

 Open access • Posted Content • DOI:10.1101/2020.10.06.327080

## **EMPress enables tree-guided, interactive, and exploratory analyses of multi-omic datasets** — [Source link](#)

Kalen Cantrell, Marcus W. Fedarko, Gibraan Rahman, Daniel McDonald ...+28 more authors

**Institutions:** University of California, Berkeley, University of Giessen, IBM, Arizona State University ...+4 more institutions

**Published on:** 08 Oct 2020 - bioRxiv (Cold Spring Harbor Laboratory)

Related papers:

- [EMPress Enables Tree-Guided, Interactive, and Exploratory Analyses of Multi-omic Data Sets.](#)
- [Phinch: An interactive, exploratory data visualization framework for metagenomic datasets](#)
- [tidyMicro: a pipeline for microbiome data analysis and visualization using the tidyverse in R](#)
- [PaintOmics 3: a web resource for the pathway analysis and visualization of multi-omics data](#)
- [Microreact: visualizing and sharing data for genomic epidemiology and phylogeography.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/empres-enables-tree-guided-interactive-and-exploratory-h5j26ubrww>

# 1 EMPress enables tree-guided, interactive, and exploratory 2 analyses of multi-omic datasets

3 Kalen Cantrell<sup>\*1,2</sup>, Marcus W. Fedarko<sup>\*1,2</sup>, Gibraan Rahman<sup>3</sup>, Daniel McDonald<sup>4</sup>, Yimeng  
4 Yang<sup>2</sup>, Thant Zaw<sup>2</sup>, Antonio Gonzalez<sup>4</sup>, Stefan Janssen<sup>5</sup>, Mehrbod Estaki<sup>4</sup>, Niina Haiminen<sup>6</sup>,  
5 Kristen L. Beck<sup>7</sup>, Qiyun Zhu<sup>4,8</sup>, Erfan Sayyari<sup>2,9</sup>, Jamie Morton<sup>10</sup>, Anupriya Tripathi<sup>4</sup>, Julia M.  
6 Gauglitz<sup>15</sup>, Clarisse Marotz<sup>4,11</sup>, Nathaniel L. Matteson<sup>13</sup>, Cameron Martino<sup>2,3,4</sup>, Jon G.  
7 Sanders<sup>16</sup>, Anna Paola Carrieri<sup>14</sup>, Se Jin Song<sup>2</sup>, Austin D. Swafford<sup>2</sup>, Pieter C. Dorrestein<sup>2,15</sup>,  
8 Kristian G. Andersen<sup>13</sup>, Laxmi Parida<sup>6</sup>, Ho-Cheol Kim<sup>7</sup>, Yoshiki Vázquez-Baeza<sup>2</sup>, Rob  
9 Knight<sup>1,2,4,12</sup>.

10

11 \*Contributed Equally.

12 <sup>1</sup>Department of Computer Science, Jacobs School of Engineering, University of California, San  
13 Diego.

14 <sup>2</sup>Center for Microbiome Innovation, Jacobs School of Engineering, University of California, San  
15 Diego.

16 <sup>3</sup>Bioinformatics and Systems Biology Program, University of California, San Diego.

17 <sup>4</sup>Department of Pediatrics, School of Medicine, University of California, San Diego.

18 <sup>5</sup>Algorithmic Bioinformatics, Justus Liebig University Giessen, Germany.

19 <sup>6</sup>IBM T. J. Watson Research Center, Yorktown Heights, New York.

20 <sup>7</sup>IBM Almaden Research Center, San Jose, California.

21 <sup>8</sup>Present Address: Biodesign Center for Fundamental and Applied Microbiomics, Arizona State  
22 University, Arizona.

23 <sup>9</sup>Department of Electrical and Computer Engineering, University of California, San Diego.

24 <sup>10</sup>Center for Computational Biology, Flatiron Institute, Simons Foundation, New York, NY, USA.

25 <sup>11</sup>Scripps Institution of Oceanography, University of California, San Diego.

26 <sup>12</sup>Department of Bioengineering, University of California, San Diego.

27 <sup>13</sup>Scripps Research Institute, San Diego, California.

28 <sup>14</sup>IBM Research, The Hartree Centre, Daresbury, UK

29 <sup>15</sup>Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and  
30 Pharmaceutical Sciences, University of California, San Diego.

31 <sup>16</sup>Cornell Institute for Host-Microbe Interaction and Disease, Cornell University, Ithaca

## 32 Abstract

33 Standard workflows for analyzing microbiomes often include the creation and curation of  
34 phylogenetic trees. Here we present EMPress, an interactive tool for visualizing trees in the  
35 context of microbiome, metabolome, etc. community data scalable beyond modern large  
36 datasets like the Earth Microbiome Project. EMPress provides novel functionality—including  
37 ordination integration and animations—alongside many standard tree visualization features, and  
38 thus simplifies exploratory analyses of many forms of ‘omic data.

39

40

## 41 Main Text

42

43 The increased availability of sequencing technologies and automation of molecular methods  
44 have enabled studies of unprecedented scale [1] prompting the creation of tools better suited to  
45 store, analyze [2], and visualize [3] studies of this magnitude. Many of these tools, such as [4, 5,  
46 6, 7], use phylogenies detailing the evolutionary relationships among features or dendrograms  
47 that organize features in a hierarchical structure (e.g. clustering of mass spectra) [8]. The  
48 challenge of enabling fully interactive analyses stems from the disconnect between feature-level  
49 tools and dataset-level tools; few can interactively integrate multiple representations of the data  
50 [9], and to our knowledge none scale to display large datasets. This is a key unresolved  
51 challenge for the field: to allow researchers to contextualize community-level patterns  
52 (groupings of samples) together with feature-level structure, i.e. which features lead to the  
53 groupings explained in a given sample set.

54

55 Here, we introduce EMPress (<https://github.com/biocore/empress>), an open-source (BSD 3-  
56 clause), interactive and scalable phylogenetic tree viewer accessible as a QIIME 2 [2] plugin.  
57 EMPress is built around the high-performance balanced parentheses tree data structure [10],  
58 and uses a hardware-accelerated WebGL-based rendering engine that allows EMPress to  
59 visualize trees with hundreds of thousands of nodes using a laptop's web browser (Methods).  
60 By integrating EMPress with the widely-used EMPeror software [3] within QIIME 2, EMPress  
61 can simultaneously visualize a phylogenetic tree of features in a study coupled with an  
62 ordination of the same study's samples. User actions in one visualization, such as selecting a  
63 set of samples in the ordination, update the other, providing context that would not be easily  
64 accessible with independent visualizations. This tight integration between displays streamlines  
65 several use-cases elaborated below that previously required manual investigation or writing  
66 custom scripts.

67

68 EMPress visualizations can be created solely from a tree, or users can provide additional  
69 metadata files and a feature table to augment the tree. Using these common data files, users  
70 can interactively configure many visual attributes in the tree (see Methods and Figures for  
71 examples).

72

73 Rather than providing a programmatic interface for the procedural generation of styled  
74 phylogenetic trees [11, 12, FigTree (<http://tree.bio.ed.ac.uk/software/figtree/>)], EMPress  
75 provides an interactive environment to support exploratory feature- and sample-level tree-based  
76 analyses. Many use-cases supported in EMPress accommodate community analysis tasks; this  
77 differs from Anvi'o [13] which is centered on the analysis of metagenomic assembled-genomes,  
78 pangenomes, etc.. PHYLOViZ [9], SigTree [14], and iTOL [15] are similar to EMPress in terms  
79 of their implementation (PHYLOViZ Online also uses WebGL), and/or use-cases (SigTree is  
80 mostly used to visualize differential abundance patterns, and iTOL supports the visualization of  
81 QIIME 2 tree artifacts). EMPress stands out in its scalability: iTOL claims trees with more than

82 10,000 tips to be “very large” (<https://itol.embl.de/help.cgi>), while EMPress readily supports trees  
83 with over hundreds of thousands of tips, as shown in Fig. 1. Many visualization customization  
84 options available in EMPeror, iTOL [15] and Anvi’o [13] are immediately accessible in EMPress’  
85 interface. Continuous feature metadata can be visualized in tip-level barplots as a color gradient  
86 and/or by adjusting the lengths of individual tips’ barplots; categorical sample metadata  
87 information can be visualized using a stacked barplot showing—for each tip—the proportion of  
88 samples containing that tip stratified by category. These options are available on the user  
89 interface and do not require programming or configuration files.

90  
91 Ordination plots computed from UniFrac distances are often used to visualize sample clustering  
92 patterns in microbiome studies. However, interpreting the patterns in these plots—and  
93 determining which features influence sample group separation—is not always straightforward.  
94 While biplots show information about influential features alongside samples, the phylogenetic  
95 relationships of these features are not immediately obvious. EMPress aids interpretation of  
96 these plots by optionally providing a unified interface where the tree and ordination  
97 visualizations are displayed side-by-side and “linked” through sample and feature identifiers  
98 [16]. This combination allows for novel exploratory data analysis tasks. For example, selecting a  
99 group of samples in the ordination highlights nodes in the tree present in those samples, and  
100 vice versa (see Methods). This integration extends to biplots: clicking feature arrows in the  
101 ordination highlights their placement in the tree. Lastly, EMPress allows visualizing longitudinal  
102 studies by simultaneously showing the tree nodes unique to groups of samples at each  
103 individual time point during an EMPeror animation (see Methods).

104  
105 Using the first data release of the Earth Microbiome Project (EMP), we demonstrate EMPress’  
106 scalability by rendering a 26,035 sample ordination and a 756,377 node tree (Figure 1A). To  
107 visualize the relative proportions of taxonomic groups at the phylum level, we use EMPress’  
108 feature metadata coloring to highlight the top 5 most prevalent phyla (see Methods). Next, we  
109 add a barplot layer showing, for each tip in the tree, the proportions of samples containing each  
110 tip summarized by level 2 of the EMP ontology (Animal, Plant, Non-Saline, and Saline). Paired  
111 visualizations allow us to click on a tip in the tree and view the samples that contain that feature  
112 in the ordination. This functionality is useful when analyzing datasets with outliers or mislabeled  
113 metadata. Tip-aligned barplots summarize environmental metadata: for example, Figure 1B  
114 shows the subset of samples (4,002) with recorded pH information and a barplot layer with the  
115 mean pH where each feature was found. The barplot reveals a relatively dark section near  
116 many Firmicutes-classified features on the tree; in concert with histograms showing mean pH  
117 for each phylum (Figure 1C), we can confirm that Firmicutes-classified features are more  
118 commonly found in higher pH environments.

119  
120 EMPress can be applied to various ‘omic datasets. To illustrate this versatility we reanalyzed a  
121 COVID-19 metatranscriptome sequencing dataset [17], a liquid chromatography mass-  
122 spectrometry (LC-MS) untargeted metabolomic food-associated dataset [8], and a 16S rRNA  
123 sequencing oral microbiome dataset [18]. Despite the vastly different natures of these datasets,  
124 EMPress provides meaningful functionality for their analysis and visualization. Supplemental  
125 Video 1 ([supplementary-video-1.mp4](#)) shows a longitudinal exploratory analysis using EMPress

126 and EMPeror representing a subset of SARS-CoV-2 genome data from GISAID. This paired  
127 visualization emphasizes the relationships in time and space among “community samples” and  
128 the convergence of locales in the United States with the outbreak in Italy (See Methods). The  
129 interactive nature of EMPress allows rapid visualization of strains observed in a collection of  
130 samples from different geographical locations.

131  
132 Figure 2A showcases Empress’ ability to identify feature clusters that are differentially abundant  
133 in COVID-19 patients compared to community-acquired pneumonia patients and healthy  
134 controls [17]. Clades showing KEGG enzyme code (EC) [19] annotations are collapsed at level  
135 two except for lyases, highlighting feature 4.1.1.20 (carboxy-lyase diaminopimelate  
136 decarboxylase) that was more abundant in COVID-19 here and in an independent  
137 metaproteomic analysis of COVID-19 respiratory microbiomes [20].

138  
139 Recent developments in cheminformatics enabled the analysis and visualization of small  
140 molecules in the context of a cladogram [8]. Using a tree that links molecules by their structural  
141 relatedness, we analyzed untargeted LC-MS/MS data from 70 food samples (see Methods).  
142 With EMPress’ sample metadata barplots, we can inspect the relationship between chemical  
143 annotations and food types. Figure 2B shows a tree where each tip is colored by its chemical  
144 super class, and where barplots show the proportion of samples in the study containing each  
145 compound by food type. This representation reveals a clade of lipids and lipid-like molecules  
146 that are well represented in animal food types and seafoods. In contrast, salads and fruits are  
147 broadly spread throughout the cladogram.

148  
149 Lastly, in Figure 2C, we compare three differential abundance methods in an oral microbiome  
150 dataset [18] as separate barplot layers on a tree. This dataset includes samples (n=32) taken  
151 before and after subjects brushed their teeth (see Methods). As observed across the three  
152 differential abundance tools’ outputs, all methods agree broadly on which features are  
153 particularly “differential” (for example, the cluster of Firmicutes-classified sequences in the  
154 bottom-right of the tree; see Methods), although there are discrepancies due to different  
155 methods’ assumptions and biases.

## 156 Conclusions

157 By providing an intuitive interface supporting both categorically new and established  
158 functionality, EMPress complements and extends the available range of tree visualization  
159 software. EMPress can perform community analyses across distinct “omics” types, as  
160 demonstrated here. Moving forward, facilitating the integration of multiple orthogonal views of a  
161 dataset at a more generalized framework level (for example, using QIIME2’s [2] visualization  
162 API) will be important as datasets continue to grow in complexity, size, and heterogeneity.

## 163 Acknowledgements

164 We thank members of the Knight Lab and IBM AIHL Bioinformatics team for feedback during  
165 code reviews and presentations. We gratefully acknowledge the following Authors from the  
166 originating laboratories responsible for obtaining the specimens and the submitting laboratories

167 where genetic sequence data were generated and shared via the GISAID Initiative, on which a  
168 portion of this research is based (Supplemental Table 1).  
169

## 170 Funding

171 This work is partly supported by IBM Research AI through the AI Horizons Network and the UC  
172 San Diego Center for Microbiome Innovation (to KC, MWF, JM, QZ, TZ, JS, ADS, SS, YVB,  
173 RK); CCF foundation #675191 (RK and PCD), U19 AG063744 01 (RK, JMG, PCD), CDC  
174 contract #75D30120C09795 (KGA and RK), U19 AI135995 (KGA).

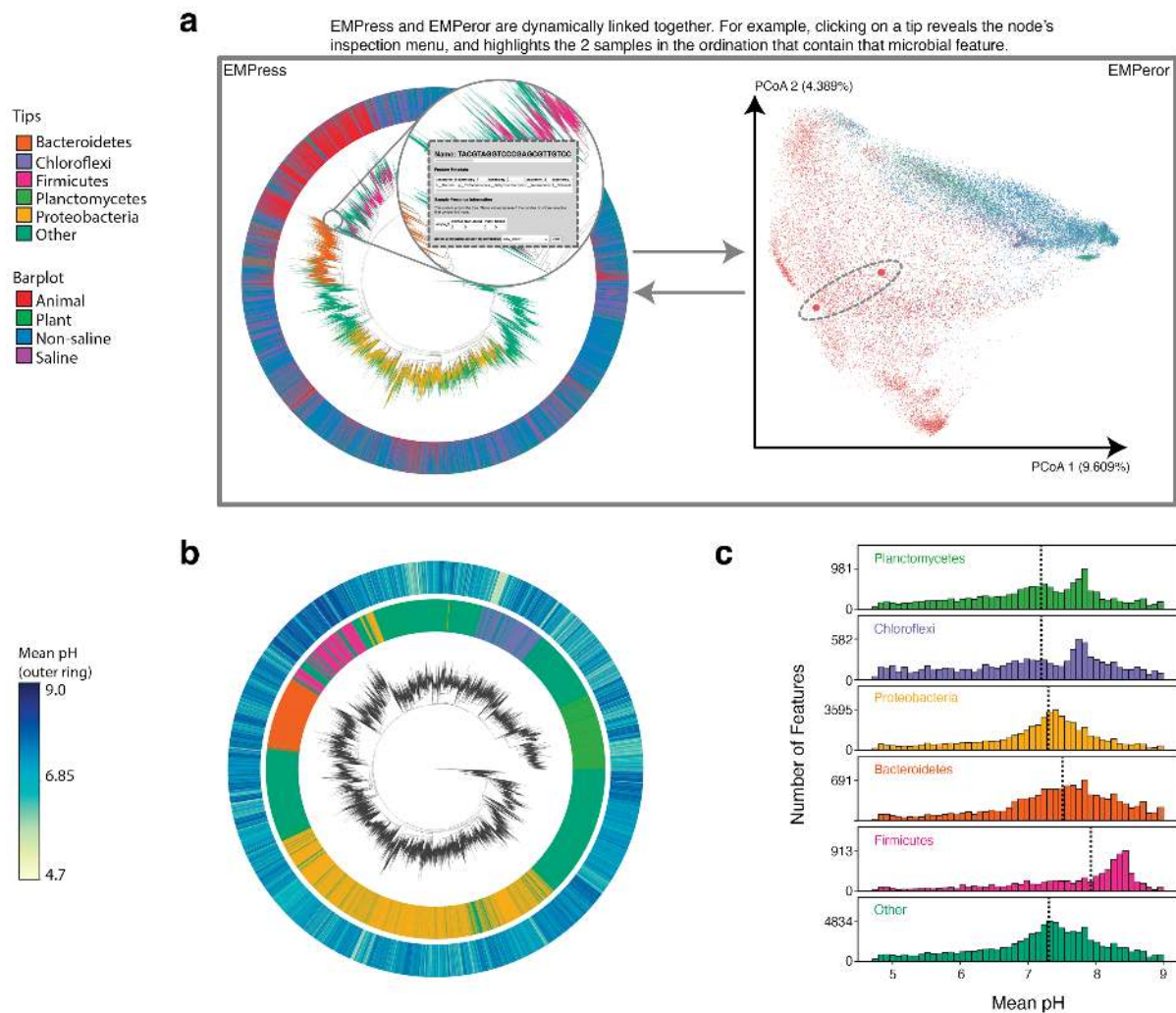
## 175 Author Contributions

176 KC, QZ, YY, JM, TZ, JS, and RK conceived the original idea for the project. KC, MWF, GR, DM,  
177 AG, SJ, ME, YY, ES, JM, TZ, QZ, YVB wrote source code and/or documentation for the project.  
178 KC, MWF, AG, YVB wrote code to facilitate integration with EMPeror. LP, HCK, SS, ADS, YVB,  
179 RK managed the project. KC, MWF, GR, NH, KLB, AT, JMG, LM, APC, NLM, CM, PCD, KGA,  
180 LP, YVB analyzed and interpreted the datasets presented in this paper. KC, MWF, GR, DM,  
181 NH, KLB, YVB contributed text to the methods section. All the authors contributed to the final  
182 version of the manuscript.

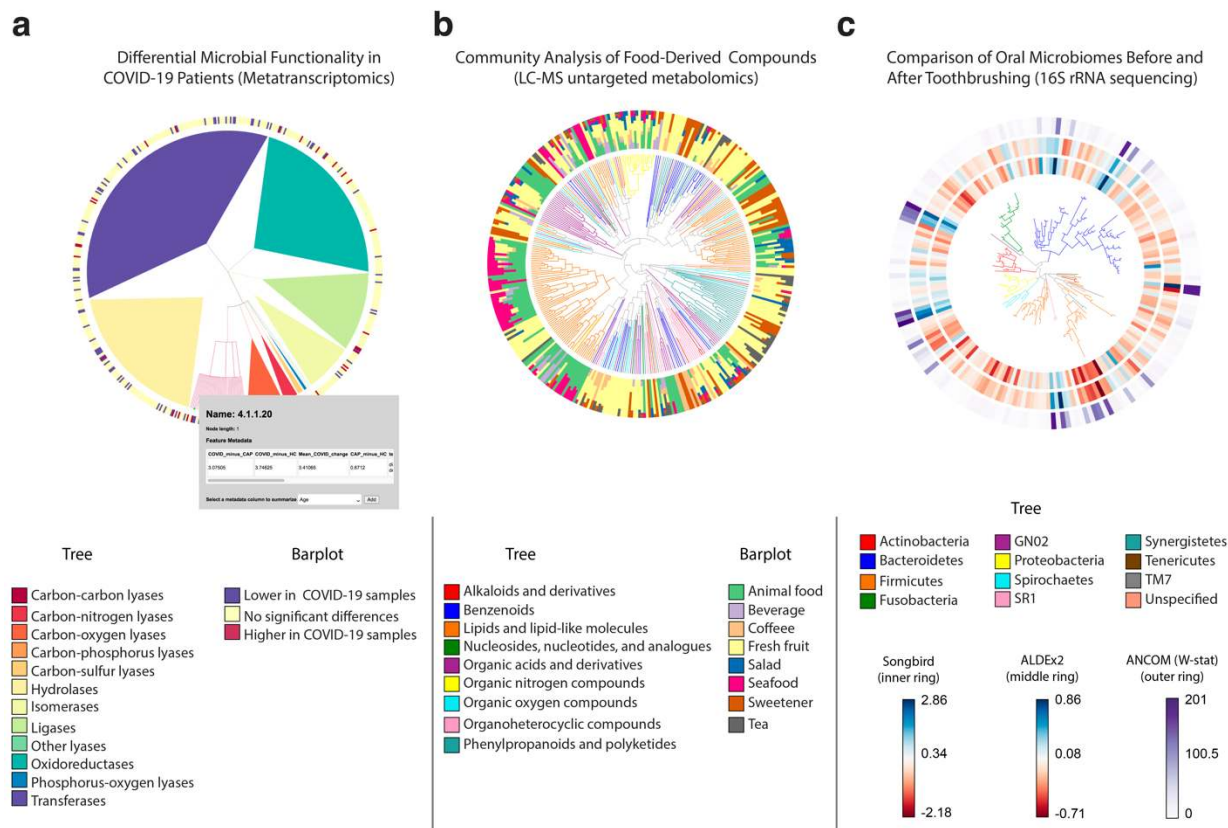
## 183 Competing Interests

184 We declare none.





185  
 186 **Figure 1. Earth Microbiome Project paired phylogenetic tree (including 756,377 nodes) and**  
 187 **unweighted UniFrac ordination (including 26,035 samples) . (a)** Graphical depiction of EMPress'  
 188 unified interface with fragment insertion tree (left), and unweighted UniFrac sample ordination (right). Tips  
 189 are colored by their phylum-level taxonomic assignment; the barplot layer is a stacked barplot describing  
 190 the proportions of samples containing each tip summarized by level 2 of the EMP ontology. Inset shows  
 191 summarized sample information for a selected feature. The ordination highlights the two samples  
 192 containing the tip selected in the tree enlarged to show their location. **(b)** Subset of EMP samples with pH  
 193 information: the inner barplot ring shows the phylum-level taxonomic assignment, and the outer barplot  
 194 ring represents the mean pH of all the samples where each tip was observed **(c)** pH distributions  
 195 summarized by phylum-level assignment with median pH indicated by dotted lines. Interactive figures can  
 196 be accessed [here](#).  
 197



198  
 199 **Figure 2. EMPress is a versatile exploratory analysis tool adaptable to various -omics data types.**  
 200 **(a)** RoDEO differential abundance scores of microbial functions from metatranscriptomic sequencing of  
 201 COVID-19 patients (n=8), community-acquired pneumonia patients (n=25), and healthy control subjects  
 202 (n=20). The tree represents the four-level hierarchy of the KEGG enzyme code. The barplot colors  
 203 significantly differentially abundant features (p<0.05) in COVID-19 patients. Clicking on a tip produces a  
 204 pop-up insert tabulating the name of the feature, its hierarchical ranks, and any feature annotations.  
 205 **(b)** Global FoodOmics Project LC-MS data. Stacked barplots indicate the proportions of samples (n=70)  
 206 (stratified by food) containing the tips in an LC-MS Qemistree of food-associated compounds, with tip  
 207 nodes colored by their chemical superclass.  
 208 **(c)** *de novo* tree constructed from 16S rRNA sequencing data from 32 oral microbiome samples. Samples  
 209 were taken before (n=16) and after (n=16) subjects (n=10) brushed their teeth; each barplot layer  
 210 represents a different differential abundance method's measure of change between before- and after-  
 211 brushing samples. The innermost layer shows estimated log-fold changes produced by Songbird; the  
 212 middle layer shows effect sizes produced by ALDEx2; and the outermost layer shows the W-statistic  
 213 values produced by ANCOM (see Methods). The tree is colored by tip nodes' phylum-level taxonomic  
 214 classifications. Interactive figures can be accessed [here](#).  
 215  
 216  
 217

## 218 References

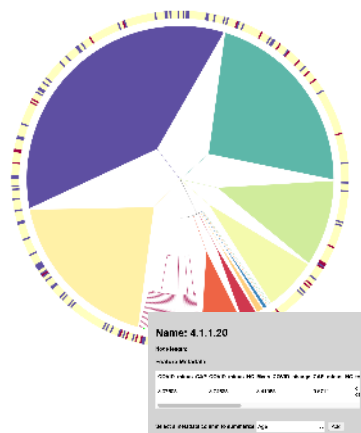
219 1. Thompson, L. R. *et al.* A communal catalogue reveals Earth's multiscale microbial  
 220 diversity. *Nature* **551**, 457–463 (2017).



- 221 2. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data  
222 science using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
- 223 3. Vázquez-Baeza, Y., Pirrung, M., Gonzalez, A. & Knight, R. EMPeror: a tool for  
224 visualizing high-throughput microbial community data. *Gigascience* **2**, 16 (2013).
- 225 4. Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for comparing microbial  
226 communities. *Appl. Environ. Microbiol.* **71**, 8228–8235 (2005).
- 227 5. Washburne, A. D. *et al.* Phylogenetic factorization of compositional data yields lineage-  
228 level associations in microbiome datasets. *PeerJ* **5**, e2969 (2017).
- 229 6. Silverman, J. D., Washburne, A. D., Mukherjee, S. & David, L. A. A phylogenetic  
230 transform enhances analysis of compositional microbiota data. *Elife* **6**, (2017).
- 231 7. Morton, J. T. *et al.* Balance Trees Reveal Microbial Niche Differentiation. *mSystems* **2**,  
232 (2017).
- 233 8. Tripathi, A. *et al.* Chemically-informed Analyses of Metabolomics Mass Spectrometry  
234 Data with Qemistree. *bioRxiv* 2020.05.04.077636 (2020)  
235 doi:[10.1101/2020.05.04.077636](https://doi.org/10.1101/2020.05.04.077636).
- 236 9. Nascimento, M. *et al.* PHYLOViZ 2.0: providing scalable data integration and  
237 visualization for multiple phylogenetic inference methods. *Bioinformatics* **33**, 128–129  
238 (2017).
- 239 10. Cordova, J. & Navarro, G. Simple and efficient fully-functional succinct trees. *Theoretical*  
240 *Computer Science* **656**, 135–145 (2016).
- 241 11. Yu, G. Using ggtree to Visualize Data on Tree-Like Structures. *Curr Protoc*  
242 *Bioinformatics* **69**, e96 (2020).
- 243 12. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and  
244 Visualization of Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
- 245 13. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data.  
246 *PeerJ* **3**, e1319 (2015).
- 247 14. Stevens, J. R., Jones, T. R., Lefevre, M., Ganesan, B. & Weimer, B. C. SigTree: A  
248 Microbial Community Analysis Tool to Identify and Visualize Significantly Responsive  
249 Branches in a Phylogenetic Tree. *Comput Struct Biotechnol J* **15**, 372–378 (2017).
- 250 15. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new  
251 developments. *Nucleic Acids Res* **47**, W256–W259 (2019).
- 252 16. Becker, R. A., Cleveland, W. S. & Wilks, A. R. Dynamic Graphics for Data Analysis.  
253 *Statist. Sci.* **2**, 355–383 (1987).
- 254 17. Shen, Z. *et al.* Genomic Diversity of Severe Acute Respiratory Syndrome–Coronavirus 2  
255 in Patients With Coronavirus Disease 2019. *Clin Infect Dis* **71**, 713–720 (2020).
- 256 18. Morton, J. T. *et al.* Establishing microbial composition measurement standards with  
257 reference frames. *Nat Commun* **10**, 2719 (2019).
- 258 19. Kanehisa, M. Enzyme Annotation and Metabolic Reconstruction Using KEGG. *Methods*  
259 *Mol. Biol.* **1611**, 135–145 (2017).
- 260 20. Maras, J. S. *et al.* Multi-Omics integration analysis of respiratory specimen characterizes  
261 baseline molecular determinants associated with COVID-19 diagnosis. *medRxiv*  
262 2020.07.06.20147082 (2020) doi:[10.1101/2020.07.06.20147082](https://doi.org/10.1101/2020.07.06.20147082).
- 263  
264

**a**

### Differential Microbial Functionality in COVID-19 Patients (Metatranscriptomics)



#### Tree

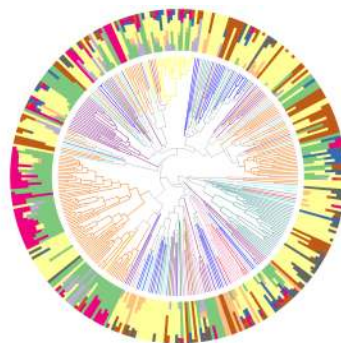
- Carbon-carbon lyases
- Carbon-nitrogen lyases
- Carbon-oxygen lyases
- Carbon-phosphorus lyases
- Carbon-sulfur lyases
- Hydrolases
- Isomerases
- Ligases
- Other lyases
- Oxidoreductases
- Phosphorus-oxygen lyases
- Transferases

#### Barplot

- Lower in COVID-19 samples
- No significant differences
- Higher in COVID-19 samples

**b**

### Community Analysis of Food-Derived Compounds (LC-MS untargeted metabolomics)



#### Tree

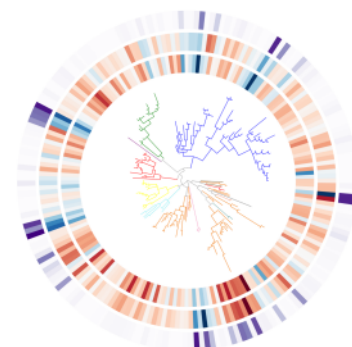
- Alkaloids and derivatives
- Benzenoids
- Lipids and lipid-like molecules
- Nucleosides, nucleotides, and analogues
- Organic acids and derivatives
- Organic nitrogen compounds
- Organic oxygen compounds
- Organoheterocyclic compounds
- Phenylpropanoids and polyketides

#### Barplot

- Animal food
- Beverage
- Coffee
- Fresh fruit
- Salad
- Seafood
- Sweetener
- Tea

**c**

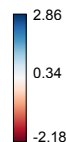
### Comparison of Oral Microbiomes Before and After Toothbrushing (16S rRNA sequencing)



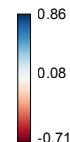
#### Tree

- Actinobacteria
- Bacteroidetes
- Firmicutes
- Fusobacteria
- GN02
- Proteobacteria
- Spirochaetes
- SR1
- Synergistetes
- Tenericutes
- TM7
- Unspecified

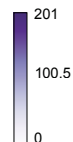
Songbird  
(inner ring)



ALDEx2  
(middle ring)

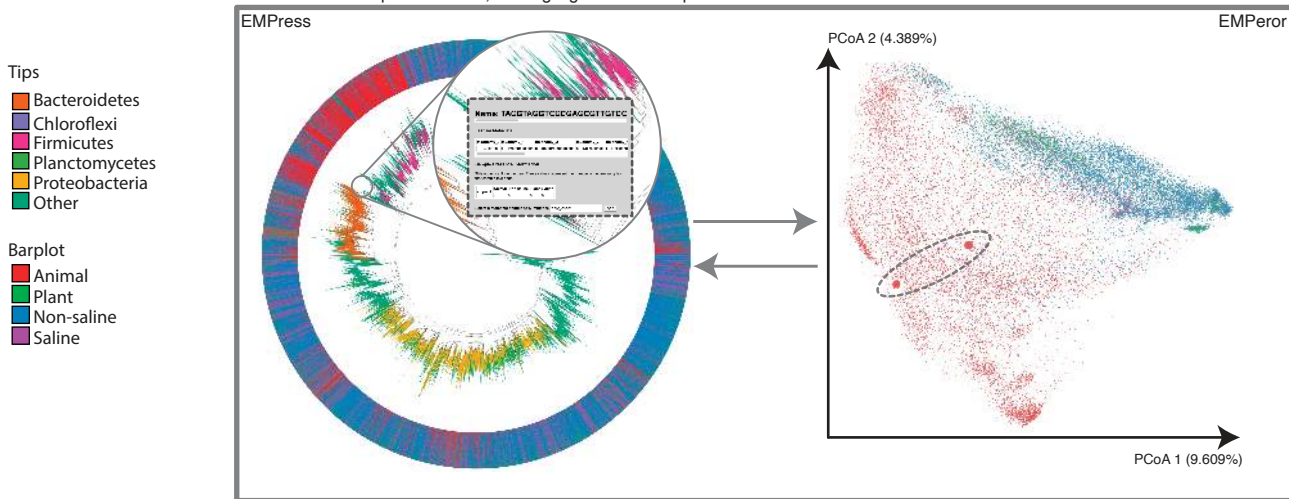
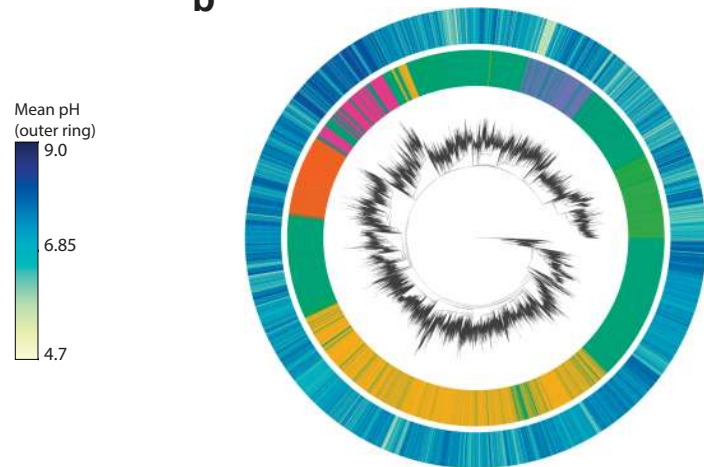


ANCOM (W-stat)  
(outer ring)



**a**

EMPress and EMPeror are dynamically linked together. For example, clicking on a tip reveals the node's inspection menu, and highlights the 2 samples in the ordination that contain that microbial feature.

**b****c**