



Efficient computation of Faith's phylogenetic diversity with applications in characterizing microbiomes

George W Armstrong, Kalen Cantrell, Shi Huang, et al.

Genome Res. published online September 3, 2021

Access the most recent version at doi:[10.1101/gr.275777.121](https://doi.org/10.1101/gr.275777.121)

P<P	Published online September 3, 2021 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

A promotional banner for Cellecta's CRISPR and RNAi Genetic Screening. The text reads "CRISPR and RNAi Genetic Screening. Your new superpower." To the right is a "LEARN MORE" button and a photograph of a person in a red superhero mask and cape. The Cellecta logo, a green molecular structure, is in the bottom right corner.

CRISPR and RNAi Genetic Screening.
Your new superpower.

LEARN MORE

CELLECTA

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Efficient computation of Faith's phylogenetic diversity with applications in characterizing**
2 **microbiomes**

3 George Armstrong^{1,2,3}, Kalen Cantrell², Shi Huang^{1,2}, Daniel McDonald¹, Niina Haiminen⁴,
4 Anna Paola Carrieri⁵, Qiyun Zhu^{6,7}, Antonio Gonzalez¹, Imran McGrath^{2,8}, Kristen L. Beck⁹,
5 Daniel Hakim^{1,3}, Aki S. Havulinna^{10,11}, Guillaume Méric^{12,13}, Teemu Niiranen^{10,14,15}, Leo Lahti¹⁶,
6 Veikko Salomaa¹⁰, Mohit Jain^{2,17,18}, Michael Inouye^{12,19}, Austin D. Swafford², Ho-Cheol Kim⁹,
7 Laxmi Parida⁴, Yoshiki Vázquez-Baeza², Rob Knight^{1,2,20,21,#}

8 ¹Department of Pediatrics, School of Medicine, University of California, San Diego, California,
9 USA; ²Center for Microbiome Innovation, Jacobs School of Engineering, University of
10 California San Diego, La Jolla, California, USA; ³Bioinformatics and Systems Biology Program,
11 University of California, San Diego, California, USA; ⁴IBM T. J. Watson Research Center,
12 Yorktown Heights, New York, USA; ⁵IBM Research Europe, The Hartree Centre, Warrington,
13 United Kingdom; ⁶School of Life Sciences, Arizona State University, Tempe, Arizona, USA;
14 ⁷Biodesign Center for Fundamental and Applied Microbiomics, Arizona State University,
15 Tempe, Arizona, USA; ⁸Division of Biological Sciences, University of California San Diego, La
16 Jolla, California, USA; ⁹IBM Almaden Research Center, San Jose, California, USA;
17 ¹⁰Department of Public Health and Welfare, Finnish Institute for Health and Welfare, Helsinki,
18 Finland; ¹¹Institute for Molecular Medicine Finland (FIMM), HiLIFE, University of Helsinki,
19 Helsinki, Finland; ¹²Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes
20 Institute, Melbourne, Victoria, Australia; ¹³Department of Infectious Diseases, Central Clinical
21 School, Monash University, Melbourne, Victoria, Australia; ¹⁴Department of Internal Medicine,
22 University of Turku, Turku, Finland; ¹⁵Division of Medicine, Turku University Hospital,
23 Finland; ¹⁶Department of Computing, University of Turku, Turku, Finland; ¹⁷Department of
24 Medicine, University of California, San Diego, California, USA; ¹⁸Department of Pharmacology,
25 University of California, San Diego, California, USA; ¹⁹Department of Public Health and
26 Primary Care, Cambridge University, Cambridge, UK; ²⁰Department of Computer Science and
27 Engineering, University of California, San Diego, La Jolla, CA, USA; ²¹Department of
28 Bioengineering, University of California, San Diego, La Jolla, CA, USA;

29 # Corresponding Author: robknight@ucsd.edu

30 Abstract

31 **The number of publicly available microbiome samples is continually growing. As dataset**
32 **size increases, bottlenecks arise in standard analytical pipelines. Faith's phylogenetic**
33 **diversity is a highly utilized phylogenetic alpha diversity metric that has thus far failed to**
34 **effectively scale to trees with millions of vertices. Stacked Faith's Phylogenetic Diversity**
35 **(SFPhD) enables calculation of this widely adopted diversity metric at a much larger scale**
36 **by implementing a computationally efficient algorithm. The algorithm reduces the amount**
37 **of computational resources required, resulting in more accessible software with a reduced**
38 **carbon footprint, as compared to previous approaches. The new algorithm produces**
39 **identical results to the previous method. We further demonstrate that the phylogenetic**
40 **aspect of Faith's PD provides increased power in detecting diversity differences between**
41 **younger and older populations in the FINRISK study's metagenomic data.**

42 Introduction

43 In microbiome research, particular attention is given to evaluating the diversity of microbes
44 within samples (McDonald et al. 2018a; The Human Microbiome Project Consortium 2012;
45 Thompson et al. 2017). Alpha diversity (within sample diversity) represents a family of summary
46 statistics that can summarize the breadth of diversity present in an environment. More recently,
47 many examples have been reported on the associations between various host or environmental
48 factors and alpha diversity of microbiomes, including country and diet in human guts (McDonald
49 et al. 2018a), disease status in humans and canines (Gevers et al. 2014; Vázquez-Baeza et al.
50 2016), the pH (Lauber et al. 2009), salinity (Thompson et al. 2017) and temperature (Zhou et al.
51 2016) of soils, among many others (Youngblut et al. 2019; Jeffery et al. 2016). A popular metric
52 that accounts for the phylogenetic relatedness of the community members, Faith's Phylogenetic

53 Diversity (Faith's PD) (Faith 1992), has been noted to be more sensitive in distinguishing disease
54 factors in the human digestive system, relative to other alpha diversity indices (Scherson and
55 Faith 2018; Youngblut et al. 2021).

56 Modern DNA sequencing instruments have enabled microbiome studies at the scale of
57 tens of thousands of samples, which presents a computational challenge for metrics that rely on a
58 phylogeny, such as Faith's PD. This metric is computed by summing the branch lengths (edge
59 weights) of the phylogeny that exclusively represents the sequences contained in a biological
60 sample. The amount of memory and number of necessary operations needed to calculate Faith's
61 PD depends on the number of edges in the phylogenetic tree, as well as the number of samples in
62 the underlying data table.

63 In today's increasingly large and sparse datasets and meta-analyses, these phylogenetic
64 trees and tables can exceed 100,000s of samples and millions of tree tips (McDonald et al.
65 2018b). Recent advances have enabled efficient computation of the UniFrac metric for beta
66 diversity. UniFrac is also a metric computed over phylogenetic trees (Lozupone and Knight
67 2005) and is mathematically related to Faith's PD (Faith et al. 2009). Specifically, Striped
68 UniFrac (McDonald et al. 2018b) improves upon previous UniFrac implementations (Hamady et
69 al. 2010) by using space- and time-efficient tree data structures (Cordova and Navarro 2016) and
70 reducing the number of vectors required to store intermediate scores in the tree.

71 Additionally, the usefulness of techniques like Faith's PD and UniFrac remains
72 underexplored for metagenomics sequencing. Recent molecular protocol optimizations, such as
73 SHOGUN (Hillmann et al. 2018), have enabled the metagenomic characterization of large
74 human cohorts (Borodulin et al. 2015; Kaplan et al. 2019; Salosensaari et al. 2021). In this
75 context, the applicability of Faith's PD has largely been limited by the technical difficulties
76 associated with constructing phylogenies from metagenomic features (Zhu et al. 2019). Efforts

77 like the Web of Life (WoL) (Zhu et al. 2019) and Genome Taxonomy Database (GTDB) (Parks
78 et al. 2018, 2020) are now addressing this issue by providing a phylogenomic tree as part of their
79 database releases that can be used for phylogeny-informed analysis.

80 Motivated by these advances in algorithms and resources for analyzing phylogenies,
81 phylogenomic trees, and sparse data, we developed a new algorithm and implementation,
82 Stacked Faith's Phylogenetic Diversity (SFPhD), for rapidly computing Faith's PD. Additionally,
83 we aim to demonstrate concrete benefits of phylogeny-informed analysis in metagenomic studies
84 where this metric is less frequently used.

85 Results

86 SFPhD is a new implementation for calculating Faith's PD. The key advances of SFPhD are
87 using a sparse matrix representation, an efficient tree structure, and partial aggregation of metric
88 constituents. Our BSD-licensed implementation of this algorithm is available in the `unifrac`
89 package (via PyPI and bioconda (Grüning et al. 2018)), which has 57,007 total conda downloads
90 and 40,434 conda downloads since the introduction of SFPhD, as of the time of writing (August
91 28, 2021). The package produces a C/C++ shared library with Python bindings and is
92 additionally linkable by any programming language (<https://github.com/biocore/unifrac>).
93 Additionally, by investigating the previously documented relationship between age and bacterial
94 richness of the gut microbiome (de la Cuesta-Zuluaga et al. 2019), we demonstrate that
95 accounting for phylogeny in metagenomic data can increase the statistical power for detecting
96 group differences (**Supplemental Code**).

97 **Stacked Faith's PD provides a faster and memory-efficient implementation over the**
98 **previous state-of-the-art algorithm.**

99 SFPhD uses the structure of microbiome data along with other practical considerations to
 100 achieve decreased time and memory requirements. An example feature table is shown in **Fig.**
 101 **1A**, with a corresponding phylogenetic tree in **Fig. 1B**. Note that for a given tree \mathcal{T} Faith's PD
 102 can be expressed as

$$PD_{\mathbf{i}} = \sum_{\mathbf{j} \in \mathcal{T}} I_{\mathbf{ij}} \times \text{branchLen}_{\mathbf{j}}(\mathcal{T})$$

103
 104 where $PD_{\mathbf{i}}$ is Faith's PD for sample \mathbf{i} , $I_{\mathbf{ij}}$ indicates if sample \mathbf{i} has any features that descend from
 105 node \mathbf{j} , and $\text{branchLen}_{\mathbf{j}}(\mathcal{T})$ indicates the length of the branch to node \mathbf{j} in the tree \mathcal{T} .

106 The previous state-of-the-art reference implementation (scikit-bio, <http://scikit-bio.org/>)
 107 computes Faith's PD for a batch of samples by first fully computing $I_{\mathbf{ij}}$. $I_{\mathbf{ij}}$ is computed by
 108 traversing the entire phylogenetic tree in a post-order traversal, where the children of a node
 109 must be visited before the node itself can be visited (the nodes in **Fig. 1B** are labeled in the order
 110 of a post-order traversal). During the traversal, when a given node \mathbf{j} is visited, all \mathbf{i} are set by
 111 determining the features present in all children of node \mathbf{j} . Subsequently, the
 112 $I_{\mathbf{ij}} \times \text{branchLen}_{\mathbf{j}}(\mathcal{T})$ for all branches is calculated. The results are obtained by summing over
 113 the branches for each sample (**Fig. 1C**). However, this approach tends to use much more space
 114 than is needed.

115 Microbiome data are known to be sparse (Martino et al. 2019; Kumar et al. 2018; Morton
 116 et al. 2017), i.e., of the entries in a data table, many are likely to be zero. This issue is
 117 exacerbated in large datasets, where many microbes are only observed in a handful of samples.
 118 In an extreme case, such a table (McDonald et al. 2018b) with 113,721 samples rarefied at 500
 119 sequences per sample, has only 0.0126% non-zero entries. Sparse representations have been used
 120 previously for storing microbiome data (McDonald et al. 2012a), and have been applied for
 121 accelerating microbiome analyses (McDonald et al. 2018b), but they have not been previously
 122 applied to Faith's PD. We identified that a major downfall of the state-of-the-art implementation

123 in scikit-bio is that it uses a full, dense table to represent all of I_{ij} in memory at once. A key
124 advancement of our approach is to use a sparse matrix implementation for storing information on
125 the taxa present for each sample and feature. Sparse matrices save space by only retaining
126 information about positions in the matrix that have non-zero values (e.g., only the grey values in
127 **Fig. 1A** and information about their positions are retained by a sparse matrix).

128 Another key advance is the partial aggregation of Faith's PD (**Fig. 1D**). Note that the
129 $I_{ij} \times \text{branchLen}_j(\mathcal{T})$, which we will call a metric constituent, can be added in any order, and
130 that I_{ij} only depends on the children of node j . Thus, if node k is a child of node j , I_{ik} is no
131 longer needed once metric constituents for node k have been computed and I_{ij} is known. As a
132 result, we can reduce the memory used to store I_{ij} by traversing the phylogeny with a post-order
133 traversal and freeing I_{ik} after they are no longer needed. Furthermore, we can reduce the storage
134 needed for the metric constituents keeping a running summation of them while traversing the
135 tree. Thus, this approach reduces the expected space complexity for storing the metrics from
136 $O(nk)$, to $O(n \log(k))$, where n is the number of samples and k is the number of vertices in the
137 tree.

138 In addition to the algorithmic improvements, we have included several practical
139 enhancements that improve the performance of the code. The topology of the phylogenetic tree
140 (**Fig. 1B**) is now represented as balanced-parentheses vector (**Fig. 1E**) that corresponds to
141 additional vectors of branch lengths and node names; this structure has a lower memory footprint
142 and a sequential memory representation which reduces the number of cache misses during a tree
143 traversal (Cordova and Navarro 2016). Finally, the software is written using C/C++ (with Python
144 extensions using Cython, <https://cython.org/>) and builds upon the foundation established by
145 Striped UniFrac (McDonald et al. 2018b). Reuse of this library facilitated our access to a much
146 faster Newick format parser, which reduces the overhead when reading a tree from disk. These

147 factors make for an improved expected and in-practice performance, despite the time complexity
148 and worst-case memory complexity remaining the same.

149 To demonstrate the scalability of SFPhD, we used a collection of 307,237 public and
150 anonymized private 16S rRNA V4 microbiome samples amounting to 1,264,796 phylogenetic
151 tree tips (after rarefaction at 500 sequences per sample). The samples were retrieved using the
152 redbiom command line interface (McDonald et al. 2019) which queried a cache of public and
153 anonymized private studies available in Qiita (Gonzalez et al. 2018). Amplicon sequence
154 variants (ASVs) were placed into the Greengenes (Gonzalez et al. 2018; McDonald et al. 2012b;
155 DeSantis et al. 2006) phylogeny using SEPP (Mirarab et al. 2012). Computing the full alpha
156 diversity vector took SFPhD 1 hour and 5 minutes wall-clock time and required a maximum
157 resident set size of less than 3 GB (see Methods for hardware details). In addition, we iteratively
158 measured runtime and memory consumption for increasingly large random subsets of samples
159 while fixing the size of the tree at 100,000 tips (**Fig. 2A, 2B, Table S1**). For the iteration with
160 20,000 samples, the memory usage of the reference implementation exceeded 150 GB and the
161 process ran for over 15 minutes. Contrastingly, with SFPhD, the process took 14 seconds to
162 execute and required less than 0.5 GB of memory. Additionally, using Green Algorithms
163 (Lannelongue et al. 2021), we estimated the carbon footprint of the scikit-bio reference
164 implementation on the 20,000 sample table to be 12.84 g CO₂e, whereas we estimated the carbon
165 footprint of SFPhD would be 0.04 g CO₂e in the United States, which is a 321-fold reduction in
166 impact on global warming.

167 **Phylogenetic diversity is a suitable metric to analyze stool metagenomic samples**

168 To demonstrate SFPhD's versatility and applicability to newer datasets, we reanalyzed
169 2,661 paired 16S rRNA and metagenomic data of stool samples from the FINRISK (Borodulin et
170 al. 2015, 2018; Salosensaari et al. 2021) study (n=1,563 aged 60 and older, n=1,098 aged 35 and

171 under). In this experiment, we select random subsets of the full sample set and compare each
172 metric's (Observed Features and Faith's PD) ability to detect differences in mean alpha diversity
173 distributions. For each step we randomly select N paired 16S and metagenomic samples, and
174 then compute the difference in mean alpha diversity between samples taken from younger adults
175 (under 35 years) and older adults (over 60 years) together with an empirical p-value. For both
176 16S and metagenomics, the alpha diversity of younger adults is lower than in older adults. In
177 metagenomics, but not in 16S sequencing, Faith's PD provides improved statistical power over
178 observed features, a phylogenetically-agnostic alternative (**Fig. 3A, B**). With 16S data, the
179 difference between the two metrics is subtle (**Fig. 3A**). In both cases, the statistical power
180 increases as the number of samples grows. With metagenomic data, the number of observed
181 features shows a weaker effect compared to Faith's PD regardless of the number of samples (**Fig.**
182 **3B**). Unlike 16S datasets (5,600 features), metagenomic datasets (1,700 features) are resolution-
183 limited by the reference databases. Whereas the nature of amplicon sequence variants (ASVs)
184 allows for a broader feature space that can capture age-differences without the need for a
185 phylogeny.

186 We investigated the difference in mean alpha diversity in metagenomic samples (**Fig.**
187 **4A**) by computing the log of the likelihood ratio of *older to younger adult* samples present for
188 each branch in the WoL phylogenomic tree (Zhu et al. 2019). We were able to identify portions
189 of the WoL tree responsible for the increase in phylogenetic diversity (**Fig. 4B**). From this
190 analysis, we found that the majority of the tree is comparably represented in young and old adult
191 samples. However, we also found two clades where *older adult* samples were more prevalent
192 than *younger adult* samples (Clade 1 has a log likelihood ratio bounded with an 80% confidence
193 interval of [1.20, 1.45] and Clade 2 has an 80% confidence interval of [0.55, 0.74]). Clade 1
194 corresponds to a majority of Lactobacillales genomes, and Clade 2 corresponds to Proteobacteria

195 genomes. The branches in Clade 1 primarily have a large log likelihood ratio, indicating that the
196 features across the entire clade are more likely to be found in samples from older adults.
197 However, the internal branches in Clade 2 additionally have low log likelihood ratios, indicating
198 that the enrichment of features in older adults is not completely consistent across the entire clade.
199 Lastly, although not confined to a few clades, there are several tips (e.g. *Staphylococcus aureus*,
200 *Bavariicoccus seileri*, *Nitratireductor indicus*, and *Campylobacter ureolyticus*) in the phylogeny
201 that are only associated with younger adults.

202 **Discussion**

203 By accounting for the relationship between features in a dataset, Faith's PD can mitigate issues
204 with sparsity and heterogeneity common to modern 'omics' datasets. Although this metric was
205 first introduced 30 years ago, the underlying algorithm for computing this metric had largely
206 remained unchanged. In this paper we demonstrated that our novel algorithm, SFPhD, performed
207 efficiently on datasets with hundreds of thousands of samples and millions of tree tips, producing
208 identical results to those of previous algorithms for computing this metric while producing a
209 speedup of up to 64× and requiring as little as 0.21% of the memory in our benchmarks.

210 An important aspect of SFPhD's underlying algorithm is substituting calculation of the
211 full presence/absence table over the phylogeny, for a tree traversal that partially aggregates
212 diversity values and frees presence/absence information when no longer needed. The result is a
213 high-performance implementation that demonstrates improved scaling with the number of
214 samples in the input dataset. Much of the engineering work here was facilitated by the balanced
215 parenthesis tree implementation provided in the UniFrac package (McDonald et al. 2018b).
216 Therefore, we believe that increasing the availability of efficient and flexible data structures for
217 phylogenetic analyses is likely to accelerate and facilitate the development of novel analytical

218 methods. In a broader sense, this is similar to the impact of NumPy's (McDonald et al. 2018b;
219 Harris et al. 2020) N-dimensional array in image processing, machine learning, neuroscience,
220 and other fields.

221 In addition, in a stool metagenomic study Faith's PD demonstrates increased statistical
222 power compared to observed features for differentiating younger from older subjects based on
223 their microbial communities. In this context, we show that Faith's PD consistently provided
224 increased statistical power for determining age-based differences in the shotgun metagenomic
225 sequencing data. While this metric was originally developed to analyze data with vastly different
226 statistical and biological properties, its use here demonstrates the versatility and applicability
227 behind measuring diversity using a tree. Furthermore, enabling efficient Faith's PD computation
228 on microbiome datasets is of particular importance when examining the impact of COVID-19 on
229 gut health (Kim et al. 2021).

230 Although we show the utility of SFPhD in large and complex microbiome studies, the
231 underlying implementation is not tied to a particular molecular technology. Thus, this
232 implementation will be relevant to fields outside of microbiology, such as conservation
233 prioritization, which inspired the original version of Faith's PD (Faith 1992) and where it
234 continues to be applied (Rosauer et al. 2017). We also envision our implementation will be
235 applicable in fields like nutrition and metabolomics research, that only recently began adopting
236 trees for analytical tasks (Tripathi et al. 2021; Johnson et al. 2019).

237 **Methods**

238 **Construction of benchmarking tables**

239 Data for the benchmarking in this study were subsampled from a BIOM table of 113,721 and
240 761,003 ASVs, which is composed of studies aggregated from several large sources of publicly
241 available microbiome data in Qiita (Amir et al. 2017; Gonzalez et al. 2018). This data table was

242 produced as previously described (McDonald et al. 2018b). The data was subset by uniformly
243 randomly sampling the desired number of ASVs and samples from the table. Ten different tables
244 were created for each number of samples and ASVs. The published insertion tree (McDonald et
245 al. 2018b) was collapsed to only contain sequences that were selected to be included in the given
246 subsampled table.

247 The table with 307,237 public and anonymized private 16S rRNA V4 microbiome
248 samples and 1,264,796 phylogenetic tree tips was also prepared as previously described
249 (McDonald et al. 2018b), but included samples with private sequencing data from Qiita.

250 **Benchmarking time and memory estimates**

251 The SFPhD implementation available in the Python package unifrac v0.10.0 was used. The
252 reference implementation uses the Faith's PD implementation from scikit-bio v0.5.4.

253 All methods were run single-threaded on shared compute nodes that were not running
254 other compute tasks. The nodes all had Intel(R) Xeon(R) CPU E5-2640 v3 @ 2.60GHz
255 processors. A job was terminated if it exceeded 6 hours of wall time or 250 GB of memory
256 (system max). Space was tracked using GNU Time. Time for both implementations was tracked
257 with a Python wrapper script. The time needed to parse data is not included in the scikit-bio
258 timings, but is included in the SFPhD timings, due to the lack of access to this information in the
259 unifrac interface. This is acceptable given that it results in a conservative estimate of the speedup
260 with SFPhD.

261 **Carbon footprint estimation**

262 The Green Algorithms interface (Lannelongue et al. 2021) was used to estimate the Carbon
263 Dioxide equivalent (CO₂e) of the benchmarked methods. The Intel(R) Xeon(R) CPU E5-2640 v3
264 CPUs used in benchmarking have a Thermal Design Power (TDP) per core of the 11.25 TDP /
265 core.

266 **FINRISK processing**

267 The 16S rRNA data were demultiplexed, quality filtered, and denoised with deblur (Amir et al.
268 2017). The Greengenes (McDonald et al. 2012b) 13.8 with a clustering level of 99% was used as
269 the reference phylogeny for open-reference feature picking with SEPP (Mirarab et al.2012).
270 ASVs with a total frequency fewer than 10 were discarded, and the table was then rarefied to a
271 sampling depth of 1000 reads/sample. The resulting table and insertion tree were used for
272 calculation of Faith's PD.

273 The shotgun metagenomic data were trimmed and quality filtered using Atropos (Didion
274 et al.2017). They were aligned to the WoL database using SHOGUN pipeline (v1.0.8) with a
275 Bowtie 2 alignment option. A table was generated from the alignments using the OGU workflow
276 (Zhu et al. 2021). OGUs with a total frequency fewer than 10 were discarded, and the table was
277 then rarefied to a sampling depth of 1000 reads/sample. The WoL phylogenomic tree (Zhu et al.
278 2021, 2019) was used for Faith's PD.

279 Both tables were filtered to include only samples from individuals 35 and younger
280 (younger criteria) or 60 and older (older criteria).

281 **Power estimation for mean difference in alpha diversity**

282 For a given N (shown on horizontal axis in **Fig. 3a,b**), the FINRISK processed samples matching
283 the younger/older criteria were sampled to this depth. On the subsampled data, the difference in
284 mean alpha diversity between younger and older adults \bar{d} , was computed. A null distribution, \hat{D} ,
285 was generated by repeating 1000 repetitions of shuffling the age category associated with an
286 alpha diversity and recomputing the difference of mean alpha diversity between the groups. The
287 p-value was computed by finding the percentile of \bar{d} in \hat{D} .

288 This test procedure was repeated for 1000 repetitions. The power for N is estimated as
 289 the proportion of tests found significant at $\alpha = 0.05$.

290 Older-younger log likelihood-ratio calculation

291 The WoL tree (Zhu et al. 2019) was pruned and filtered to only include the OGU's (Zhu et al.
 292 2021) belonging to the FINRISK samples with age ≤ 35 and ≥ 60 . For each node $t \in \mathcal{T}$ in the
 293 tree,

$$294 \quad \logLikelihoodRatio_t = \log \left(\frac{|Samples_{older}(Descendants(t))|}{|Samples_{younger}(Descendants(t))|} \right) - \log \left(\frac{|Samples_{older}(\mathcal{T})|}{|Samples_{younger}(\mathcal{T})|} \right)$$

295

296 where $Descendants(t)$ is the set of descendants of t in \mathcal{T} , and for a set of nodes \mathcal{N} ,

297 $Samples_{group}(\mathcal{N})$ is the set of samples that contain any features in \mathcal{N} .

298 Phylogenetic Visualization

299 Tree was visualized using EMPress (Cantrell et al. 2021). A node in the tree was considered old
 300 if its $age_{log} > 0$ and young if its $age_{log} < 0$.

301 Software Availability

302 The data used for benchmarking Faith's PD timing and memory usage are available as per the
 303 Striped UniFrac paper (McDonald et al. 2018b). The code for the benchmarking is available on
 304 GitHub (<https://github.com/biocore/faiths-pd-benchmarking>). The data and code needed for
 305 benchmarking the FINRISK metagenomics data are also available on GitHub. The SFPhD code
 306 is available in the unifrac Python package (<https://github.com/biocore/unifrac>). All of the
 307 software is also available in the **Supplemental Code**.

308 Competing Interests

309 We declare that we have no competing interests.

310 **Acknowledgements**

311 This work was supported in part by IBM Research AI through the AI Horizons Network, the
312 Center for Microbiome Innovation at UC San Diego, the Academy of Finland grant 321351 and
313 the Emil Aaltonen Foundation (to T.N.), the National Institutes of Health grant R01ES027595 (to
314 M.J.), the Academy of Finland grants 321356 and 335525 (A.S.H), the Academy of Finland
315 grant 295741 (L.L.). MI was supported by the Munz Chair of Cardiovascular Prediction and
316 Prevention. VS was supported by the Finnish Foundation for Cardiovascular Research.

317 **Figure Legends**

318 **Figure 1. Partially aggregating branch lengths reduces the space complexity of the**
319 **algorithm.** a) Faith's PD calculation depends on the representation of features present in
320 samples. In the table, the letters (R, O, B, K) represent samples and the numbers (0, 1, 2, 4, 6, 9,
321 10) represent features. A 1 in an entry indicates the presence of a feature in the sample. SFPhD
322 uses sparse table data structures, which reduce memory by only keeping track of the non-zero
323 values in a matrix (highlighted in gray). b) A mock reference phylogenetic tree is shown, with
324 the features from (a) as tips. Labels for the samples from (a) are located next to tips that they
325 contain. The nodes are labeled by their order in a post-order traversal of the tree. c) Graphic
326 depiction of the reference implementation's calculation of Faith's PD by first aggregating the
327 presence/absence information for each branch in the tree, followed by multiplication by the
328 branch lengths to get the metric constituents, and finally a sum over the entire branch \times metric
329 constituent table. d) Graphic representation of the execution of SFPhD. On the left, the stack of
330 presence/absence information is shown at three points during the algorithm's execution (i, ii, iii).
331 Each of these times shows the stack immediately before memory is freed. On the right, the state
332 of the partially aggregated phylogenetic diversity (PD) is shown after each node is added to the

333 stack. Each row represents the vector after a step in the algorithm. In practice, there is only one
334 such vector. e) The balanced parentheses representation for the phylogenetic tree from (b).

335 **Figure 2. SFPhD outperforms the reference implementation in terms of runtime and**
336 **memory usage.** a) Runtime in seconds for computing Faith's PD on datasets with thousands of
337 samples and 100,000 tips in the phylogeny. Data is independently sub-sampled from a collection
338 of 113,721 public samples in Qiita (Zhu et al. 2019; Gonzalez et al. 2018) as previously
339 processed (McDonald, Vázquez-Baeza, et al. 2018). Mean of n=10 repetitions with 95% CI error
340 bars. b) Memory usage for the same experiment as in (a). For both a and b jobs were terminated
341 if they exceeded 250 GB of memory.

342 **Figure 3. Phylogenetic diversity provides increased statistical power to differentiate age**
343 **groups in shotgun metagenomics but not in 16S rRNA sequencing.** a) Statistical power to
344 differentiate young adults from old adults in two alpha diversity metrics at different sample sizes
345 using 16S rRNA sequencing in the FINRISK cohort. b) Same as (a) but for shallow shotgun
346 metagenomic sequencing.

347 **Figure 4. Phylogenetic tree colored by age-group log of the likelihood ratio of older to**
348 **younger adults per node.** a) Distribution of Faith's PD by age group on the full dataset. b) Web
349 of Life (WoL) Phylogenetic tree with branches colored by the log of likelihood ratio of old adults
350 compared to young adults in descendants of the branch, for the FINRISK dataset. The inner
351 circle is colored by the log of likelihood ratio of older adults compared to younger adults in the
352 tips of the tree. The outer circle is colored by the phylum of the taxon represented by each tree
353 tip. Red ellipses mark two clades enriched for samples from older individuals.

354 **References**

- 355 Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Xu ZZ, Kightley EP,
356 Thompson LR, Hyde ER, Gonzalez A, Knight R. 2017. Deblur Rapidly Resolves Single-
357 Nucleotide Community Sequence Patterns. *mSystems* **2** (2).
358 <https://doi.org/10.1128/mSystems.00191-16>.
- 359 Borodulin K, Tolonen H, Jousilahti P, Jula A, Juolevi A, Koskinen S, Kuulasmaa K, Laatikainen
360 T, Mannisto S, Peltonen M et al. 2018. Cohort Profile: The National FINRISK Study. *Int J*
361 *Epidemiol.* **47** (3): 696-696i. doi: 10.1093/ije/dyx239.
- 362 Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T, Männistö S,
363 Salomaa V, Sundvall J, Puska P. 2015. Forty-Year Trends in Cardiovascular Risk Factors in
364 Finland. *Eur J Public Health.* **25** (3): 539–46. doi: 10.1093/eurpub/cku174.
- 365 Cantrell K, Fedarko MW, Rahman G, McDonald D, Yang Y, Zaw T, Gonzalez A, et al. 2021.
366 EMPress Enables Tree-Guided, Interactive, and Exploratory Analyses of Multi-Omic Data
367 Sets. *mSystems* **6** (2). doi: 10.1128/mSystems.01216-20.
- 368 Cordova J, and Navarro G. 2016. Simple and Efficient Fully-Functional Succinct Trees. *Theor.*
369 doi: 10.1016/j.tcs.2016.04.031.
- 370 de la Cuesta-Zuluaga J, Kelley ST, Chen Y, Escobar JS, Mueller NT, Ley RE, McDonald D,
371 Huang S, Swafford AD, Knight R, Thackray VG. 2019. Age- and Sex-Dependent Patterns
372 of Gut Microbial Diversity in Human Adults. *mSystems* **4** (4). doi:
373 10.1128/mSystems.00261-19.
- 374 DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber T, Dalevi D, Hu P,
375 Andersen GL. 2006. Greengenes, a Chimera-Checked 16S rRNA Gene Database and
376 Workbench Compatible with ARB. *Appl Environ Microbiol.* **72** (7): 5069–72. doi:
377 10.1128/AEM.03006-05.

- 378 Didion JP, Martin M, Collins FS. 2017. Atropos: Specific, Sensitive, and Speedy Trimming of
379 Sequencing Reads. *PeerJ*. **5**: e3720. doi: 10.7717/peerj.3720.
- 380 Faith DP. 1992. Conservation Evaluation and Phylogenetic Diversity. *Biol Conserv* 61 (1): 1–10.
- 381 Faith DP, Lozupone CA, Nipperess D, and Knight R. 2009. The Cladistic Basis for the
382 Phylogenetic Diversity (PD) Measure Links Evolutionary Features to Environmental
383 Gradients and Supports Broad Applications of Microbial Ecology’s ‘Phylogenetic Beta
384 Diversity’ Framework. *Int J Mol Sci* **10** (11). doi: 10.3390/ijms10114723.
- 385 Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, Schwager E,
386 Knights D, Song SJ, Yassour M et al. 2014. The Treatment-Naive Microbiome in New-
387 Onset Crohn’s Disease. *Cell Host Microbe* **15** (3): 382–92. doi:
388 10.1016/j.chom.2014.02.005.
- 389 Gonzalez A, Navas-Molina JA, Kosciolk T, McDonald D, Vázquez-Baeza Y, Ackermann G,
390 DeReus J, Janssen S, Swafford AD, Orchanian SB et al. 2018. Qiita: Rapid, Web-Enabled
391 Microbiome Meta-Analysis. *Nat Methods* **15** (10): 796–98. doi: 10.1038/s41592-018-0141-
392 9.
- 393 Grüning B, The Bioconda Team, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH,
394 Valieris R, Köster J. 2018. Bioconda: Sustainable and Comprehensive Software Distribution
395 for the Life Sciences. *Nat Methods*. doi: 10.1038/s41592-018-0046-7.
- 396 Hamady M, Lozupone C, Knight R. 2010. Fast UniFrac: Facilitating High-Throughput
397 Phylogenetic Analyses of Microbial Communities Including Analysis of Pyrosequencing
398 and PhyloChip Data. *ISME J* **4** (1): 17–27. doi: 10.1038/ismej.2009.97.
- 399 Harris CR., Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E,
400 Taylor J, Berg S, Smith NJ et al. 2020. Array Programming with NumPy. *Nature* **585**
401 (7825): 357–62. doi: 10.1038/s41586-020-2649-2.

- 402 Hillmann B, Al-Ghalith GA, Shields-Cutler RR, Zhu Q, Gohl DM, Beckman KB, Knight R,
403 Knights D. 2018. Evaluating the Information Content of Shallow Shotgun Metagenomics.
404 *mSystems* **3** (6). doi: 10.1128/mSystems.00069-18.
- 405 Jeffery IB., Lynch DB, O’Toole PW. 2016. Composition and Temporal Stability of the Gut
406 Microbiota in Older Persons.” *ISME J* **10** (1): 170–82.
- 407 Johnson AJ., Vangay P, Al-Ghalith GA, Hillmann BM, Ward TL, Shields-Cutler RR, Kim AD,
408 Shmagel AK, Syed AN, Personalized Microbiome Class Students et al. 2019. Daily
409 Sampling Reveals Personalized Diet-Microbiome Associations in Humans. *Cell Host*
410 *Microbe* **25** (6): 789–802.e5.
- 411 Kaplan RC., Wang Z, Usyk M, Sotres-Alvarez D, Daviglius ML, Schneiderman N, Talavera GA,
412 Gellman MD, Thyagarajan B, Moon JY et al. 2019. Gut Microbiome Composition in the
413 Hispanic Community Health Study/Study of Latinos Is Shaped by Geographic Relocation,
414 Environmental Factors, and Obesity. *Genome Biology* **20**. doi: 10.1186/s13059-019-1831-z.
- 415 Kim HN , Joo EJ, Lee CW, Ahn KS, Kim HL, Park DI, Park SK. 2021. Reversion of Gut
416 Microbiota during the Recovery Phase in Patients with Asymptomatic or Mild COVID-
417 19: Longitudinal Study. *Microorganisms* **9** (6): 1237. doi:
418 10.3390/microorganisms9061237.
- 419 Kumar MS, Slud EV, Okrah K, Hicks SC, Hannenhalli S, Bravo HC. 2018. Analysis and
420 Correction of Compositional Bias in Sparse Sequencing Count Data. *BMC Genomics* **19** (1):
421 799.
- 422 Lannelongue L, Grealey J, Inouye M. 2021. Green Algorithms: Quantifying the Carbon
423 Footprint of Computation. *Adv Sci* **2100707**. doi: 10.1002/advs.202100707.
- 424 Lauber CL, Hamady M, Knight R, Fierer N. 2009. Pyrosequencing-Based Assessment of Soil pH
425 as a Predictor of Soil Bacterial Community Structure at the Continental Scale. *Appl Environ*

- 426 *Microbiol.* doi: 10.1128/aem.00335-09.
- 427 Lozupone C, Knight R. 2005. UniFrac: A New Phylogenetic Method for Comparing Microbial
428 Communities. *Appl Environ Microbiol* **71** (12): 8228–35.
- 429 Martino C, Morton JT, Marotz CA, Thompson LR, Tripathi A, Knight R, Zengler K. 2019. A
430 Novel Sparse Compositional Technique Reveals Microbial Perturbations. *mSystems* **4** (1).
431 doi: 10.1128/mSystems.00016-19.
- 432 McDonald D, Clemente JC, Kuczynski J, Rideout JR, Stombaugh J, Wendel D, Wilke A, Huse
433 S, Hufnagle J, Meyer F et al. 2012a. The Biological Observation Matrix (BIOM) Format or:
434 How I Learned to Stop Worrying and Love the Ome-Ome. *GigaScience* **1** (1): 7. doi:
435 10.1186/2047-217X-1-7.
- 436 McDonald D, Hyde E, Debelius JW, Morton JT, Gonzalez A, Ackermann G, Aksenov AA,
437 Behsaz B, Brennan C, Chen Y et al. 2018a. American Gut: An Open Platform for Citizen
438 Science Microbiome Research. *mSystems* **3** (3). doi: 10.1128/mSystems.00031-18.
- 439 McDonald D, Kaehler B, Gonzalez A, DeReus J, Ackermann G, Marotz C, Huttley G, Knight R.
440 2019. Redbiom: A Rapid Sample Discovery and Feature Characterization System. *mSystems*
441 **4** (4). doi: 10.1128/mSystems.00215-19.
- 442 McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, Andersen GL,
443 Knight R, Hugenholtz P. 2012b. An Improved Greengenes Taxonomy with Explicit Ranks
444 for Ecological and Evolutionary Analyses of Bacteria and Archaea. *ISME J* **6** (3): 610–18.
- 445 McDonald D, Vázquez-Baeza Y, Koslicki D, McClelland J, Reeve N, Xu Z, Gonzalez A, Knight
446 R. 2018b. Striped UniFrac: Enabling Microbiome Analysis at Unprecedented Scale. *Nat*
447 *Methods* **15** (11): 847–48.
- 448 Mirarab S, Nguyen N, Warnow T. 2012. SEPP: SATé-Enabled Phylogenetic Placement. *Pac*
449 *Symp*, 247–58. doi: 10.1142/9789814366496_0024.

- 450 Morton JT, Toran L, Edlund A, Metcalf JL, Lauber C, Knight R. 2017. Uncovering the
451 Horseshoe Effect in Microbial Analyses. *mSystems* **2** (1). doi:10.1128/mSystems.00166-16.
- 452 Parks DH, Chuvochina M, Chaumeil PA, Rinke C, Mussig AJ, Hugenholtz P. 2020. A Complete
453 Domain-to-Species Taxonomy for Bacteria and Archaea. *Nat Biotechnology* **38** (9): 1079–
454 86.
- 455 Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, Hugenholtz P.
456 2018. A Standardized Bacterial Taxonomy Based on Genome Phylogeny Substantially
457 Revises the Tree of Life. *Nat Biotechnology* **36** (10): 996–1004.
- 458 Rosauer DF, Pollock LJ, Linke S, Jetz W. 2017. Phylogenetically Informed Spatial Planning Is
459 Required to Conserve the Mammalian Tree of Life. *Proc R Soc. B* **284** (20170627).
460 doi:10.1098/rspb.2017.0627.
- 461 Salosensaari A, Laitinen V, Havulinna AS, Meric G, Cheng S, Perola M, Valsta L, Alfthan G,
462 Inouye M, Watrous JD et al. 2021. Taxonomic Signatures of Cause-Specific Mortality Risk
463 in Human Gut Microbiome. *Nat Communications* **12** (1): 2671.
- 464 Scherson RA, Faith DP. 2018. In *Phylogenetic Diversity: Applications and Challenges in*
465 *Biodiversity Science*. Springer International Publishing, New York, New York.
- 466 The Human Microbiome Project Consortium. 2012. Structure, Function and Diversity of the
467 Healthy Human Microbiome. *Nature*. doi:10.1038/nature11234.
- 468 Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, Prill RJ, Tripathi A,
469 Gibbons SM, Ackerman G et al. 2017. A Communal Catalogue Reveals Earth’s Multiscale
470 Microbial Diversity. *Nature* **551** (7681): 457–63.
- 471 Tripathi A, Vázquez-Baeza Y, Gauglitz JM, Wang M, Dührkop K, Nothias-Esposito M,
472 Acharya DD, Ernst M, van der Hoof JJJ, Zhu Q et al. 2021. Chemically Informed Analyses
473 of Metabolomics Mass Spectrometry Data with Qemistree. *Nat Chem Biol* **17** (2): 146–51.

- 474 Vázquez-Baeza Y, Hyde ER, Suchodolski JS, Knight R. 2016. Dog and Human Inflammatory
475 Bowel Disease Rely on Overlapping yet Distinct Dysbiosis Networks. *Nat Microbiol*. doi:
476 10.1038/nmicrobiol.2016.177.
- 477 Yatsunenkov T, Rey FE, Manary MJ, Trehan I, Dominguez-Bello MG, Contreras M, Magris M,
478 Hidalgo G, Baldassano RN, Anokhin AP et al. 2012. Human Gut Microbiome Viewed across
479 Age and Geography. *Nature* **486** (7402): 222–27.
- 480 Youngblut ND, de la Cuesta-Zuluaga J, Ley RE. 2021. Incorporating Genome-Based Phylogeny
481 and Functional Similarity into Diversity Assessments Helps to Resolve a Global Collection
482 of Human Gut Metagenomes. *bioRxiv*. doi: 10.1101/2020.07.16.207845.
- 483 Youngblut ND, Reischer GH, Walters W, Schuster N, Walzer C, Stalder G, Ley RE, Farnleitner
484 AH. 2019. Host Diet and Evolutionary History Explain Different Aspects of Gut
485 Microbiome Diversity among Vertebrate Clades. *Nat Commun* **10** (1): 2200.
- 486 Zhou J, Deng Y, Shen L, Wen C, Yan Q, Ning D, Qin Y, Xue K, Wu L, He Z et al. 2016.
487 Temperature Mediates Continental-Scale Diversity of Microbes in Forest Soils. *Nat*
488 *Commun* **7** (July): 12083.
- 489 Zhu Q, Huang Q, Gonzalez A, McGrath I, McDonald D, Haiminen N, Armstrong G, Vazquez-
490 Baeza Y, Yu J, Kuczynski J et al. 2021. OGU's Enable Effective, Phylogeny-Aware Analysis
491 of Even Shallow Metagenome Community Structures. *bioRxiv*. doi:
492 10.1101/2021.04.04.438427.
- 493 Zhu Q, Mai U, Pfeiffer W, Janssen S, Asnicar F, Sanders JG, Belda-Ferre P, Al-Ghalith GA,
494 Kopylova E, McDonald D et al. 2019. Phylogenomics of 10,575 genomes reveals
495 evolutionary proximity between domains Bacteria and Archaea. *Nat Commun* **10** (1): 5477.

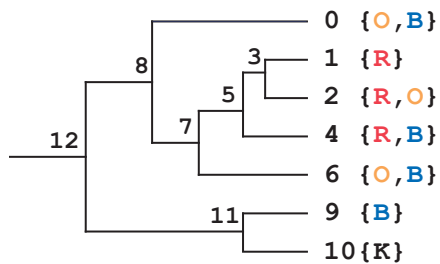
A

Feature Table

	0	1	2	4	6	9	10
R	0	1	1	1	0	0	0
O	1	0	1	0	1	0	0
B	1	0	0	1	1	1	0
K	0	0	0	0	0	0	1

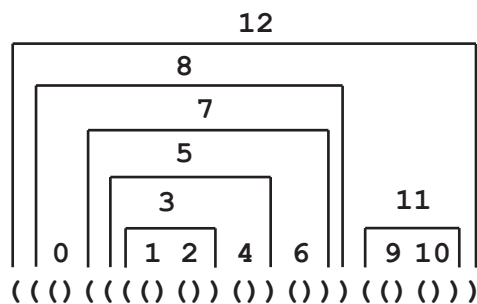
B

Phylogenetic Tree



E

Balanced Parentheses



C

Reference Implementation

Branch	R	O	B	K	Length	R*L	O*L	B*L	K*L
0	0	1	1	0	0.7	0	0.7	0.7	0
1	1	0	0	0	0.2	0.2	0	0	0
2	1	1	0	0	0.2	0.2	0.2	0	0
4	1	0	1	0	0.3	0.3	0	0.3	0
6	0	1	1	0	0.5	0	0.5	0.5	0
9	0	0	1	0	0.3	0	0	0.3	0
10	0	0	0	1	0.3	0	0	0	0.3
3	1	1	0	0	0.1	0.1	0.1	0	0
5	1	1	1	0	0.2	0.2	0.2	0.2	0
7	1	1	1	0	0.2	0.2	0.2	0.2	0
8	1	1	1	0	0.3	0.3	0.3	0.3	0
11	0	0	1	1	0.7	0	0	0.7	0.7
12	1	1	1	1	0.3	0.3	0.3	0.3	0.3

Final Faith's PD = sum = 1.8 2.5 3.5 1.3

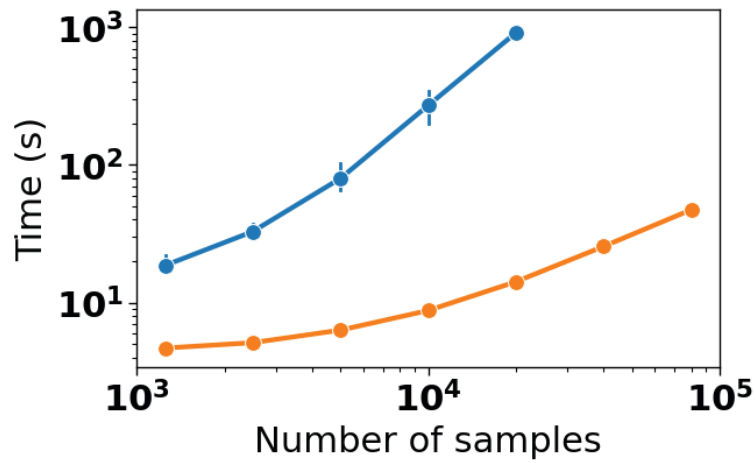
D

Stacked Faith's PD Implementation

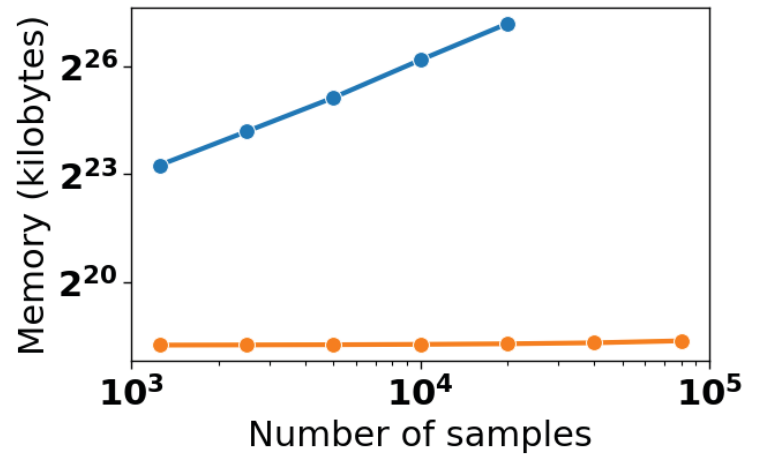
Branch	R	O	B	K	Length	Partially Aggregated PD			
0	0	1	1	0	0.7	0.0	0.7	0.7	0.0
i. 1	1	0	0	0	0.2	0.2	0.7	0.7	0.0
2	1	1	0	0	0.2	0.4	0.9	0.7	0.0
0	0	1	1	0	0.7	0.5	1.0	0.7	0.0
ii. 3	1	1	0	0	0.1	0.8	1.0	1.0	0.0
4	1	0	1	0	0.3	1.0	1.2	1.2	0.0
0	0	1	1	0	0.7	1.0	1.7	1.7	0.0
iii. 5	1	1	1	0	0.2				
6	0	1	1	0	0.5				

Final Faith's PD = 1.8 2.5 3.5 1.3

A



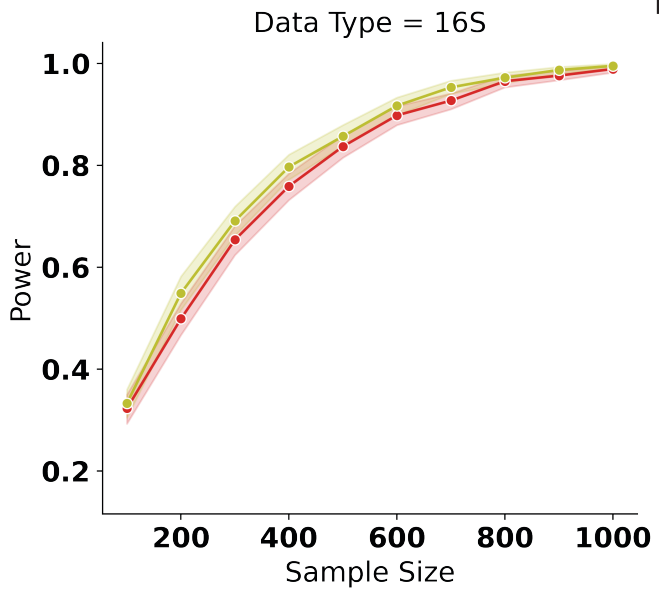
B



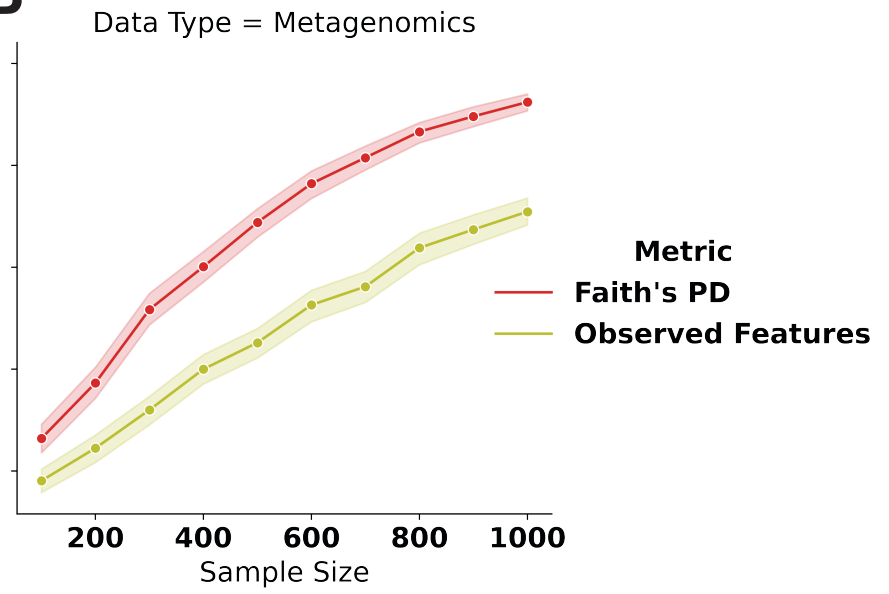
 **Stacked Faith's PD**

 **Reference Implementation**

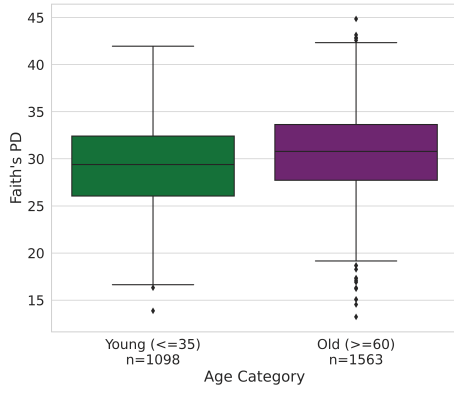
A



B



A



B

