

1 **Authors:** Daniel McDonald, Yoshiki Vazquez Baeza, David Koslicki, Jason McClelland, Nicolai  
2 Reeve, Zhenjiang Xu, Antonio Gonzalez, Rob Knight

3  
4 **Title:** Striped UniFrac: enabling microbiome analysis at unprecedented scale.

5  
6 **Abstract:** UniFrac is synonymous with microbiome research, yet it no longer scales to large  
7 datasets. We propose a new algorithm, Striped UniFrac, which produces identical results,  
8 reduces space and time complexities by >10x and exhibits near linear parallel scaling. We  
9 highlight it by computing UniFrac on 113,721 samples in 48 hours using 256 CPUs. A BSD-  
10 licensed implementation is available that produces a C shared library linkable by any  
11 programming language.

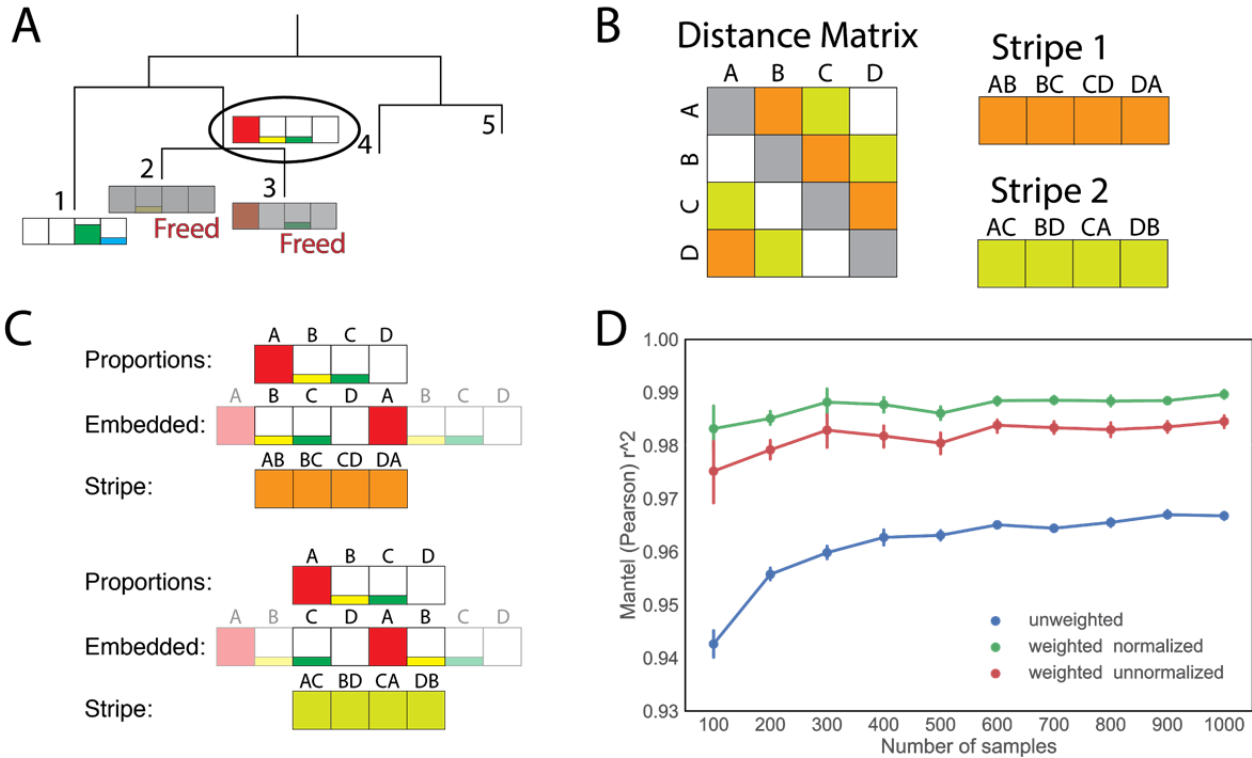
12  
13 **[main text: presently 1389 words]**

14  
15 The study of the microbiome has rapidly expanded, largely because of the insight afforded by  
16 UniFrac <sup>1</sup>. UniFrac is a phylogenetic measure of beta-diversity that assesses differences  
17 between pairs of microbiome profiles, and provides the method underlying some of the field's  
18 most iconic insights <sup>2-4</sup>. UniFrac is central to microbial community studies because it accounts  
19 for evolutionary relationships between microbes present within a sample, whereas other  
20 distance metrics such as Euclidean distance, Bray-Curtis, and Jaccard make the unrealistic  
21 implicit assumption that all organisms are equally related (see <sup>5</sup> for more detail), leading to  
22 statistical artifacts on the resulting sparse data matrices.

23  
24 Microbiome studies have recently transitioned from experimental designs with a few dozen  
25 samples to designs spanning tens of thousands of samples. Large-scale studies, such as the  
26 Earth Microbiome Project <sup>6</sup>, afford the statistical crucial for untangling the many factors that  
27 influence microbial community composition. Unfortunately, these large studies make prohibitive  
28 demands on our algorithms and data structures. In the summer of 2015, we set out to apply  
29 UniFrac to the Earth Microbiome Project, a dataset spanning over 25,000 samples. We  
30 discovered that Fast UniFrac <sup>7</sup>, irrespective of implementation, could not process this dataset  
31 even with months running on specialized hardware. Here we describe a novel algorithm with a  
32 shared library implementation usable by any programming language that can process the EMP  
33 dataset on a laptop in <24 hours.

34  
35 Five important advances were made in the design of Striped UniFrac. First was to reduce the  
36 average case space requirement from by performing a postorder aggregation of the sample  
37 proportions to compute UniFrac at each internal node (fig. 1A; proof in supplemental text),  
38 similar to EMDUniFrac <sup>8</sup>, avoiding the necessity of a dense matrix representation of sample  
39 proportions at all vertices. The second advance was to orient the pairwise sample comparisons  
40 along diagonals, or stripes, of the resulting distance matrix. This transform allows for substantial  
41 vectorization through single instruction multiple data (SIMD) operations, because groups of  
42 pairwise distances can be computed in a single instruction (fig. 1BC). Memory locality is also  
43 preserved by representing proportions and stripes as contiguous C-style arrays. The third  
44 advance was to allow independent execution of distance matrix stripes allowing considerable

45 task-level parallelism (either by threads or processes). Fourth, in a bifurcating tree,  
 46 approximately 50% of the vertices are represented by the tips; we reasoned that collapsing the  
 47 phylogeny slightly by disregarding UniFrac computations at the tips would yield a highly  
 48 correlated result (e.g., similar to the minimal differences caused by moving from 99% to 97%  
 49 OTUs) (fig. 1D). This heuristic is optional. Last, we represent the phylogeny using balanced  
 50 parentheses<sup>9</sup>, a succinct data structure that supports rapid tree reductions and traversals, and  
 51 has minimal memory requirements.  
 52



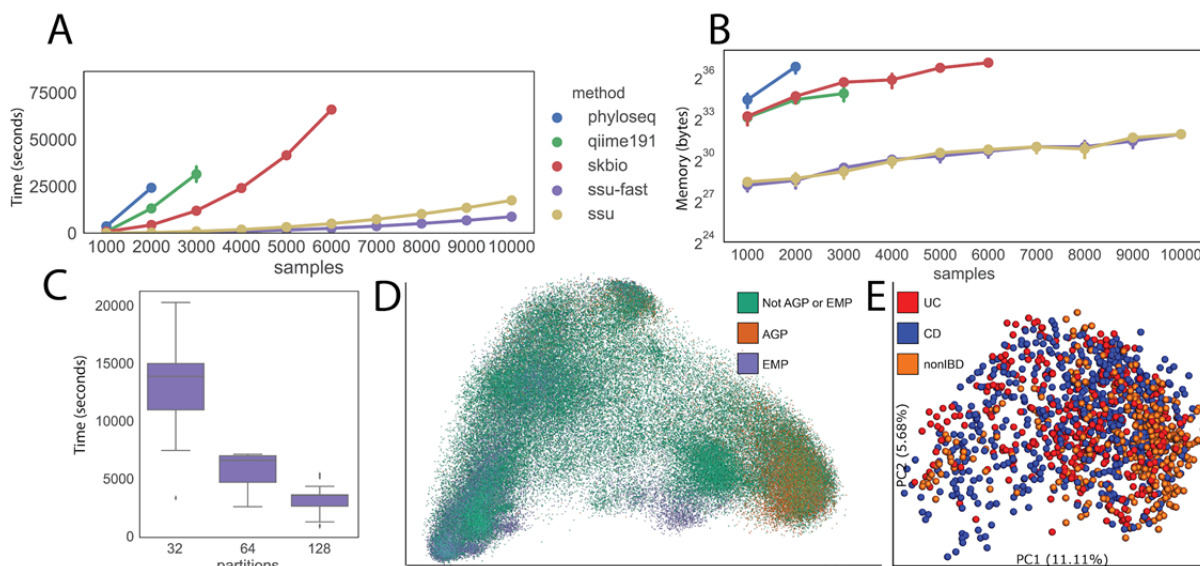
53 **Figure 1.** Algorithm highlights. **(A)** A depiction of the fourth vertex evaluated in a postorder  
 54 traversal and the resulting sample proportions (circled). This vertex is the parent of tips “2” and  
 55 “3”. The sample proportions for this vertex represent the aggregate (sum) of the sample  
 56 proportions of the features which descend. The memory for features “2” and “3” is no longer  
 57 needed, and can be freed. **(B)** A schematic of the two stripes in a four sample logical distance  
 58 matrix; the labels above the stripes denote the pairwise comparison represented (e.g., “AB”  
 59 indicates the distance between samples “A” and “B”). **(C)** A nodes sample proportions are  
 60 embedded, duplicating the proportion information. This duplication allows the sample  
 61 proportions to be slid along the embedded proportions, allowing comparison of linear blocks of  
 62 memory for all pairwise combinations of samples. **(D)** Mantel tests (Pearson) between Strided  
 63 State UniFrac in exact mode, which produces identical results to UniFrac vs. fast mode in which  
 64 the UniFrac distances are not computed at the tips of the tree during traversal. Each data point  
 65 represents 10 random subsets of the Earth Microbiome Project Deblur 90nt dataset, with the  
 66 median R<sup>2</sup> value depicted. Error bars are 95% CI.  
 67  
 68

69 In order to benchmark this new algorithm, we randomly sampled the Earth Microbiome Project  
 70 Deblur 90 nucleotide dataset at increasing numbers of samples with ten iterations at a given

71 number of samples. For each sOTU table, the EMP phylogeny was sheared to only the features  
 72 contained to avoid benchmarking different phylogenetic data structures. For each table and tree  
 73 pair, we computed unweighted UniFrac, weighted normalized UniFrac and weighted  
 74 unnormalized UniFrac with Striped UniFrac, Fast UniFrac as used by QIIME1<sup>10</sup> (the reference  
 75 implementation as implemented in PyCogent<sup>11</sup>), Fast UniFrac in QIIME2 (a Cython'd  
 76 implementation in scikit-bio, and Fast UniFrac in phyloseq (a independent R implementation)<sup>12</sup>.  
 77 For each execution, runtime and space were tracked using GNU Time (fig. 2AB). To explore  
 78 parallelism, we executed Striped UniFrac on the full EMP table (25,145 samples) using shared  
 79 computational nodes with each process using two threads, and tracked the per process time  
 80 and memory (fig. 2C) for walltime distributions, (fig. S1) for memory distributions.

81  
 82 We then obtained 120,790 Illumina 16S rRNA V4 samples processed by Deblur<sup>13</sup> from Qiita  
 83 (Qiita study IDs, sample counts and study titles in table S1). The 5,522,523 fragments were  
 84 inserted into Greengenes<sup>14</sup> 13\_8 99% tree using SEPP<sup>15,16</sup>. Rarefaction to 500 sequences per  
 85 sample reduced the total number of samples to 113,721. By distributing the computation over  
 86 256 processors, we computed unweighted UniFrac in under 48 hours walltime, using 7,977  
 87 CPU hours with a peak resident memory during the distributed computation at 1.3GB. We then  
 88 performed Principal Coordinates Analysis using centered FSVD<sup>17</sup> on the resulting distances,  
 89 and visualized them through EMPeror<sup>18</sup> (fig. 2D). This level of integration shows the dramatic  
 90 range of diversity associated with the American Gut Project (in press), which includes samples  
 91 from only a single host type (humans) suggesting a high degree of “unique” microbiomes remain  
 92 to be found by deeper sampling in other host types.

93



94  
 95 **Figure 2.** Space and time complexities. (A-B) Time and space comparisons between phyloseq,  
 96 QIIME v1.9.1, scikit-bio, Striped UniFrac in exact mode, and Striped UniFrac in fast mode. Each  
 97 datapoint represents 10 random subsets of the Earth Microbiome Project 90nt dataset at  
 98 increasing numbers of samples. All methods were run single threaded on non-shared compute  
 99 nodes which were not running other compute tasks. A job was killed if it exceeded 24 hours  
 100 walltime or 256GB of memory (system max). (C) Walltime distributions of independent

101 processes operating on the full Earth Microbiome Project dataset (over 25,000 samples)  
102 executing on shared compute nodes. An individual partition represents a single independent  
103 process, and each process was run with two threads; 32 partitions indicates 32 processes using  
104 two threads each. A higher partition count means each individual process is doing less work.  
105 Given sufficient available resources, the maximum for a distribution represents a near upper  
106 bound on walltime. (D) An ordination plot of unweighted UniFrac distances over 113,721  
107 samples sourced from Qiita. (E) Unweighted UniFrac applied to the metagenomic data from the  
108 integrated Human Microbiome Project data using the taxonomy as the tree. UC is Ulcerative  
109 Colitis, CD is Crohn's Disease and nonIBD is non-Inflammatory Bowel Disease.

110  
111 For extensibility, we implemented the algorithm in C++ under the BSD open source license, and  
112 provide a shared C library with examples of interfacing to it using C and R, and a  
113 comprehensive Python using the C-API (<https://github.com/biocore/unifrac>). In addition, we  
114 implemented kernels to support other variants of UniFrac such as Generalized UniFrac<sup>19</sup> and  
115 Variance Adjusted UniFrac<sup>20</sup>. To facilitate broader adoption, a generalized x86-64 build of  
116 library is now part of QIIME2's "q2-diversity" plugin, the same plugin used by Qiita.

117  
118 The design of Striped UniFrac allows UniFrac to scale well into the future, with demonstrated  
119 operational capability of >100,000 samples. Its parallel model and empirical scaling suggests  
120 application to datasets an order of magnitude larger than the EMP, and already benefits users  
121 operating on studies of "merely" thousands of samples such as the Integrated Human  
122 Microbiome Project dataset, in which UniFrac can be computed in 0.5 seconds on a laptop,  
123 allowing interactive analyses rather than batch-mode (fig. 2E). These reductions are critical for  
124 methods assay variability introduced by rarefaction, such as jackknifed beta diversity. Similarly,  
125 these reductions continue the democratization of analysis, bringing projects at the scale of the  
126 Earth Microbiome Project from a supercomputer to your laptop.

## 127 128 **Methods**

### 129 *Postorder traversal memory reduction*

130 At initialization, a stack is created to store sample proportion vectors of type double and length  
131  $N$  where  $N$  is the number of samples. The vectors in the stack are used to represent sample  
132 proportions across the tree. The stack is used to avoid reallocation of sample proportion vectors  
133 over the tree; allocated memory is reused after a vertex has been evaluated with more memory  
134 allocated to the stack only as needed. The stack is combined with a hash table indexed by a  
135 unique node identifier so specific sample proportions can be retrieved if presently stored.

136  
137 Over a postorder traversal of the input tree, if the vertex examined does not have children (i.e.,  
138 is a leaf), then a proportion vector is popped from the stack (allocating memory if the stack is  
139 empty), and the proportions in this vector are set to the sample proportions for the observed  
140 feature from the input BIOM table. An entry is then added into a hash table mapping the index of  
141 the node to the address of the vector. If the vertex evaluated instead has children (i.e., is not a  
142 leaf), then a vector is popped from the stack (allocating memory if the stack is empty), the  
143 sample proportions for each child of the vertex are obtained from the hash table, and the  
144 sample proportions of the children are summed to create the proportion vector for the vertex

145 under evaluation. The sample proportions of the children are then pushed on to the onto the  
146 stack, and the hash table entry for their vertex identifiers is erased. Because the traversal is  
147 performed in postorder, the children of a vertex are always evaluated first, which ensures the  
148 sample proportions of the children are present in the hash table.

149

#### 150 *Stripe aggregation*

151 A matrix of size  $K \times N$ , where  $K$  is the number of stripes and  $N$  is the number of samples is  
152 allocated. The number of stripes for a full distance matrix  $(N + 1) / 2$ , resulting in an allocation of  
153  $((N + 1) / 2) * N$  elements. For an odd number of samples, the number of elements in this  
154 matrix is exactly the number of elements in the upper triangle of the distance matrix. For an  
155 even number of samples,  $(N / 2)$  space is replicated. The elements of these stripes span both  
156 the upper and lower triangles of the logical distance matrix.

157

158 For each vertex, the sample proportions  $P$  for the vertex are embedded into a vector  $E$  of length  
159  $2N$ , such that the first  $N$  elements are a copy of the sample proportions, and the  $N$  to  $2N$   
160 elements are a copy of the sample proportions. This embedding allows for bulk pairwise  
161 comparisons between samples by allowing the execution of a distance kernel  $D$  (e.g.,  
162 unweighted UniFrac) into a single stripe  $k$  with  $D(P, E[k:k+N])$ . This approach results in the  
163 comparison of two linear and  $C$  contiguous blocks of memory of length  $N$ , storing into a linear  
164 and  $C$  contiguous block of memory of length  $N$ .

165

#### 166 *Distance kernels and parallelism*

167 Each individual distance metric (e.g., unweighted UniFrac, Generalized UniFrac, etc) is  
168 expressed as a compute kernel. These compute kernels operate per vertex during the tree  
169 traversal, compute a metric at every node within the tree, adding the computed distances into  
170 the stripes. Profiling suggests the vast bulk (>99%) of the computational time is expended in  
171 these kernels.

172

173 The compute of any diagonal (stripe) in the resulting distance matrix is independent of any other  
174 diagonal. This property allows the computation to be expressed as a map-reduce problem,  
175 where we *map* stripe sets to processing engines, and *reduce* by consolidating the stripes into a  
176 logical distance matrix or condensed form matrix.

177

#### 178 *Data sets*

179 The full Deblur 90nt table rarefied at 1000 sequences per sample, and corresponding  
180 phylogenetic tree were obtained from the Earth Microbiome Project. Sets of samples were  
181 randomly pulled from this table at from 1000 to 10000 samples in steps of 1000. At each sample  
182 size, 10 iterations were performed producing 10 tables. For a given set of samples, the subset  
183 of the tree vertices and edges ancestral to the features in a sample set were retained, all other  
184 vertices and edges were pruned out. Any single descendant nodes were collapsed, aggregating  
185 the branch length toward the root.

186

187 The iHMP metagenomic data were obtained from the iHMP Qiita portal (<https://ihmp.ucsd.edu>),  
188 study 10001. Preparation ID 5 of the metagenomic data were downloaded as a BIOM table.

189 Only observations at the species level were retained. A tree was constructed using the lineage  
190 information embedded in the IDs. Unweighted UniFrac was computed on this table using the  
191 Python API for Striped UniFrac using 4 threads.

192

193 *Pseudocode representation of the algorithm with unweighted UniFrac*

```
194 function unweighted(embedded_props, stripes, stripe_totals)
195     n_samples = number_of_columns(stripes)
196     for stripe_index in stripes
197         start = stripe_index
198         end = start + n_samples
199         stripe_props = embedded_props[start:end]
200
201         unique = embedded_props[:n_samples] XOR stripe_props
202         total = embedded_props[:n_samples] OR stripe_props
203
204         # inplace operation
205         stripes[stripe_index] += unique
206         stripe_totals[stripe_index] += total
207
208
209 function initialize_child_proportions()
210     stack_of_proportions = empty()
211     hashmap_of_proportions = empty()
212     return (stack_of_proportions, hashmap_of_proportions)
213
214
215 function merge_child_proportions(child_props, tree, node)
216     node_props = empty()
217     for child in children(tree, node)
218         child_prop = child_props.hashmap_of_props.pop(child)
219         node_props += child_prop
220         child_props.stack_of_props.push(child_prop)
221     return node_props
222
223
224 function associate_node_props(child_props, node, props)
225     child_props.hashmap_of_props[node] = props
226
227
228 function get_prop_vector(child_props)
229     if child_props.stack_of_props.empty()
230         return empty_vector_of_length_number_of_samples
231     else
232         return child_props.stack_of_props.pop()
233
```

```

234
235 function embed_props(vector)
236     embedded = zeros(length(vector) * 2)
237     embedded[:length(vector)] = vector
238     embedded[length(vector):] = vector
239     return embedded
240
241
242 function get_and_set_leaf_vector(table, node, props)
243     for index, value in getleaf(table, node)
244         props[index] = value
245
246
247 function deconvolute(stripes)
248     n_samples = length(stripes[0])
249     matrix = zeros(n_samples, n_samples)
250
251     for index, stripe in enumerate(stripes)
252         k = 0
253         row = 0
254         col = index + 1
255         while row < n_samples
256             if(col < n_samples):
257                 matrix[row, col] = stripe[k]
258                 matrix[col, row] = stripe[k]
259             else
260                 matrix[col % n_samples][row] = stripe[k]
261                 matrix[row][col % n_samples] = stripe[k]
262             row = row + 1
263             col = col + 1
264             k = k + 1
265     return matrix
266
267
268 function unfrac(table, tree, kernel)
269     n_samples = number_of_samples(table)
270     n_stripes = ceil((n_samples^2 - n_samples) / 2)
271     stripes = zeros(n_stripes, n_samples)
272     stripe_totals = zeros(n_stripes, n_samples)
273     child_props = initialize_child_props()
274
275     for node in postorder(tree)
276         if isleaf(node)
277             props = get_prop_vector(child_props)
278             get_and_set_leaf_vector(table, node, props)

```

```

279         associate_node_props(child_props, node, props)
280     else
281         props = merge_child_props(child_props, tree, node)
282
283         embedded_props = embed_props(props)
284         kernel(embedded_props, length(node), stripes, stripe_totals)
285
286     if kernel is normalized
287         for stripe_index in stripes:
288             unique = striped[stripe_index]
289             total = stripe_totals[stripe_index]
290             stripes[stripe_index] = unique / total
291
292     return deconvolute(stripes)
293

```

- 294 1. Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for comparing microbial  
295 communities. *Appl. Environ. Microbiol.* **71**, 8228–8235 (2005).
- 296 2. Smith, M. I. *et al.* Gut microbiomes of Malawian twin pairs discordant for kwashiorkor.  
297 *Science* **339**, 548–554 (2013).
- 298 3. Yatsunencko, T. *et al.* Human gut microbiome viewed across age and geography. *Nature*  
299 **486**, 222–227 (2012).
- 300 4. Ley, R. E., Lozupone, C. A., Hamady, M., Knight, R. & Gordon, J. I. Worlds within worlds:  
301 evolution of the vertebrate gut microbiota. *Nat. Rev. Microbiol.* **6**, 776–788 (2008).
- 302 5. Kuczynski, J. *et al.* Microbial community resemblance methods differ in their ability to detect  
303 biologically relevant patterns. *Nat. Methods* **7**, 813–819 (2010).
- 304 6. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale microbial diversity.  
305 *Nature* **551**, 457–463 (2017).
- 306 7. Hamady, M., Lozupone, C. & Knight, R. Fast UniFrac: facilitating high-throughput  
307 phylogenetic analyses of microbial communities including analysis of pyrosequencing and  
308 PhyloChip data. *ISME J.* **4**, 17–27 (2010).
- 309 8. McClelland, J. & Koslicki, D. EMDUniFrac: exact linear time computation of the UniFrac  
310 metric and identification of differentially abundant organisms. *J. Math. Biol.* (2018).



311 doi:10.1007/s00285-018-1235-9

- 312 9. Cordova, J. & Navarro, G. Simple and efficient fully-functional succinct trees. *Theor.*  
313 *Comput. Sci.* **656**, 135–145 (2016).
- 314 10. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing  
315 data. *Nat. Methods* **7**, 335–336 (2010).
- 316 11. Knight, R. *et al.* PyCogent: a toolkit for making sense from sequence. *Genome Biol.* **8**,  
317 R171 (2007).
- 318 12. McMurdie, P. J. & Holmes, S. phyloseq: an R package for reproducible interactive analysis  
319 and graphics of microbiome census data. *PLoS One* **8**, e61217 (2013).
- 320 13. Amir, A. *et al.* Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns.  
321 *mSystems* **2**, (2017).
- 322 14. McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological  
323 and evolutionary analyses of bacteria and archaea. *ISME J.* **6**, 610–618 (2012).
- 324 15. Mirarab, S., Nguyen, N. & Warnow, T. SEPP: SATé-enabled phylogenetic placement. *Pac.*  
325 *Symp. Biocomput.* 247–258 (2012).
- 326 16. Janssen, S. *et al.* Phylogenetic Placement of Exact Amplicon Sequences Improves  
327 Associations with Clinical Information. *mSystems* **3**, e00021–18 (2018).
- 328 17. Halko, N., Martinsson, P., Shkolnisky, Y. & Tygert, M. An Algorithm for the Principal  
329 Component Analysis of Large Data Sets. *SIAM J. Sci. Comput.* **33**, 2580–2594 (2011).
- 330 18. Vázquez-Baeza, Y., Pirrung, M., Gonzalez, A. & Knight, R. EMPeror: a tool for visualizing  
331 high-throughput microbial community data. *Gigascience* **2**, 16 (2013).
- 332 19. Chen, J. *et al.* Associating microbiome composition with environmental covariates using  
333 generalized UniFrac distances. *Bioinformatics* **28**, 2106–2113 (2012).
- 334 20. Chang, Q., Luan, Y. & Sun, F. Variance adjusted weighted UniFrac: a powerful beta  
335 diversity measure for comparing communities based on phylogeny. *BMC Bioinformatics* **12**,  
336 118 (2011).



Consider a rooted tree  $T$  with root  $\rho$  on  $n$  nodes as an  $n$  dimensional vector spaces over the real numbers. Identify the subtrees of  $T$  with the nodes of  $T$ , that is subtree  $i$  is the subtree which does not contain  $\rho$  formed by choosing node  $i$  as a root. The subtree corresponding to  $\rho$  is  $T$ . Choose a basis  $\{v_i | 1 \leq i \leq n\}$  for  $T$  such that  $v_i$  is the indicator function for subtree  $i$ , that is  $v_i(j) = 1$  for those nodes  $j$  in subtree  $i$  and zero otherwise. Label the  $n - 1$  edges of  $T$  such that edge  $e_i$  is the unique edge adjacent to node  $i$  on a path from node  $i$  to the root  $\rho$ .

Let  $W$  be the  $n \times n$  matrix whose rows correspond to the basis vectors  $v_i$  scaled by the corresponding edge weight  $l(e_i)$ . Consider probability distributions  $P$  and  $Q$  on  $T$  as column vectors, ordered such that entry  $i$  corresponds to the root of subtree  $i$ .

Then

$$\text{UniFrac}(P, Q) = \|W(P - Q)\|_{L_1}$$

Thus follows as we note that

$$\|W(P - Q)\|_{L_1} = \sum_{i=1}^n \sum_{j=1}^n l(e_i) v_i(j) |P(v_j) - Q(v_j)|$$

This summation is the discrete integral over  $T$  of the distributions  $P$  and  $Q$  with respect to measure formed from the branch lengths of  $T$ . By [1] the value of this integral is equivalent to earth mover's distance between  $P$  and  $Q$  with respect to the tree  $T$ , and thus is equal to the UniFrac distance between  $P$  and  $Q$ .

## References

- [1] Evans SN, Matsen FA. *The phylogenetic Kantorovic Rubinstein metric for environmental sequence samples*. Journal of the Royal Statistical Society Series B, Statistical methodology. 2012