

1 **Scalable power analysis and effect size exploration of microbiome community**

2 **differences with Evident**

3

4 Gibraan Rahman<sup>1,2</sup>, Daniel McDonald<sup>1</sup>, Antonio Gonzalez<sup>1</sup>, Yoshiki Vázquez-Baeza<sup>3</sup>, Lingjing

5 Jiang<sup>4</sup>, Climent Casals-Pascual<sup>5</sup>, Shyamal Peddada<sup>6</sup>, Daniel Hakim<sup>1,2</sup>, Amanda Hazel

6 Dilmore<sup>1,7</sup>, Brent Nowinski<sup>8</sup>, Rob Knight<sup>1,9,10,\*</sup>

7

8 1. Department of Pediatrics, School of Medicine, University of California, San Diego,  
9 California, USA

10 2. Bioinformatics and Systems Biology Program, University of California, San Diego,  
11 California, USA

12 3. BiomeSense Inc, Chicago, IL, USA

13 4. Janssen Research & Development, Spring House, PA, USA

14 5. Department of Microbiology, CDB, Hospital Clinic, University of Barcelona, Barcelona,  
15 Spain

16 6. Biostatistics and Bioinformatics Branch, Eunice Kennedy Shriver National Institute of  
17 Child Health and Human Development (NICHD), NIH, Bethesda, MD, USA

18 7. Biomedical Sciences Program, University of California, San Diego, La Jolla, California,  
19 USA

20 8. Center for Microbiome Innovation, Jacobs School of Engineering, University of California  
21 San Diego, La Jolla, California, USA

22 9. Department of Computer Science and Engineering, University of California, San Diego,  
23 La Jolla, CA, USA

24 10. Department of Bioengineering, University of California, San Diego, La Jolla, CA, USA

25

26 (\*) Corresponding author

27

## 28 Abstract

29 Differentiating microbial communities among samples is a major objective in biomedicine.  
30 Quantifying the effect size of these differences allows researchers to understand the factors  
31 most associated with communities and to optimize the design and clinical resources required to  
32 address particular research questions. Here, we present Evident, a package for effect size  
33 calculations and power analysis on microbiome data and show that Evident scales to large  
34 datasets with numerous metadata covariates.

## 35 Main text

36 The microbiome has been implicated as a crucial factor in a broad range of health and disease  
37 outcomes. Differences in microbial communities have been linked to differential metabolic  
38 regulation, often resulting in drastic phenotypic changes. One of the key computational methods  
39 for quantifying these community changes is diversity analysis. Alpha diversity measures the  
40 overall breadth of microbial features represented in a single sample, while beta diversity  
41 quantifies the pairwise community differences between samples via some choice of distance  
42 metric. Determining the magnitude of diversity differences among groups of samples is one of  
43 the objectives of computational microbiome analysis.

44

45 Evaluating the putative differences among groups is most often performed through null  
46 hypothesis significance testing (NHST). Under this framework, researchers quantify the  
47 probability that an observed difference (or one more extreme) would be observed due to chance

48 (p-value). This value is often used as a measure of significance of diversity differences.

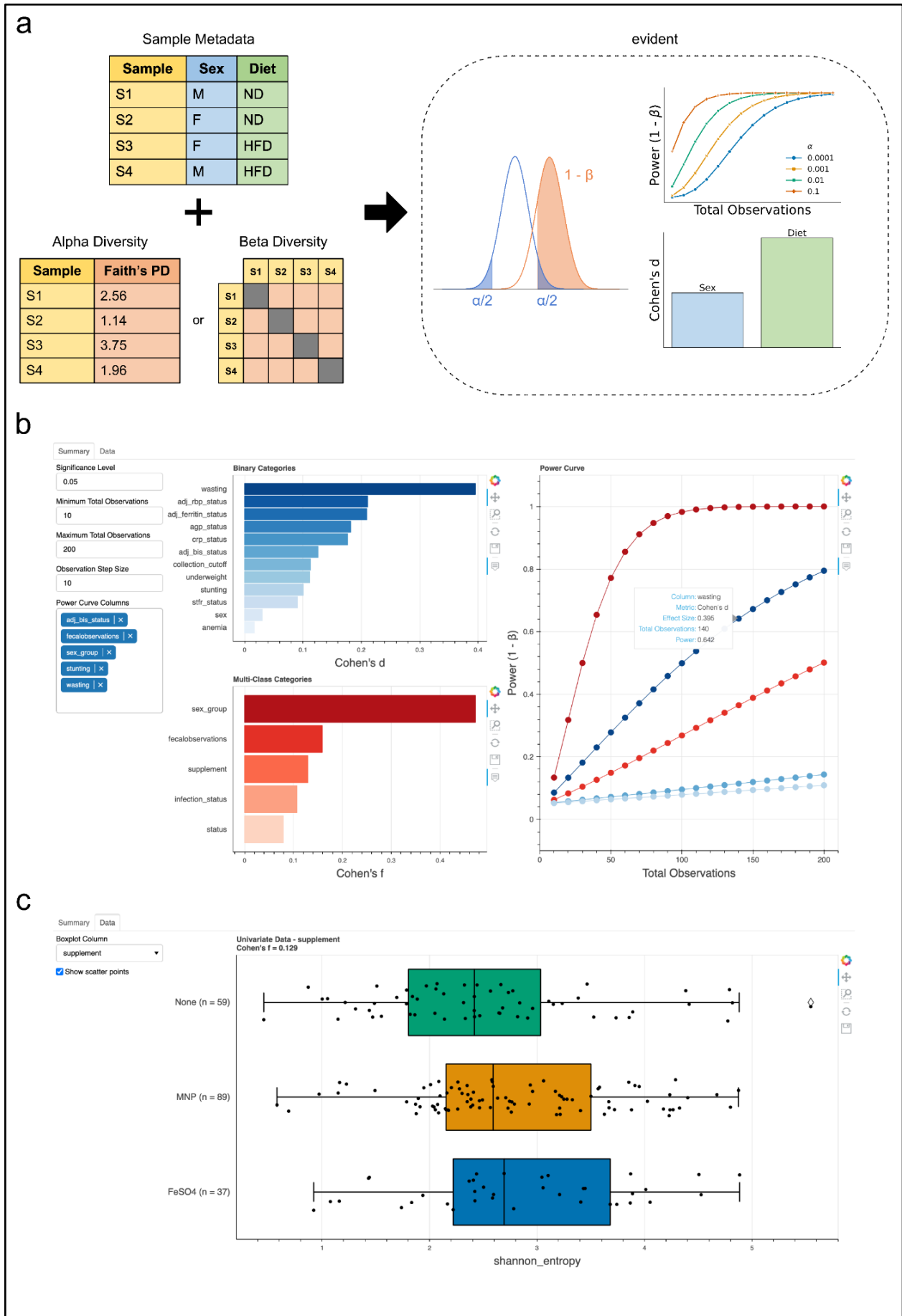
49 However, p-values by themselves are not enough. A significant p-value on its own does not  
50 provide any information about the magnitude of a given effect<sup>1</sup>.

51  
52 In addition to p-values, we suggest reporting the effect size of microbial community differences.  
53 Effect sizes quantify differences among sample groups and can be used to describe the  
54 magnitude of biological effects. Importantly, standardized effect sizes are dimensionless and  
55 can be compared between datasets and experiments<sup>2</sup>. Additionally, effect sizes allow  
56 researchers to determine the statistical power of new experimental designs. Statistical power is  
57 used to quantify the probability of making a Type II error (false negative) given that the  
58 alternative hypothesis is true<sup>3</sup>. Provided an effect size, desired significance level, and sample  
59 size, researchers can calculate the statistical power of experimental designs.

60  
61 Large scale microbiome data collection efforts, such as the American Gut Project (AGP)<sup>4</sup>,  
62 TEDDY<sup>5</sup>, and FINRISK<sup>6</sup> provide a unique opportunity to explore effect sizes across a wide  
63 variety of biological factors such as age, obesity, etc. These datasets contain dozens or even  
64 hundreds of metadata categories. Determining which covariates contribute the most to microbial  
65 diversity is crucial for prioritization of resources. Researchers interested in designing new  
66 experiments can keep these effect sizes in mind to efficiently allocate resources to maximize the  
67 chances of finding significant biological signal<sup>7</sup>.

68  
69 Here we introduce Evident, a new open source tool for efficient effect size and power  
70 calculations of microbiome data. Evident is available both as a standalone Python package as  
71 well as a QIIME 2 plugin<sup>8</sup>. With Evident, researchers can seamlessly explore the effect size of  
72 community differences in dozens of metadata columns at once.

73



75 **Fig 1: Evident workflow and interactive visualizations**

76 **a**, Graphical overview of Evident usage. Sample metadata with categorical groups are used to  
77 determine differences among samples. Effect size calculation can be performed and used to  
78 generate power curves at multiple statistical significance levels and sample sizes. **b,c**  
79 Screenshots of interactive webpage for dynamic exploration of effect sizes and power analysis.  
80 Summarized effect sizes of all columns can be used to inform interactive power analysis on  
81 multiple groups (b). The underlying grouped data can be visualized with boxplots and,  
82 optionally, the raw data as scatter plots (c). Data shown is from McClorry et al. (Qiita study ID:  
83 11402)<sup>9</sup>.

84  
85 Figure 1a shows an overview of the Evident workflow. As input, Evident takes a sample  
86 metadata file and a data file. Both univariate data (in vector form such as alpha diversity) and  
87 multivariate data (as a distance matrix such as beta diversity) are supported. For univariate  
88 measures, the differences in means among groups are considered. For multivariate measures,  
89 the difference in means among within-group pairwise distances are considered. It is important to  
90 note that Evident does not perform formal hypothesis testing of community differences, only  
91 effect size calculations. While we highlight diversity differences in this work, we note that  
92 Evident can also be used for other sample-level quantitative metrics such as log-ratios<sup>10</sup>.

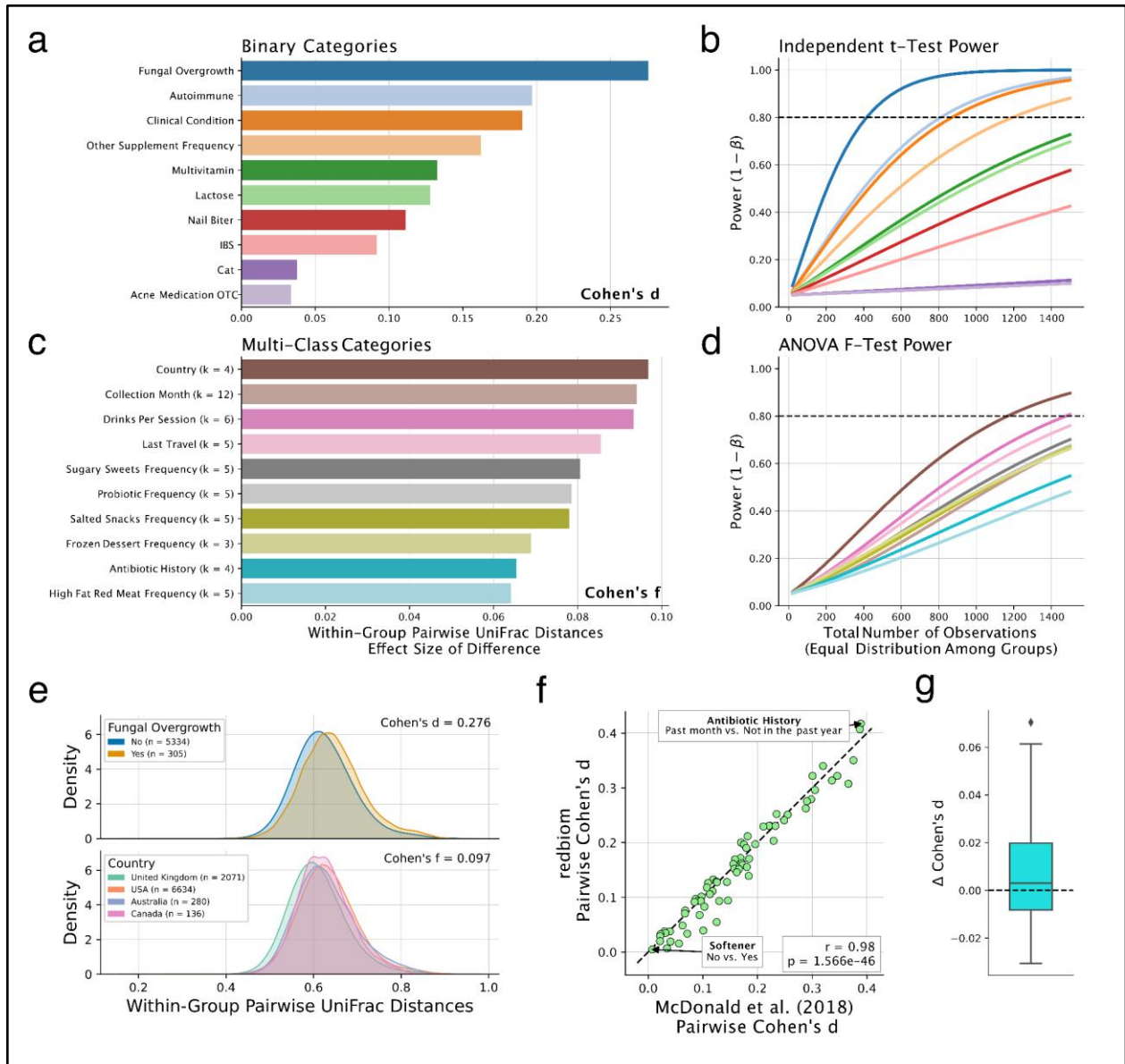
93  
94 Evident supports both binary categories and multi-class categories. For binary categories,  
95 Cohen's d is calculated between the two levels. For multi-class categories, Cohen's f is  
96 calculated among the levels<sup>11</sup>. Users can specify pairwise effect size calculations between  
97 levels of a multi-class category rather than a single group-wise effect size. Effect size  
98 calculations can be performed on multiple categories at once with simple parallelization by  
99 providing the number of CPUs to use. This architecture allows us to decrease the runtime of

100 effect size calculations for 9495 samples comprising 69 categories from over 12 minutes to 3.5  
101 minutes using 4 CPUs in parallel.

102  
103 Evident also provides an interactive component by which users can dynamically explore sample  
104 groupings. In Figure 1b,c, we show a screenshot of a web app users can access with Evident.  
105 Metadata categories are pre-sorted by effect size, allowing efficient determination of interesting  
106 categories. Power analysis is implemented dynamically - multiple categories can be visualized  
107 simultaneously for a specified significance level and number of observations.

108  
109 As a demonstration of Evident, we reprocessed 9495 samples from the AGP to compare the  
110 published effect sizes with those from a new analysis with Evident<sup>4</sup>. We downloaded the same  
111 samples from the original paper and reprocessed the data and metadata in the same manner,  
112 focusing on within-group UniFrac<sup>12</sup> distances. First, we compute the group-wise effect sizes for  
113 all valid metadata categories. The top ten binary categories and multi-class effect sizes are  
114 shown in Figure 2a,c, respectively. Using these effect sizes, we performed power analyses for  
115 each category at a significance level of 0.05 for a range of sample sizes from 20 to 1500 (Figure  
116 2b,d). We plot the distribution of the highest effect size binary and multi-class categories in  
117 Figure 2e. Finally, we compute the pairwise effect sizes as performed in the original paper to  
118 verify that Evident returns the same values. Figure 2f shows that the effect sizes map extremely  
119 closely between the published data and the newly reprocessed data. The values of effect size  
120 differences in Figure 2g are distributed around 0, indicating that there is very little difference  
121 between effect size calculations.

122



123

124

**Figure 2: Analysis of American Gut Project data**

125

**a**, Top 10 binary categories by group-wise effect size. **b**, Two-sample independent t-test power

126

analysis of selected binary category effect sizes for significance level of 0.05. **c**, Top 10 multi-

127

class categories by group-wise effect size. **d**, One-way ANOVA F-test power analysis of

128

selected multi-class category effect sizes at significance level of 0.05. **e**, Distributions of within-

129

group pairwise UniFrac distances for highest effect size binary category (top) and multi-class

130

category (bottom). **f**, Comparison of pairwise effect sizes between reprocessed data from

131 redbiom and published effect sizes from McDonald et al. **g**, Boxplot of differences in effect sizes  
132 between published and reprocessed effect sizes.

133  
134 We encourage microbiome researchers to incorporate Evident into their workflows for both  
135 reporting effect sizes of microbial community differences and planning experimental designs. In  
136 the future, we hope to enhance flexibility by including quantitative metadata categories and  
137 variable sample size power analyses.

## 138 Methods

### 139 Overview of Evident

140 Evident requires a diversity file and its associated sample metadata for a microbiome  
141 sequencing experiment. Both univariate and multivariate data are supported as an input pandas  
142 Series and scikit-bio DistanceMatrix respectively<sup>13</sup>. When evaluating multivariate differences  
143 among sample groups, Evident calculates the difference among pairwise within-group sample  
144 distances.

### 145 Effect size calculations

146 Effect size calculations are available for both binary metadata categories (e.g. Yes vs. No) or  
147 multi-class categories (e.g. diet 1, diet 2, diet 3). For binary categories, Evident calculates  
148 Cohen's d; for multi-class categories, Evident calculates Cohen's f according to a one-way  
149 ANOVA. For both types of categories, the pooled standard deviation ( $\sigma_p$ ) is computed as follows  
150 for  $G$  groups where  $n_i$  and  $s_i$  are the sample size and sample variance of group  $i$ , respectively:

151



$$\sigma_p = \sqrt{\frac{\sum_{i=1}^G (n_i - 1) s_i^2}{\sum_{i=1}^G (n_i - 1)}}$$

152

153

154 Cohen's d is calculated by the difference in means of the two groups divided by the pooled  
155 standard deviation by the following equation:

156

$$d = \frac{\mu_1 - \mu_2}{\sigma_p}$$

157

158

159 As a rule of thumb, Cohen's d values of 0.2, 0.5, and 0.8 are generally considered "small",  
160 "medium", and "large" effect sizes, respectively<sup>11</sup>.

161

162 Cohen's f is calculated by the following equation, where  $N$  is the total number of samples, and  
163  $\mu_w$  is the weighted average of all groups by sample size:

164

$$f = \frac{\sqrt{\sum_{i=1}^G (n_i/N)(\mu_i - \mu_w)^2}}{\sigma_p}$$

165

166

167 Cohen's f values of 0.1, 0.25, and 0.4 are generally considered "small", "medium", and "large"  
168 effect sizes, respectively<sup>11</sup>. For the case of two groups of equal sample size, Cohen's f is equal  
169 to Cohen's d divided by two.

170

171 In Evident, the determination of which effect size measure to use is performed automatically  
172 given the number of groups within the chosen metadata column(s). Effect size calculations are  
173 performed using NumPy<sup>14</sup> and SciPy<sup>15</sup>. Calculations in Evident assume that quantitative data is  
174 normally distributed and that population variances are homogenous among groups<sup>16</sup>.

175

176 Evident allows users to specify the maximum number of levels in a category for the category to  
177 be considered. This is useful for ignoring categories with many levels that may not be of interest  
178 (e.g. sample identifier). Additionally, one can provide a minimum number of samples in a  
179 category level so that rare groups are excluded.

## 180 Power analysis

181 With computed effect sizes, Evident is able to compute statistical power given significance  
182 level(s) and sample size(s). Additionally, Evident allows users to input a target difference (effect  
183 size numerator) to use in lieu of automatic computation. This is useful for researchers interested  
184 in designing experiments with specific effect sizes in mind<sup>7</sup>. Power analysis assumes that the  
185 number of samples is the same in each group for both of these statistical tests. These power  
186 analyses are calculated using the statsmodels package in Python<sup>17</sup>.

187

188 Notably, Evident is designed for flexibility in power analysis. Users can compute either number  
189 of observations, effect size, or statistical power given the other two variables. Additionally,  
190 Evident is designed with generating power curves in mind. For example, with a Cohen's d of  
191 0.4, a user can specify significance levels of 0.1, 0.05, and 0.01 from 20 total observations to  
192 100 in increments of 10. The statistical power will then be evaluated at each entry in the  
193 Cartesian product of these two parameter sets. These results can be directly plotted as a power  
194 curve using Evident, delineating the curves from different significance levels.

## 195 Repeated measures

196 Evident supports a limited implementation of repeated measures analysis. For datasets in which  
197 the same subject is measured more than once, statistical analysis must be performed with this

198 in mind to account for variation due to between and within subjects. Univariate data such as  
199 alpha diversity can be analyzed with repeated measures in mind by providing a mapping of  
200 sample to subject.

201  
202 For repeated measures datasets, Evident calculates eta squared ( $\eta^2$ ). This effect size can be  
203 used to calculate statistical power of a balanced repeated measures ANOVA. The number of  
204 subjects and number of measurements per subject are used to determine the total sample size.  
205 In addition to significance level, sphericity (variance between pairs of treatments<sup>18</sup>) and sample  
206 correlation parameters are used to calculate statistical power as described previously<sup>19,20</sup>. For  
207 convenience, subjects with missing values are removed.

## 208 Interactive exploration of community differences

209 The interactive visualization provided in Evident is created with Bokeh. Given microbiome data  
210 and sample metadata, Evident creates a Bokeh app that dynamically calculates effect sizes and  
211 power analysis for the chosen parameters. This view also shows the raw data values as  
212 boxplots with optional scatter points.

## 213 Analysis of AGP data

214 A sample ID list was generated from the original distance matrix used in the AGP study. 100  
215 nucleotide 16S-V4 data for these samples were downloaded from the AGP study on Qiita (study  
216 ID: 10317) using redbiom<sup>21,22</sup>. Both preparation and sample metadata were also retrieved with  
217 redbiom. Due to multiple preparations of some samples, we performed disambiguation by  
218 keeping the samples with the highest sequencing depth.

219

220 We then processed the feature table and metadata according to the original study. The original  
221 workflow used the default parameters in Deblur to remove features with fewer than 10  
222 occurrences in the data<sup>23</sup>. Because Qiita does not perform this filtering by default, we performed  
223 this filtering manually. To remove sequences associated with sample bloom, we performed  
224 bloom filtering<sup>24</sup>. We then rarefied the feature table to 1250 sequences as in the original  
225 analysis.

226  
227 We processed the sample metadata in accordance with the original study. Because of  
228 differences in self-reporting protocols from 2018, metadata categories associated with reported  
229 Vioscreen responses as well as those associated with alcohol consumption were removed. The  
230 following categories were removed due to mismatches in sample metadata: roommates,  
231 allergies, age\_cat, bmi\_cat, longitude, latitude, elevation, height\_cm, collection\_time, and  
232 center\_project\_name. Only the top four annotated countries were considered - US, UK,  
233 Australia and Canada. All other countries were ignored. Overall, 69 metadata categories  
234 common to both the original data and redbiom data were used for further analysis.

235  
236 Sequences from the feature table were placed into a 99% Greengenes<sup>25</sup> insertion reference tree  
237 using SEPP<sup>26</sup>. We then used unweighted UniFrac to generate a sample-by-sample distance  
238 matrix<sup>27</sup>. This distance matrix was used as input to Evident along with the disambiguated,  
239 processed sample metadata.

240  
241 We used effect\_size\_by\_category to calculate the whole-group effect sizes for each column in  
242 the metadata and pairwise\_effect\_size\_by\_category to calculate the group-pairwise effect sizes  
243 for multi-class categories. For each whole-group effect size, we computed a power analysis for  
244 alpha values of 0.01, 0.05, and 0.1. Power was calculated on total sample size values from 20  
245 to 1500 in increments of 40 samples. Evident analyses were performed in parallel on a high

246 performance computing environment. Group-wise and pairwise effect size calculations both took  
247 under 4 minutes for 94 metadata categories on 9495 samples using 4 CPUs (we note the AGP  
248 paper used n=9511 but operated at 125nt; we observe a slightly reduced number of samples at  
249 100nt). We also benchmarked group-wise effect size calculations using only a single CPU as  
250 comparison; this process took 12.4 minutes - meaning the parallelization decreased runtime by  
251 approximately 3.5x. Power analysis calculation took 2 minutes for 94 categories using 8 CPUs  
252 in parallel.

## 253 Code availability

254 The latest version of Evident is available at <https://github.com/biocore/evident> under the BSD-3  
255 license. Evident is installable from PyPI both as a standalone Python 3 package and a QIIME 2  
256 plugin. The scripts used to download and analyze AGP data as well as the processed Evident  
257 results are available at <https://github.com/knightlab-analyses/evident-analyses>. Analysis of AGP  
258 data in this study was performed with Evident version 0.2.0.

## 259 Data availability

260 Data for the demonstration in Figure 1 were downloaded from Qiita (study ID: 11402)<sup>9</sup> at 90  
261 nucleotides using the deblur<sup>23</sup> pipeline. AGP data were downloaded from Qiita (study ID: 10317)  
262 using redbiom with context "Deblur\_2021.09-Illumina-16S-V4-100nt-50b3a2". The original  
263 pairwise effect sizes, sample metadata, and unweighted UniFrac distance matrix were  
264 downloaded from the original McDonald et al. study for comparison.

## 265 Contributions

266 G.R., A.G, D.M., & R.K. conceived the idea for the software and study. G.R., A.G., D.M., Y.V.B  
267 & L.J. developed the software. G.R., D.M., & A.G. conducted the analysis of the AGP data.  
268 B.N. assisted with AGP metadata curation. A.H.D., Y.V.B., D.M., & D.H. reviewed the software  
269 code and provided valuable feedback and bug reports. C.C., L.J., & S.P. contributed to the  
270 statistical computation used in Evident. All authors contributed to and reviewed the final  
271 manuscript.

## 272 Acknowledgements

273 We would like to thank the members of the Knight Lab for feedback on the scope and details of  
274 Evident. We thank Jamie Morton for valuable discussions about effect size. This work was  
275 supported in part by the Alfred P. Sloan foundation (G-2017-9838), NIH-NIDDK  
276 (P01DK078669), NIH-NCI (U24CA248454), and NIH (1DP1AT010885, U19AG063744).  
277 Research of S.P. was funded by the intramural research program of *Eunice Kennedy Shriver*  
278 National Institute of Child Health and Human Development, NIH.

## 279 References

- 280 1. Sullivan, G. M. & Feinn, R. Using Effect Size—or Why the P Value Is Not Enough. *J. Grad.*  
281 *Med. Educ.* **4**, 279–282 (2012).
- 282 2. Baguley, T. Standardized or simple effect size: What should be reported? *Br. J. Psychol.* **100**,  
283 603–617 (2009).
- 284 3. Cohen, J. Statistical Power Analysis. *Curr. Dir. Psychol. Sci.* **1**, 98–101 (1992).
- 285 4. McDonald, D. *et al.* American Gut: an Open Platform for Citizen Science Microbiome

- 286 Research. *mSystems* **3**, e00031-18.
- 287 5. The Environmental Determinants of Diabetes in the Young (TEDDY) Study. *Ann. N. Y. Acad.*  
288 *Sci.* **1150**, 1–13 (2008).
- 289 6. Vartiainen, E. *et al.* Cardiovascular risk factor changes in Finland, 1972–1997. *Int. J.*  
290 *Epidemiol.* **29**, 49–56 (2000).
- 291 7. Casals-Pascual, C. *et al.* Microbial Diversity in Clinical Microbiome Studies: Sample Size and  
292 Statistical Power Considerations. *Gastroenterology* **158**, 1524–1528 (2020).
- 293 8. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data science  
294 using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
- 295 9. McClorry, S. *et al.* Anemia in infancy is associated with alterations in systemic metabolism  
296 and microbial structure and function in a sex-specific manner: an observational study. *Am. J.*  
297 *Clin. Nutr.* **108**, 1238–1248 (2018).
- 298 10. Morton, J. T. *et al.* Establishing microbial composition measurement standards with  
299 reference frames. *Nat. Commun.* **10**, 2719 (2019).
- 300 11. Cohen, J. *Statistical power analysis for the behavioral sciences*. (L. Erlbaum Associates,  
301 1988).
- 302 12. McDonald, D. *et al.* Striped UniFrac: enabling microbiome analysis at unprecedented  
303 scale. *Nat. Methods* **15**, 847–848 (2018).
- 304 13. McKinney, W. *Data Structures for Statistical Computing in Python*. 6 (2010).
- 305 14. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (2020).
- 306 15. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python.  
307 *Nat. Methods* **17**, 261–272 (2020).
- 308 16. Li, J. C.-H. Effect size measures in a two-independent-samples case with nonnormal  
309 and nonhomogeneous data. *Behav. Res. Methods* **48**, 1560–1574 (2016).
- 310 17. Seabold, S. & Perktold, J. *Statsmodels: Econometric and Statistical Modeling with*  
311 *Python*. in 92–96 (2010). doi:10.25080/Majora-92bf1922-011.

- 312 18. Quené, H. & van den Bergh, H. On multi-level modeling of data from repeated measures  
313 designs: a tutorial. *Speech Commun.* **43**, 103–121 (2004).
- 314 19. Guo, Y., Logan, H. L., Glueck, D. H. & Muller, K. E. Selecting a sample size for studies  
315 with repeated measures. *BMC Med. Res. Methodol.* **13**, 100 (2013).
- 316 20. Vonesh, E. F. & Schork, M. A. Sample Sizes in the Multivariate Analysis of Repeated  
317 Measurements. *Biometrics* **42**, 601–610 (1986).
- 318 21. Gonzalez, A. *et al.* Qiita: rapid, web-enabled microbiome meta-analysis. *Nat. Methods*  
319 **15**, 796–798 (2018).
- 320 22. McDonald, D. *et al.* redbiom: a Rapid Sample Discovery and Feature Characterization  
321 System. *mSystems* (2019) doi:10.1128/mSystems.00215-19.
- 322 23. Amir, A. *et al.* Deblur Rapidly Resolves Single-Nucleotide Community Sequence  
323 Patterns. *mSystems* **2**, e00191-16 (2017).
- 324 24. Amir, A. *et al.* Correcting for Microbial Blooms in Fecal Samples during Room-  
325 Temperature Shipping. *mSystems* **2**, e00199-16.
- 326 25. McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological  
327 and evolutionary analyses of bacteria and archaea. *ISME J.* **6**, 610–618 (2012).
- 328 26. Mirarab, S., Nguyen, N. & Warnow, T. SEPP: SATé-enabled phylogenetic placement.  
329 *Pac. Symp. Biocomput. Pac. Symp. Biocomput.* 247–258 (2012)  
330 doi:10.1142/9789814366496\_0024.
- 331 27. Lozupone, C. & Knight, R. UniFrac: a New Phylogenetic Method for Comparing Microbial  
332 Communities. *Appl. Environ. Microbiol.* **71**, 8228–8235 (2005).