments were processed using Woltka, resulting in an OGU  
163 table with 6,220 features (reference genomes) (Fig. S1A). Beta diversity analysis using the weighted  
164 UniFrac metric with the WoL reference phylogeny was performed on the OGU table (Fig. 2). For  
165 comparison, we analyzed the dataset using the currently adopted method (CAM) (e.g., (17)): using Bray-  
166 Curtis on a species-level taxonomic profile. We exemplified the CAM by using the profile inferred by

167 Bracken (18) on the same WoL database (Fig. 2), but also tested and reported the results of SHOGUN  
168 (19), Centrifuge (20), and MetaPhlAn (21) (Fig. S1).

169 Principal Coordinates Analysis (PCoA) of OGU (Figs. 2A and S2A), with the first three axes  
170 explaining 71.01% of community structure variance (Figs. 2C and S1B), revealed that microbiomes  
171 were clustered mainly by the body site from which they were sampled, which overshadowed clustering  
172 by host sex, if any. This pattern is largely consistent with the previous report (17). The PCoA plot by  
173 CAM (Figs. 2B and S2B, also see S3), although with less explained variance (46.30%) (Figs. 2C and  
174 S1B), also displayed a clustering-by-site pattern. However, it is notable from the plot that sample  
175 clusters are aligned diagonally—a typical pattern indicating the saturation of distances caused by the  
176 inadequacy of shared features (species) among body sites (22) (Figs. 2B and S2B). This characteristic  
177 limits the power of resolving community diversity.

178 Permutational multivariate analysis of variance (PERMANOVA) of the beta diversity distance matrices  
179 suggested that all methods were able to clearly differentiate samples by body site ( $p=0.001$ ), with OGU  
180 generating the strongest statistic (Figs. 2E and S1C) (OGU:  $F=77.82$ ; CAM:  $F=42.36$ ). The distinction  
181 by host sex was less obvious. Only OGU was able to distinguish microbiome by sex ( $F=3.011$ ,  
182  $p=0.013$ ), whereas CAM failed to distinguish sex with statistical significance ( $F=1.692$ ,  $p=0.086$ ) (Figs.  
183 2F and S1E-F). This demonstrated the power of the OGU method in capturing subtle but relevant trends,  
184 even when another primary factor (body site) is driving most of the community diversity. Three of the  
185 seven body sites are located in the oral environment: tongue, teeth and buccal mucosa (Fig. 1A, B).  
186 They together indicate weaker differentiation by sex (OGU:  $F=1.905$ ,  $p=0.099$ ; CAM:  $F=1.610$ ,  
187  $p=0.130$ ) (Figs. 2F and S1G-H). In parallel, we reason that sites sharing the same environment likely  
188 have higher microbial connections. To test this effect, we calculated the relative distance between the  
189 three oral sites versus oral sites to non-oral sites. This distance is significantly smaller with OGU (0.699

190  $\pm 0.098$ , mean and std. dev., same below) than with CAM ( $0.808 \pm 0.051$ ) (two-tailed paired  $t=-14.398$ ,  
191  $p=2.57e-26$ ) (Figs. 2D and S1D), suggesting that OGU is more effective at relating subgroups of  
192 samples with shared properties.

193 The OGU table plus the WoL tree further enabled differential abundance analysis using the phylogenetic  
194 factorization method (23) (Figs. S4-5). The result was visualized and analyzed using the recently  
195 released massive tree visualizer EMPress (24) (Fig. 2G). It revealed that the phylogenetic clade  
196 separated by Factor 1 represents the genus *Lactobacillus*, contained in predominantly posterior fornix  
197 samples from female hosts, which is expected (25). Meanwhile, Factor 2 (genus *Neisseria*), Factor 3  
198 (genus *Capnocytophaga*) and Factor 4 (species *Leptotrichia buccalis*) are more frequently observed in  
199 the oral sites of male hosts. For comparison, we applied the tree-free method ANCOM (26) on the  
200 taxonomic profiles generated by alternative methods (Table S2). At genus level, all four methods were  
201 able to capture only *Lactobacillus*, consistent with our Factor 1. However, at species and OGU levels,  
202 results were discordant between methods and no method reported any *Lactobacillus* sp., again showing  
203 the limitations of confining analyses to taxonomic ranks without phylogenetic information.

204 Finally, we assessed the efficacy of OGUs along a gradient of decreasing sampling depths. The  
205 correlation between the original OGU table (from one million paired-end reads) and each of the  
206 subsampled OGU tables was consistently high. A Pearson's  $r$  of  $0.961 \pm 0.0726$  (mean and std. dev.,  
207 same below) was retained even at the sampling depth of 200 (Fig. S6A). The PCoA clustering pattern  
208 largely remained the same at all sampling depths (Fig. S7). The oral-vs-other relative distance (see  
209 above) retained a Pearson's  $r$  of  $0.971 \pm 0.00613$  when sampling depth was 200 (Fig. S6B). The  
210 PERMANOVA  $F$ -statistics calculated based on 10 replicates of random subsampling were close to the  
211 original statistic and largely stable down to very low sampling depths. The mean difference from the  
212 original statistic was still within 5% at the sampling depth of 1,000 for body site ( $3.349 \pm 1.361$ , unit:

213 percentage of the original statistic, same below), or 500 for host sex ( $2.680 \pm 5.473$ ) (Fig. S6C-D).

214 These findings suggest that the OGU method remains valid even on very shallow metagenomic samples,  
215 including those that would otherwise be considered unusable for typical metagenomic analyses.

## 216 **OGUs improve prediction of host age from the gut microbiome**

217 We next analyzed 6,430 stool samples collected through a random sampling of the Finnish population  
218 using both 16S rRNA gene amplicon sequencing and shallow shotgun metagenomic sequencing. This  
219 “FINRISK” study (27) provides an opportunity to explore the dependency of feature sets (e.g.  
220 taxonomic levels and data source: 16S rRNA amplicon vs. shotgun metagenomic data) on the prediction  
221 accuracy of a machine learning model on the targeted phenotype (e.g., age). We quantitatively examined  
222 the impact of taxonomic level of microbiome features on the empirical error (mean absolute error, or  
223 MAE) in predicting human chronological age using a Random Forests regressor (28), constructed using  
224 5-fold cross-validation.

225 Our results (Fig. 3A) showed the prediction accuracy continued to improve, resulting in lower absolute  
226 errors with finer microbial feature classification levels. Shotgun data outperformed 16S data at all levels,  
227 and was able to reduce MAE to less than 10 years at the genus level or below. At the lower limit of both  
228 16S and shotgun data, we achieved an MAE of  $9.581 \pm 0.116$  years (mean and std. dev., same below)  
229 with OGUs (Fig. 3B), whereas ASVs, the highest possible resolution allowed by 16S data, resulted in a  
230 higher MAE of  $10.110 \pm 0.103$  years (two-tailed  $t=-7.25$ ,  $p=8.81e-5$ ). Meanwhile, using the species-  
231 level profile inferred by Bracken, we also obtained a higher MAE of  $10.273 \pm 0.089$  years (vs. OGU:  
232 two-tailed  $t=-10.59$ ,  $p=5.53e-6$ ) (Fig. S8). Decreasing sequencing depth did not reduce the age  
233 prediction accuracy for individual samples (Fig. S9). For example, samples with 320-366k metagenomic  
234 sequences (2nd bin from low end in the figure) had an MAE of  $9.290 \pm 6.378$  years, whereas samples  
235 with 1,386-1,931k sequences (2nd bin from high end) had an MAE of  $10.118 \pm 6.086$  years, which were

236 not significantly different (two-tailed  $t=-1.37$ ,  $p=0.170$ ). We then explored which OGU contributed to  
237 the superior performance in age prediction as compared to 16S rRNA ASVs. Therefore, we identified a  
238 reduced set ( $n=128$ ) of the most important OGUs that can maximize the prediction accuracy via a  
239 recursive feature elimination approach (Fig. S10). Among these important features, a few gut microbial  
240 strains increased in abundance with aging, such as multiple strains from *Streptococcus mutans*,  
241 *Eubacterium sp.* (Figs. 3C, S11-12). Remarkably, those *Streptococcus spp.* are typically located in the  
242 oral cavity yet can be over-represented in the gut of elderly individuals, suggesting potential microbial  
243 transmissions between oral and gut microbiomes related to typical aging in a large population (29, 30).  
244 Next, we also identified a few microbial OGUs that were under-represented in the elderly, such as  
245 *Anaerostipes hadrus* DSM 3319 and members of *Bifidobacterium*, including *B. longum* NCC2705 and  
246 *B. saguini* DSM 23967 Bifsag. Many of these important taxonomic features were not identified in the  
247 16S data, putatively because the partial sequences of a 16S rRNA gene cannot provide sufficient  
248 resolution to distinguish species or strains. For example, a few 16S rRNA ASVs annotated with  
249 *Lachnospiraceae* have been associated with aging and were identified in either this or past studies (31),  
250 whereas our method identified several OGUs (*Anaerostipes hadrus* DSM 3319) within the family of  
251 *Lachnospiraceae* that exhibited strong predictive powers for discriminating aging.

## 252 Discussion

253 The OGU method introduced in this article provides a way to maximize the resolution of feature tables  
254 by directly considering reference genomes without the reliance on taxonomic classification in shotgun  
255 metagenomics studies. Although the strategy of taxonomy-free community structure analysis has been  
256 widely adopted in 16S data analysis (e.g., ASV or *de novo* OTU clustering), it remains underexplored in  
257 metagenomics, largely due to the difficulties in defining and quantifying “features” without using an *a*

258 *priori* classification system. Our study shows that sequence alignment hits to individual reference  
259 genomes can be used as the minimum unit for features, referred to as OGU.

260 Through comparative analysis of OGU and alternative methods using a synthetic case study and two  
261 real-world microbiome studies, we demonstrated that classical high-dimensional statistics and machine  
262 learning methods developed and matured in the field of 16S rRNA gene amplicon analysis can be  
263 directly applied to OGUs to provide biologically relevant insights. The OGU results often are superior to  
264 currently adopted metagenomic classification methods and ASV analysis of the 16S rRNA data.  
265 Meanwhile, we showed that the use of taxonomic units as features, as many researchers have been  
266 practicing to date, has conceptual and performance limitations compared with the OGU method,  
267 particularly at higher taxonomic ranks due to the loss of resolution.

268 The independence from taxonomy further enables the utilization of explicit phylogenetic trees. A  
269 researcher can choose from pre-computed reference phylogenies, such as the one we introduced in the  
270 “Web of Life” (WoL) project (10), or custom phylogenomic trees computed from *de novo* construction  
271 or placement, through tools such as PhyloPhlAn3 (32) and DEPP (33), which are scalable to large  
272 numbers of genomes. This connects evolutionary biologists’ efforts in updating the tree of life (e.g., (10,  
273 11, 34)), computational biologists’ efforts in forging phylogeny-aware methods (e.g., UniFrac and  
274 PhyloFactor), and microbiome scientists’ pursuits of relating high-dimensional microbiome data with  
275 biology.

276 Taxonomy, despite being relatively coarse-grained and error-prone as a classification system, may serve  
277 as an implicit replacement of phylogeny if the latter is not available. We tested this idea by applying  
278 UniFrac to an artificial taxonomic tree with constant branch lengths between ranks (analogous to (35)).  
279 Although this treatment is controversial, because taxonomic ranks do not directly indicate evolutionary  
280 distances, we did observe improvement compared to not using a tree (Fig. S13). Although there have

281 been remarkable efforts for curating taxonomy using phylogenetics, however, the number of taxonomic  
282 ranks is limited (typically 7 to 8), and can constrain the topology for an ever-growing number of  
283 sequenced genomes. For example, the current release (R95) of GTDB (36) has 31,910 species clusters,  
284 constituting a taxonomy tree of 45,502 vertices, whereas NCBI RefSeq and GenBank host 977,729  
285 unique genomes as of March 30, 2021, and a fully resolved phylogenetic tree of them can theoretically  
286 have 1,955,456 vertices. The history of 16S rRNA studies (7) is repeating itself in whole-genome  
287 studies, such that building a phylogeny is not only advantageous but often more feasible than defining  
288 taxonomy, and the OGU method powerfully provides an analogous extension to shotgun sequencing  
289 studies. As a new notion to microbiome research, OGU's properties in statistical analyses has yet to be  
290 characterized in a large number of studies, as was done for 16S rRNA ASVs. Unique challenges in  
291 shotgun metagenomics may impact analyses that were designed for 16S rRNA data. For example, very-  
292 low-abundance false positive assignments, which are prevalent from typical metagenomic classifiers,  
293 may impair the accuracy of the recovered community composition (37). A typical treatment is to only  
294 consider features with relative abundance above a given threshold in each sample (37). While we  
295 provide this function in Woltka to facilitate user's preferences, our tests suggested that the result of an  
296 OGU analysis is highly stable against a wide range of filtering thresholds when using abundance-based  
297 metrics (weighted UniFrac and Bray-Curtis), as compared with presence/absence-based metrics  
298 (unweighted UniFrac and Jaccard) (Fig. S14). This observation implies the OGU method is robust to  
299 noise commonly introduced into metagenomic datasets from many low abundance observations.

300 The robustness of an OGU analysis is only limited by the comprehensiveness of the reference. Despite  
301 that available genomic data have grown to an enormous volume, the size of a reference genome database  
302 that can be realistically used in a metagenomic analysis with typical computing facilities is  
303 circumscribed, limiting the increase of resolution beyond sub-species levels. Balancing alignment  
304 accuracy and database content is therefore an important consideration in designing the analytical

305 strategy. The algorithm we previously designed and used in the WoL database to maximize the covered  
306 biodiversity given a fixed number of genomes (10) may be beneficial in this situation, but its efficacy  
307 needs to be further tested in the background of various biospecimens and biological questions.  
308 Leaderboard sequencing may also be a useful strategy for iteratively augmenting the reference database  
309 with the common genomes in each sample (38). In the long run, efforts to improve algorithms, increase  
310 database coverage, and improve computing efficiency are all needed to facilitate effective advances in  
311 the field of metagenomics, and the OGU method provides an important step forward in that direction.

312



## 313 Materials and Methods

### 314 Protocol details

315 The OGU method is flexible to the type of sequence alignment. The recommended protocol, which is  
316 also the protocol demonstrated and benchmarked in this article, is as follows: Shotgun metagenomic  
317 sequencing data were aligned against the WoL reference genome database using SHOGUN v1.0.8 (19),  
318 with Bowtie2 v2.4.1 (39) as the backend. This process is equivalent to a Bowtie2 run with the following  
319 parameters:

```
320 --very-sensitive -k 16 --np 1 --mp "1,1" --rdg "0,1" --rfg "0,1" --score-min "L,0,-0.05"
```

321 The sequence alignment is treated as a mapping from queries (sequencing data) to subjects (reference  
322 genomes). It is possible that one sequence is mapped to multiple genomes (up to 16 using the  
323 aforementioned Bowtie2 command). In this scenario, each genome is counted  $1/k$  times ( $k$  is the  
324 number of genomes to which this sequence is mapped. The frequencies of individual genomes were  
325 summed after the entire alignment was processed, and rounded to the nearest even integer. Therefore,  
326 the sum of OGU frequencies per sample is nearly (considering rounding) equal to the number of aligned  
327 sequences in the dataset. The output feature table has columns as sample IDs, rows as feature IDs  
328 (OGUs), and cell values as the frequency of each OGU in each sample. This table is ready to be  
329 analyzed using software packages such as QIIME 2 (1).

### 330 Implementation

331 The OGU method is implemented in the bioinformatics tool Woltka (Web of Life Toolkit App), under  
332 the BSD-3-Clause open-source license. The program is written in Python 3, following high-quality  
333 software engineering standards. Its unit test coverage is 100%. The source code is hosted in the GitHub

334 repository: <https://github.com/qiyunzhu/woltka>, together with instructions, tutorials, command-line  
335 references, and test datasets. The program has been included in the Python Package Index (PyPI). In  
336 addition to the standalone Woltka program, a QIIME 2 (1) plugin is included in the software package.  
337 Woltka automatically recognizes and parses multiplexed or per-sample sequence alignment files, either  
338 original or compressed using Gzip, Bzip2 or LZMA algorithms. It supports three alignment file formats:  
339 1) SAM (Sequence Alignment Map) (40), which is supported by multiple short read alignment  
340 programs, such as Bowtie2 (39), BWA (41) and Minimap2 (42); 2) the standard BLAST (43) tabular  
341 output format (“-outfmt 6”), which is supported by multiple sequence alignment programs, such as  
342 BLAST, VSEARCH (44) and DIAMOND (45); 3) A plain mapping of query sequences to subject  
343 genomes, which is customizable to adopt other tools and pipelines.

344 In addition to OGU table generation, Woltka supports summarizing features into higher-level groups.  
345 This enables taxonomic classification, for comparison purposes. The output of Woltka’s classification  
346 function and that of SHOGUN’s “assign\_taxonomy” function are identical. Woltka supports three  
347 formats of classification systems: 1) the Greengenes-style lineage strings (supported by programs such  
348 as QIIME 2 (1), MetaPhlAn (21) and GTDB-tk (46)); 2) The NCBI-style taxonomy database (47) (a.k.a.  
349 “taxdump”, supported by programs such as Kraken 2 (48), Centrifuge (20) and DIAMOND (45)); 3)  
350 One or multiple plain mappings of child-to-parent classification units.

## 351 **Deployment**

352 The Woltka program has been incorporated in the Qiita web analysis platform (<https://qiita.ucsd.edu/>)  
353 (13), as part of the standard operating procedure for analyzing shotgun metagenomic data (qp-woltka,  
354 code hosted at: <https://github.com/qiita-spots/qp-woltka>). It can be directly launched from the graphic  
355 user interface. A job array system is used to parallelize analyses on a per-sample base to maximize

356 processing speed. Each process uses eight cores of an Intel E5-2640 v3 CPU and 90 GB DDR4 memory.  
357 Two reference genome databases are available for user choice: 1) The “Web of Life” (WoL) database  
358 (10), with 10,575 bacterial and archaeal genomes that were evenly sampled through an algorithm. 2) The  
359 reference and representative genomes of microbes defined in NCBI RefSeq release 200 (8). The  
360 subsequent community ecology analyses based on the OGU table are also available from Qiita. The  
361 WoL reference phylogeny is available for choice for phylogenetic analyses (such as UniFrac (4)).  
  
362 This system allowed us to re-analyze all metagenomic datasets hosted on Qiita (totaling 143 studies and  
363 57,063 samples, as of Mar 3, 2021) to generate OGU tables as well as tables at multiple taxonomic  
364 ranks, which are ready for subsequent meta-analysis by Qiita users. Although runtime varies by sample  
365 size, the average wall clock time for analyzing one metagenomic sample (including sequence alignment  
366 against WoL using Bowtie2 and feature table generation using Woltka) was 13.8 minutes in this large  
367 effort.

### 368 **The HMP dataset**

369 The Human Microbiome Project (HMP) (17) dataset was downloaded from the official website  
370 (<https://www.hmpdacc.org/hmp/>). It contains 241 samples of 100 bp paired-end whole genome  
371 sequencing (WGS) reads. The sequencing data were already processed to remove human contamination  
372 and low-quality regions. We dropped samples with less than 1M paired-end reads, leaving 210 samples.  
373 They were randomly subsampled to 1M paired-end reads per sample. These samples represent both male  
374 ( $n=138$ ) and female ( $n=72$ ) human subjects. They represent seven body sites: stool ( $n=78$ ), tongue  
375 dorsum ( $n=42$ ), supragingival plaque ( $n=33$ ), buccal mucosa ( $n=28$ ), retroauricular crease ( $n=13$ ),  
376 posterior fornix ( $n=10$ ), and anterior nares ( $n=6$ ).

## 377 **Taxonomic profiling**

378 In comparison with the OGU method, we performed taxonomic profiling on the shotgun metagenomic  
379 data using four existing methods, specified as below. The default parameters were used for all programs.  
380 To maximize comparability, we used the WoL reference genome database (10) for all methods, except  
381 for MetaPhlAn (because it uses a special marker gene database which is difficult to customize).

- 382 1. SHOGUN: SHOGUN v1.0.8 (19), which calls Bowtie2 v2.4.1 to perform sequence alignment.
- 383 2. Bracken: Bracken v2.5 (18) on the results of Kraken v2.0.8 (48).
- 384 3. Centrifuge: Centrifuge v1.0.3 (20).
- 385 4. MetaPhlAn: MetaPhlAn v2.6.0 (21) with its database (mpa\_v20\_m200). Results (relative  
386 abundances) were normalized to counts per million sequences.

## 387 **Beta diversity analysis**

388 Beta diversity analysis of the HMP dataset was performed using QIIME 2 (1), following recommended  
389 protocols (49). Specifically, beta diversity distance matrices were constructed using the “qiime  
390 diversity beta” command with Jaccard and Bray-Curtis metrics, and using the “qiime diversity  
391 beta-phylogenetic” command (50) with unweighted UniFrac and weighted UniFrac metrics, based on  
392 the WoL reference phylogeny. Principal coordinates analysis (PCoA) was performed using the “qiime  
393 diversity pcoa” command. The correlation between biological factors (body site and host sex) and beta  
394 diversity was assessed using the PERMANOVA test, through the command “qiime diversity adonis”,  
395 with 999 permutations (the default setting).

## 396 **Site clustering by environment**

397 In the HMP study, we quantified the proximity of the three oral sites (tongue dorsum, supragingival  
398 plaque, and buccal mucosa) as compared with the four non-oral sites (stool, retroauricular crease,

399 posterior fornix, and anterior nares) as follows: For each sample in the three oral sites, we calculated the  
400 beta diversity distance to all samples in all but the current site. We then separated these distances into  
401 oral (i.e., the two oral sites other than the current one) and non-oral (i.e., the four non-oral sites). We  
402 calculated the ratio of the mean distance of the former versus the latter. Finally we reported the  
403 distribution of the mean ratios of all oral samples.

#### 404 **Phylogenetic factorization**

405 We performed phylogenetic factorization as implemented in Phylofactor v0.0.1 to infer phylogenetic  
406 clades (“factors”) that are differentially abundant between male and female subjects. Two samples with  
407 less than 100,000 OGU counts were excluded from the analysis. OGUs with relative abundance below  
408 0.01% were dropped from each sample, and OGUs present in fewer than two samples were also  
409 excluded. We built an explained variance-maximizing (the choice parameter was set to “var”)   
410 Phylofactor model using the OGU table and the WoL phylogeny. We specified the model to return 20  
411 factors. They were labeled by the taxonomic annotation of the corresponding phylogenetic clades as  
412 provided in the WoL database. The results were visualized with EMPress. In each factor, we tested the  
413 differences in male vs female subjects by comparing the ILR-transformed vectors corresponding to each  
414 sample group using a two-tailed independent samples *t*-test.

#### 415 **Subsampling of OGU tables**

416 To assess the impact of sampling depth on analysis results, we randomly subsampled the OGU tables to  
417 lower depths (sum of OGU frequencies per sample). This process mimicked lower sequencing depths in  
418 the original data, because the sum of OGU frequencies is nearly equal to the number of aligned  
419 sequences (see above). This process further considered the unaligned part of the sequencing data. For  
420 example, if  $m$  out of  $n$  sequences in a sample were aligned to at least one reference genome (therefore

421 the sum of OGU frequencies was  $m$ ), we added an extra “unaligned” feature of frequency of  $n - m$  to the  
422 OGU table, prior to random subsampling, and removed this feature after sampling.

### 423 **The FINRISK 2002 datasets**

424 The FINRISK 2002 is a large, well-phenotyped, and representative cohort based on a stratified random  
425 sample of the population aged 25 to 74 years from specific geographical areas of Finland (27). All  
426 volunteer participants took a self-administered questionnaire, physical measurements and collection of  
427 blood and stool samples. The microbiome data and metadata that support the findings of this study are  
428 available from the THL Biobank based on a written application and following relevant Finnish  
429 legislation. Details of the application process are described in the website of the Biobank:  
430 <https://thl.fi/en/web/thl-biobank/for-researchers>.

431 Paired 16S rRNA gene amplicon sequencing data and shotgun metagenomic sequencing data are  
432 available for 6,430 stool samples. The 16S rRNA data were demultiplexed, quality filtered, and denoised  
433 with deblur v1.1.0 (51), resulting in an average ASV frequency of 8,787 per sample, followed by  
434 normalization to 10,000 per sample. Taxonomic classification was performed using a pre-trained Naive  
435 Bayes classifier against the Greengenes 13\_8 database at an OTU clustering level of 99%. Feature tables  
436 were rarefied to a sampling depth of 10,000. The shotgun metagenomic data were trimmed and quality  
437 filtered using Atropos v1.1.25 (52), resulting in an average of 1.07 million paired-end sequences per  
438 sample. They were aligned to the WoL database using SHOGUN v1.0.8. An OGU table was generated  
439 using the current approach. As a comparison, Bracken v2.5 with Kraken v2.0.8 were used to infer  
440 taxonomic profiles using the same WoL database. These analyses were the same as the corresponding  
441 analyses of the HMP shotgun metagenomic dataset, as described above.

## 442 **Supervised regression for age prediction**

443 We performed machine learning analysis of microbial profiles derived from both 16S amplicon  
444 sequencing and shotgun metagenomics sequencing, at distinct levels of resolution. These included  
445 taxonomic ranks (phylum, class, order, family, genus and species) for both 16S rRNA and shotgun  
446 metagenomic data (the latter of which were inferred by either SHOGUN or Bracken), ASV for 16S  
447 rRNA data, and OGU for shotgun metagenomic data (inferred by SHOGUN with Woltka). In each  
448 profile, features with a study-wide prevalence less than 0.001 were excluded. Random Forest regressors  
449 for predicting chronological age were trained based on each profile with tuned hyperparameters with a  
450 stratified 5-fold cross-validation approach using R package ranger v0.12.1 (53). Each dataset was split  
451 into five groups with similar age distributions, and we trained the classifier on 80% of the data, and  
452 made predictions on the remaining 20% of the data in each fold iteration. We next evaluated the  
453 performance of age prediction using mean absolute error (MAE), which calculated as  $MAE = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$ ,  
454 where  $y$  denotes the predicted age,  $x$  denotes the chronological age, and  $n$  is the total number of samples.  
455 Based on the MAE evaluation, we next determined the most predictive taxonomic levels derived from  
456 both 16S and shotgun metagenomics.

457 To identify the most important taxonomic features that contributed to the age prediction, we visualized  
458 the top-128 ranked important features by built-in Random Forest importance scores and their  
459 phylogenetic relationships using EMPress (54). We next performed the feature selection analysis to  
460 identify a set of important microbial features that can maximize the model performance. We built age  
461 regressors using a series of reduced sets ( $n = 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024$ , and the number of  
462 all features) of the most predictive taxonomic features (namely, OGU) and compared their performance.  
463 The rationale is to observe the trough in MAE when additional features are added into the regression  
464 model.

465 **Statistics statement**

466 All data analysis was performed using QIIME 2 release 2020.6. PERMANOVA was performed using  
467 the “adonis” command (which wraps the “adonis” function in vegan v2.5-6). Paired *t*-test was performed  
468 using the “ttest\_rel” function in SciPy v1.4.1.

469 **Acknowledgements**

470 We are grateful to Gabriel Al-Ghalith, Zachary Burcham, Jeff DeReus, Marcus Fedarko, Shalisa  
471 Hansen, Stefan Janssen, Emily Kobayashi, Evguenia Kopylova, Tomasz Kosciolk, Holly Lutz,  
472 Cameron Martino, Siavash Mirarab, James Morton, Oriane Moyne, Wayne Pfeiffer, Daniel Roush and  
473 Se Jin Song for valuable testing of the methodology, insightful discussions on this study and additional  
474 assistance.

475 This work is supported in part by an Arizona State University start-up grant (to Q.Z.), Sloan Foundation  
476 G-2017-9838, IBM Artificial Intelligence for Healthy Living A1770534, DARPA JUMP/CRISP, NIH  
477 P30DK120515, DP1AT010885, U19AG063744, U24CA248454, Emerald Foundation Distinguished  
478 Investigator Award, Crohn’s and Colitis Foundation 675191, NSF RAPID 2038509, IBM Research AI  
479 through the AI Horizons Network and the UC San Diego Center for Microbiome Innovation (to S.H.,  
480 I.M., Y.V.-B., and R.K.). G.D.S.-P. is supported by a fellowship from the National Institutes of Health  
481 (F30 CA243480). T.N. was funded by the Emil Aaltonen Foundation, the Finnish Medical Foundation,  
482 the Finnish Foundation for Cardiovascular Disease, and the Academy of Finland (grant 321351). V.S.  
483 was supported by the Finnish Foundation for Cardiovascular Research. This work used the Comet  
484 supercomputer at the San Diego Supercomputer Center through allocation BIO150043 through the  
485 Extreme Science and Engineering Discovery Environment (XSEDE).



486 Q.Z. and R.K. conceived the project. Q.Z. led the development of the methodology and software. S.H.  
487 and Q.Z. led the analysis and interpretation of the datasets presented in this article. S.H., A.G., D.M. and  
488 Y.V.-B. contributed to the design of the method. A.G., D.M. and G.A. contributed to the development of  
489 the software. G.D.S.-P., A.D.S., P.D., F.L. contributed to the test of the method. P.B.-F., A.S.H., G.M.,  
490 T.N., L.L., V.S. and M.J. contributed to data curation. A.G., I.M., J.Y., Y.V.-B. and J.K. contributed to  
491 data analysis. N.H., G.D.S.-P., A.S.H., G.M., T.N., L.L., V.S., H.-C.K., M.J., M.I., J.A.G. and R.K.  
492 contributed to result interpretation. R.K. and Q.Z. managed the project. All the authors contributed to the  
493 composition and discussion of the manuscript.

494 We declare that we have no competing interests.

495

## 496 References

- 497 1. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H, Alm EJ,  
498 Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A, Brislawn CJ, Brown CT,  
499 Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope EK, Da Silva R, Diener C, Dorrestein PC,  
500 Douglas GM, Durall DM, Duvallet C, Edwardson CF, Ernst M, Estaki M, Fouquier J, Gauglitz JM,  
501 Gibbons SM, Gibson DL, Gonzalez A, Gorlick K, Guo J, Hillmann B, Holmes S, Holste H,  
502 Huttenhower C, Huttley GA, Janssen S, Jarmusch AK, Jiang L, Kaehler BD, Kang KB, Keefe CR,  
503 Keim P, Kelley ST, Knights D, Koester I, Kosciulek T, Kreps J, Langille MGI, Lee J, Ley R, Liu  
504 Y-X, Lofthfield E, Lozupone C, Maher M, Marotz C, Martin BD, McDonald D, McIver LJ, Melnik  
505 AV, Metcalf JL, Morgan SC, Morton JT, Naimey AT, Navas-Molina JA, Nothias LF, Orchanian  
506 SB, Pearson T, Peoples SL, Petras D, Preuss ML, Pruesse E, Rasmussen LB, Rivers A, Robeson  
507 MS 2nd, Rosenthal P, Segata N, Shaffer M, Shiffer A, Sinha R, Song SJ, Spear JR, Swafford AD,  
508 Thompson LR, Torres PJ, Trinh P, Tripathi A, Turnbaugh PJ, Ul-Hasan S, van der Hooft JJJ,  
509 Vargas F, Vázquez-Baeza Y, Vogtmann E, von Hippel M, Walters W, Wan Y, Wang M, Warren J,  
510 Weber KC, Williamson CHD, Willis AD, Xu ZZ, Zaneveld JR, Zhang Y, Zhu Q, Knight R,  
511 Caporaso JG. 2019. Reproducible, interactive, scalable and extensible microbiome data science  
512 using QIIME 2. *Nat Biotechnol* 37:852–857.
- 513 2. Breitwieser FP, Lu J, Salzberg SL. 2019. A review of methods and databases for metagenomic  
514 classification and assembly. *Brief Bioinform* 20:1125–1136.
- 515 3. Callahan BJ, McMurdie PJ, Holmes SP. 2017. Exact sequence variants should replace operational  
516 taxonomic units in marker-gene data analysis. *ISME J* 11:2639–2643.
- 517 4. Lozupone C, Knight R. 2005. UniFrac: a new phylogenetic method for comparing microbial

- 518 communities. *Appl Environ Microbiol* 71:8228–8235.
- 519 5. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, Andersen GL, Knight  
520 R, Hugenholtz P. 2012. An improved Greengenes taxonomy with explicit ranks for ecological and  
521 evolutionary analyses of bacteria and archaea. *ISME J* 6:610–618.
- 522 6. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. 2013. The  
523 SILVA ribosomal RNA gene database project: improved data processing and web-based tools.  
524 *Nucleic Acids Res* 41:D590–6.
- 525 7. Janssen S, McDonald D, Gonzalez A, Navas-Molina JA, Jiang L, Xu ZZ, Winker K, Kado DM,  
526 Orwoll E, Manary M, Mirarab S, Knight R. 2018. Phylogenetic Placement of Exact Amplicon  
527 Sequences Improves Associations with Clinical Information. *mSystems* 3.
- 528 8. O’Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B, Robbertse B,  
529 Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O, Brover V, Chetvernin  
530 V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D, Hatcher E, Hlavina W,  
531 Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, McGarvey KM, Murphy MR, O’Neill K,  
532 Pujar S, Rangwala SH, Rausch D, Riddick LD, Schoch C, Shkeda A, Storz SS, Sun H, Thibaud-  
533 Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C, Webb D, Wu W, Landrum MJ, Kimchi A,  
534 Tatusova T, DiCuccio M, Kitts P, Murphy TD, Pruitt KD. 2016. Reference sequence (RefSeq)  
535 database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids*  
536 *Res* 44:D733–45.
- 537 9. Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, Pollard KS, Sakharova E,  
538 Parks DH, Hugenholtz P, Segata N, Kyrpides NC, Finn RD. 2021. A unified catalog of 204,938  
539 reference genomes from the human gut microbiome. *Nat Biotechnol* 39:105–114.

- 540 10. Zhu Q, Mai U, Pfeiffer W, Janssen S, Asnicar F, Sanders JG, Belda-Ferre P, Al-Ghalith GA,  
541 Kopylova E, McDonald D, Kosciolk T, Yin JB, Huang S, Salam N, Jiao J-Y, Wu Z, Xu ZZ,  
542 Cantrell K, Yang Y, Sayyari E, Rabiee M, Morton JT, Podell S, Knights D, Li W-J, Huttenhower C,  
543 Segata N, Smarr L, Mirarab S, Knight R. 2019. Phylogenomics of 10,575 genomes reveals  
544 evolutionary proximity between domains Bacteria and Archaea. *Nat Commun* 10:5477.
- 545 11. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A, Hugenholtz P.  
546 2018. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree  
547 of life. *Nat Biotechnol* 36:996–1004.
- 548 12. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, Knight R,  
549 Knights D. 2018. Evaluating the Information Content of Shallow Shotgun Metagenomics.  
550 *mSystems* 3.
- 551 13. Gonzalez A, Navas-Molina JA, Kosciolk T, McDonald D, Vázquez-Baeza Y, Ackermann G,  
552 DeReus J, Janssen S, Swafford AD, Orchanian SB, Sanders JG, Shorenstein J, Holste H, Petrus S,  
553 Robbins-Pianka A, Brislawn CJ, Wang M, Rideout JR, Bolyen E, Dillon M, Caporaso JG,  
554 Dorrestein PC, Knight R. 2018. Qiita: rapid, web-enabled microbiome meta-analysis. *Nat Methods*  
555 15:796–798.
- 556 14. Poore GD, Kopylova E, Zhu Q, Carpenter C, Fraraccio S, Wandro S, Kosciolk T, Janssen S,  
557 Metcalf J, Song SJ, Kanbar J, Miller-Montgomery S, Heaton R, Mckay R, Patel SP, Swafford AD,  
558 Knight R. 2020. Microbiome analyses of blood and tissues suggest cancer diagnostic approach.  
559 *Nature* 579:567–574.
- 560 15. Gauglitz JM, Morton JT, Tripathi A, Hansen S, Gaffney M, Carpenter C, Weldon KC, Shah R,  
561 Parampil A, Fidgett AL, Swafford AD, Knight R, Dorrestein PC. 2020. Metabolome-Informed

- 562 Microbiome Analysis Refines Metadata Classifications and Reveals Unexpected Medication  
563 Transfer in Captive Cheetahs. *mSystems* 5.
- 564 16. Ha CWY, Martin A, Sepich-Poore GD, Shi B, Wang Y, Gouin K, Humphrey G, Sanders K,  
565 Ratnayake Y, Chan KSL, Hendrick G, Caldera JR, Arias C, Moskowitz JE, Ho Sui SJ, Yang S,  
566 Underhill D, Brady MJ, Knott S, Kaihara K, Steinbaugh MJ, Li H, McGovern DPB, Knight R,  
567 Fleshner P, Devkota S. 2020. Translocation of Viable Gut Microbiota to Mesenteric Adipose Drives  
568 Formation of Creeping Fat in Humans. *Cell* 183:666–683.e17.
- 569 17. Turnbaugh PJ, Qin J, D N Fredricks T L Fiedler, Costello EK, Grice EA, Ravel J, Segata N,  
570 Gillespie JJ, Sharpton TJ, Sokol H, JA. Aas, BJ. Paster, LN. Stokes, I. Olsen, FE. Dewhirst, Medini  
571 D, S K Mazmanian J L Round, Goodman AL, Kuehnert MJ, Caporaso JG, M. Kanehisa, S. Goto,  
572 M. Furumichi, M. Tanabe, M. Hidakawa, H. Li RD, Giannoukos G, MG. Langille FSB. 2012.  
573 Structure, function and diversity of the healthy human microbiome. *Nature* 486:207–214.
- 574 18. Lu J, Breitwieser FP, Thielen P, Salzberg SL. 2017. Bracken: estimating species abundance in  
575 metagenomics data. *PeerJ Comput Sci* 3:e104.
- 576 19. Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Knight R, Knights D. 2020. SHOGUN: a  
577 modular, accurate and scalable framework for microbiome quantification. *Bioinformatics* 36:4088–  
578 4090.
- 579 20. Kim D, Song L, Breitwieser FP, Salzberg SL. 2016. Centrifuge: rapid and sensitive classification of  
580 metagenomic sequences. *Genome Res* 26:1721–1729.
- 581 21. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower C,  
582 Segata N. 2015. MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nat Methods*

- 583 12:902–903.
- 584 22. Morton JT, Toran L, Edlund A, Metcalf JL, Lauber C, Knight R. 2017. Uncovering the Horseshoe  
585 Effect in Microbial Analyses. *mSystems* 2.
- 586 23. Washburne AD, Silverman JD, Leff JW, Bennett DJ, Darcy JL, Mukherjee S, Fierer N, David LA.  
587 2017. Phylogenetic factorization of compositional data yields lineage-level associations in  
588 microbiome datasets. *PeerJ* 5:e2969.
- 589 24. Cantrell K, Fedarko MW, Rahman G, McDonald D, Yang Y, Zaw T, Gonzalez A, Janssen S, Estaki  
590 M, Haiminen N, Beck KL, Zhu Q, Sayyari E, Morton JT, Armstrong G, Tripathi A, Gauglitz JM,  
591 Marotz C, Matteson NL, Martino C, Sanders JG, Carrieri AP, Song SJ, Swafford AD, Dorrestein  
592 PC, Andersen KG, Parida L, Kim H-C, Vázquez-Baeza Y, Knight R. 2021. EMPress Enables Tree-  
593 Guided, Interactive, and Exploratory Analyses of Multi-omic Data Sets. *mSystems* 6.
- 594 25. Ma B, Forney LJ, Ravel J. 2012. Vaginal microbiome: rethinking health and disease. *Annu Rev*  
595 *Microbiol* 66:371–389.
- 596 26. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD. 2015. Analysis of  
597 composition of microbiomes: a novel method for studying microbial composition. *Microb Ecol*  
598 *Health Dis* 26:27663.
- 599 27. Borodulin K, Tolonen H, Jousilahti P, Jula A, Juolevi A, Koskinen S, Kuulasmaa K, Laatikainen T,  
600 Männistö S, Peltonen M, Perola M, Puska P, Salomaa V, Sundvall J, Virtanen SM, Vartiainen E.  
601 2018. Cohort Profile: The National FINRISK Study. *Int J Epidemiol* 47:696–696i.
- 602 28. Breiman L. 2001. Random Forests. *Mach Learn* 45:5–32.
- 603 29. Zhang X, Zhong H, Li Y, Shi Z, Ren H, Zhang Z, Zhou X, Tang S, Han X, Lin Y, Yang F, Wang

- 604 D, Fang C, Fu Z, Wang L, Zhu S, Hou Y, Xu X, Yang H, Wang J, Kristiansen K, Li J, Ji L. 2021.  
605 Sex- and age-related trajectories of the adult human gut microbiota shared across populations of  
606 different ethnicities. *Nature Aging* 1:87–100.
- 607 30. Schmidt TS, Hayward MR, Coelho LP, Li SS, Costea PI, Voigt AY, Wirbel J, Maistrenko OM,  
608 Alves RJ, Bergsten E, de Beaufort C, Sobhani I, Heintz-Buschart A, Sunagawa S, Zeller G, Wilmes  
609 P, Bork P. 2019. Extensive transmission of microbes along the gastrointestinal tract. *Elife* 8.
- 610 31. Huang S, Haiminen N, Carrieri A-P, Hu R, Jiang L, Parida L, Russell B, Allaband C, Zarrinpar A,  
611 Vázquez-Baeza Y, Belda-Ferre P, Zhou H, Kim H-C, Swafford AD, Knight R, Xu ZZ. 2020.  
612 Human Skin, Oral, and Gut Microbiomes Predict Chronological Age. *mSystems* 5.
- 613 32. Asnicar F, Thomas AM, Beghini F, Mengoni C, Manara S, Manghi P, Zhu Q, Bolzan M, Cumbo F,  
614 May U, Sanders JG, Zolfo M, Kopylova E, Pasolli E, Knight R, Mirarab S, Huttenhower C, Segata  
615 N. 2020. Precise phylogenetic analysis of microbial isolates and genomes from metagenomes using  
616 PhyloPhlAn 3.0. *Nat Commun* 11:2500.
- 617 33. Jiang Y, Balaban M, Zhu Q, Mirarab S. 2021. DEPP: Deep Learning Enables Extending Species  
618 Trees using Single Genes. Cold Spring Harbor Laboratory.
- 619 34. Castelle CJ, Banfield JF. 2018. Major New Microbial Groups Expand Diversity and Alter our  
620 Understanding of the Tree of Life. *Cell* 172:1181–1197.
- 621 35. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, Gregor I, Majda S, Fiedler J,  
622 Dahms E, Bremges A, Fritz A, Garrido-Oter R, Jørgensen TS, Shapiro N, Blood PD, Gurevich A,  
623 Bai Y, Turaev D, DeMaere MZ, Chikhi R, Nagarajan N, Quince C, Meyer F, Balvočiūtė M, Hansen  
624 LH, Sørensen SJ, Chia BKH, Denis B, Froula JL, Wang Z, Egan R, Don Kang D, Cook JJ, Deltel

- 625 C, Beckstette M, Lemaitre C, Peterlongo P, Rizk G, Lavenier D, Wu Y-W, Singer SW, Jain C,  
626 Strous M, Klingenberg H, Meinicke P, Barton MD, Lingner T, Lin H-H, Liao Y-C, Silva GGZ,  
627 Cuevas DA, Edwards RA, Saha S, Piro VC, Renard BY, Pop M, Klenk H-P, Göker M, Kyrpides  
628 NC, Woyke T, Vorholt JA, Schulze-Lefert P, Rubin EM, Darling AE, Rattei T, McHardy AC. 2017.  
629 Critical Assessment of Metagenome Interpretation-a benchmark of metagenomics software. *Nat*  
630 *Methods* 14:1063–1071.
- 631 36. Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. 2020. A complete  
632 domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol* 38:1079–1086.
- 633 37. Ye SH, Siddle KJ, Park DJ, Sabeti PC. 2019. Benchmarking Metagenomics Tools for Taxonomic  
634 Classification. *Cell* 178:779–794.
- 635 38. Sanders JG, Nurk S, Salido RA, Minich J, Xu ZZ, Zhu Q, Martino C, Fedarko M, Arthur TD, Chen  
636 F, Boland BS, Humphrey GC, Brennan C, Sanders K, Gaffney J, Jepsen K, Khosroheidari M, Green  
637 C, Liyanage M, Dang JW, Phelan VV, Quinn RA, Bankevich A, Chang JT, Rana TM, Conrad DJ,  
638 Sandborn WJ, Smarr L, Dorrestein PC, Pevzner PA, Knight R. 2019. Optimizing sequencing  
639 protocols for leaderboard metagenomics by combining long and short reads. *Genome Biol* 20:226.
- 640 39. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9:357–  
641 359.
- 642 40. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,  
643 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and  
644 SAMtools. *Bioinformatics* 25:2078–2079.
- 645 41. Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.



- 646        Bioinformatics 25:1754–1760.
- 647    42. Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34:3094–3100.
- 648    43. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.  
649        *Journal of Molecular Biology*.
- 650    44. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. 2016. VSEARCH: a versatile open source tool  
651        for metagenomics. *PeerJ* 4:e2584.
- 652    45. Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using DIAMOND. *Nat*  
653        *Methods* 12:59–60.
- 654    46. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. 2019. GTDB-Tk: a toolkit to classify genomes  
655        with the Genome Taxonomy Database. *Bioinformatics*  
656        <https://doi.org/10.1093/bioinformatics/btz848>.
- 657    47. Federhen S. 2011. The NCBI Taxonomy database. *Nucleic Acids Res* 40:D136–D143.
- 658    48. Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2. *Genome Biol*  
659        20:257.
- 660    49. Estaki M, Jiang L, Bokulich NA, McDonald D, González A, Kosciółek T, Martino C, Zhu Q,  
661        Birmingham A, Vázquez-Baeza Y, Dillon MR, Bolyen E, Gregory Caporaso J, Knight R. 2020.  
662        QIIME 2 Enables Comprehensive End-to-End Analysis of Diverse Microbiome Data and  
663        Comparative Studies with Publicly Available Data. *Current Protocols in Bioinformatics*.
- 664    50. McDonald D, Vázquez-Baeza Y, Koslicki D, McClelland J, Reeve N, Xu Z, Gonzalez A, Knight R.  
665        2018. Striped UniFrac: enabling microbiome analysis at unprecedented scale. *Nat Methods* 15:847–

666 848.

- 667 51. Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley EP,  
668 Thompson LR, Hyde ER, Gonzalez A, Knight R. 2017. Deblur Rapidly Resolves Single-Nucleotide  
669 Community Sequence Patterns. *mSystems* 2.
- 670 52. Didion JP, Martin M, Collins FS. 2017. Atropos: specific, sensitive, and speedy trimming of  
671 sequencing reads. *PeerJ* 5:e3720.
- 672 53. Wright MN, Ziegler A. 2017. ranger: A Fast Implementation of Random Forests for High  
673 Dimensional Data in C++ and R. *J Stat Softw* 77:1–17.
- 674 54. Cantrell K, Fedarko MW, Rahman G, McDonald D, Yang Y, Zaw T, Gonzalez A, Janssen S, Estaki  
675 M, Haiminen N, Beck KL, Zhu Q, Sayyari E, Morton J, Tripathi A, Gauglitz JM, Marotz C,  
676 Matteson NL, Martino C, Sanders JG, Carrieri AP, Song SJ, Swafford AD, Dorrestein PC,  
677 Andersen KG, Parida L, Kim H-C, Vázquez-Baeza Y, Knight R. 2020. EMPress enables tree-  
678 guided, interactive, and exploratory analyses of multi-omic datasets. Cold Spring Harbor  
679 Laboratory.
- 680

681 Figure Legends

682 **Figure 1. Feature resolution impacts community structure analysis even in small conceptual**  
683 **examples. A.** A synthetic dataset involving three microbial communities, each of which having 12  
684 unique read hits, as represented by black circles in the frequency table, to a total of 10 reference  
685 genomes (OGUs), classified under five species, three genera and one family, as noted to the left. A  
686 phylogenetic tree of the 10 genomes is shown on the right. In this simplified case, the phylogeny is not  
687 much more complex than the taxonomy (with three more edges); however, the taxonomic assignment  
688 and the phylogenetic placement of genome O5 are not consistent. **B.** Beta diversity of the dataset. The  
689 three samples (circles) are connected by edges representing the pairwise distances calculated by Bray-  
690 Curtis (BC) or weighted UniFrac (WU) on the frequency table. For the latter measure, either the  
691 taxonomy or the phylogeny was used to quantify the hierarchical relationships among OGUs, as noted in  
692 the parentheses. The edge lengths were normalized so that their sum is equal in each graph. This  
693 synthetic case study demonstrates that different resolutions of features and feature structures can lead to  
694 very different conclusions regarding sample relationships.

695

696 **Figure 2. Analysis of the HMP metagenomes reveals clustering by body environment and**  
697 **differentiation by host sex.** Beta diversity analysis was performed on 210 samples subsampled to one  
698 million paired-end shotgun reads each. **A.** PCoA by the method proposed in this study (OGU): weighted  
699 UniFrac metric calculated with the WoL reference phylogeny based on the OGU table. Samples (dots)  
700 are colored by body site and shaped by host sex. **B.** PCoA using the current adopted method (CAM):  
701 Bray-Curtis calculated on species-level taxonomic units identified by Bracken, which shows a diagonal  
702 pattern that aligns all samples of the four non-oral body sites in one plane (also see Figs. S2B and S3).  
703 **C.** Proportions of community structure variance explained by the first three axes of PCoA. **D.** Mean

704 ratio of the beta diversity distances from any oral sample to a sample of the two other oral sites versus to  
705 that of non-oral body sites. The lower the mean ratio is, the more similar communities of the three oral  
706 sites are to each other in the background of multiple body environments. The bold line in each box  
707 represents the median. The whiskers represent 1.5 IQR. **E** and **F**. PERMANOVA pseudo- $F$  statistics  
708 indicating the differentiation of community structures by body site (**E**) and by host sex (**F**). The larger  $F$   
709 is, the more distinct the community structures are between groups versus within groups. The y-axis is  
710 aligned to  $F=1.0$  which indicates no difference. For **E**, all statistics have a  $p$ -value of 0.001. For **F**, an  
711 asterisk (\*) indicates  $p$ -value  $\leq 0.05$ . **G**. Differentially abundant phylogenetic clades by host sex inferred  
712 using PhyloFactor and visualized using EMPress on the WoL reference phylogeny. The tree was  
713 subsetted to only include OGU detected in the dataset. The top 20 clades by effect size are colored (full  
714 details provided in Figs. S4-5). The top five clades are numbered 1 through 5 by decreasing effect size,  
715 circled, and labeled with corresponding taxonomic annotations. The small color ring represents phylum-  
716 level annotations. The inner and outer barplot rings indicate the OGU counts split by body site (using the  
717 same color scheme as in A and B) and by host sex, respectively.

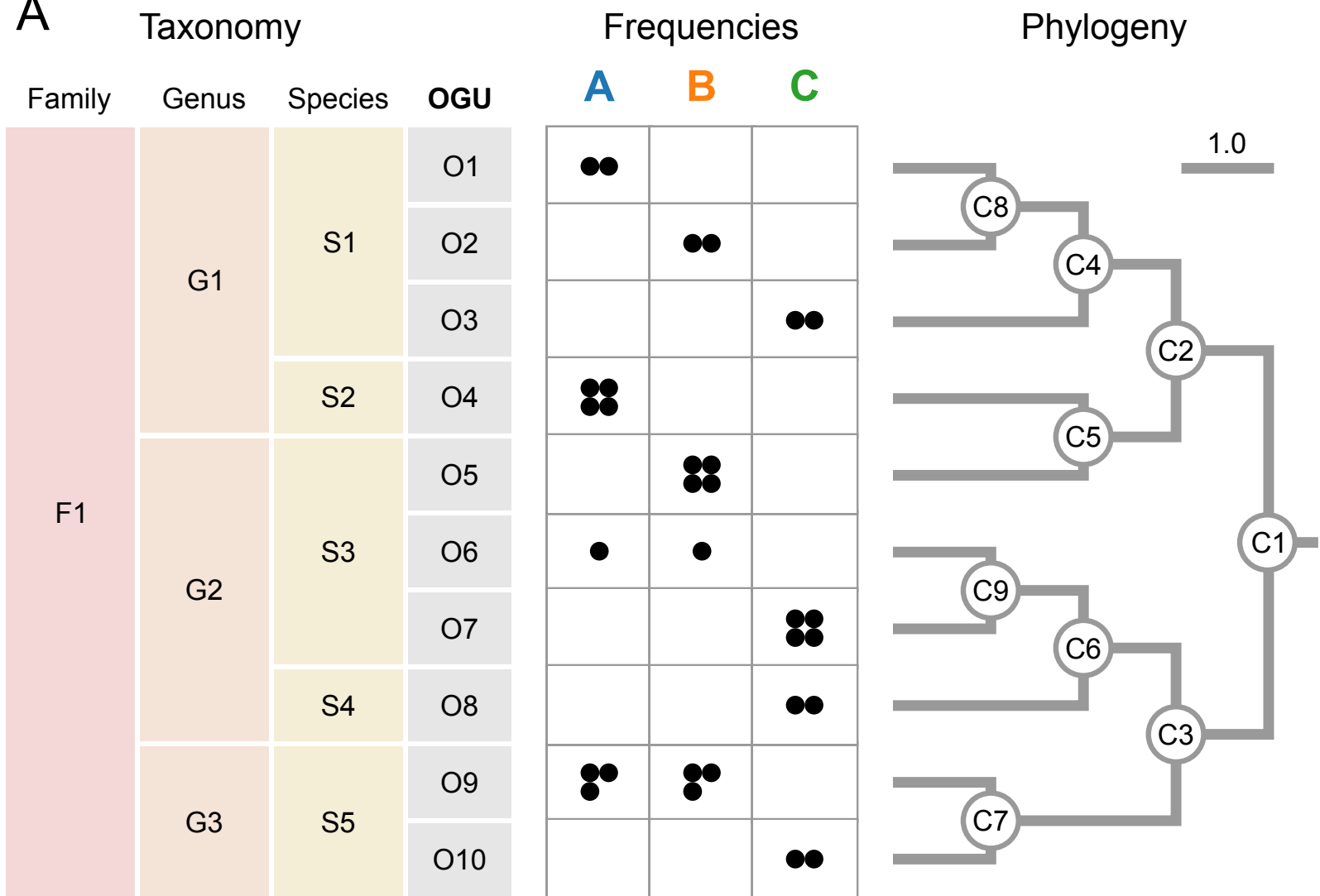
718

719 **Figure 3. Analysis of the FINRISK metagenomes showing superior prediction accuracy over**  
720 **taxonomic units and 16S rRNA data. A.** The empirical error (mean absolute error, MAE) in  
721 predicting host chronological age using microbiome features at distinct taxonomic ranks in paired 16S  
722 rRNA amplicon and shotgun metagenomics data with a Random Forests regressor. “None” represents  
723 the taxonomy-free, finest-possible level (ASV for 16S, OGU for shotgun). Small circles indicate MAEs  
724 in all iterations of five-fold cross validation. Large circles and error bars indicate the mean and standard  
725 deviations of the five MAEs. **B.** Scatter plot of the actual age vs. the predicted age by the best-  
726 performing model with OGU features in the five-fold cross-validation. The black line was generated  
727 using ggplot2’s local polynomial regression fitting. **C.** Phylogenomic tree of 169 OGUs with importance

728 score  $\geq 0.1$  in the prediction model. The tree was subsampled based on the WoL reference phylogeny,  
729 and drawn to scale (branch lengths represent mutations per site). Branch colors indicate the mean  
730 importance score of all descendants of the clade. Taxonomic labels are displayed where needed. Circles  
731 and lines with stops are displayed where needed to assist location of taxonomic labels to target branches  
732 or clades.

733

734

**A****B Beta diversity**