## Supplementary Methods

**Ovalbumin (OVA)-driven CD4$^+$ T cell cultures and intracellular cytokine staining**

Mesenteric lymph nodes cells from reconstituted and OVA-sensitized GF mice were labeled with the Violet CellTrace proliferative dye (Invitrogen; Grand Island, NY) and cultured with 200µg/ml OVA and 250pg/ml IL-2 for 72 hours. During the last 4 hours, cultured cells were stimulated with PdBU (500 ng/ml; Sigma-Aldrich, St. Louis MO) and Ionomycin (500 ng/ml; Sigma-Aldrich) in the presence of Brefeldin A (1µg/ml; BD Biosciences – San Jose, CA). Cells were stained with the following conjugated antibodies: CD3 (145-2C11), CD4 (RM4-5), IL-4 (11B11) and IFN-$\gamma$ (XMG1.2) (eBioscience, San Diego, CA). Intracellular cytokines were detected in Violet CellTrace$^+$ proliferating CD3$^+$CD4$^+$ T cells by using Cytofix/Cytoperm (BD Biosciences) buffers, according to the manufacturer's instructions. Stained cells were analyzed on a LSRII Fortessa cytometer (BD Biosciences) and data processed using Flowjo (Tree Star; Ashland, OR).

**PhyloChip$^{TM}$ data analysis**

*Pre-processing and Data Reduction.* Fluorescent images were captured with the GeneChip Scanner 3000 7G (Affymetrix, Santa Clara, CA). An individual array feature occupied approximately 8x8 pixels in the image file corresponding to a single probe 25mer on the surface. To calculate the summary intensity for each feature on each array, the central 9 pixels of individual features were ranked by intensity and the 75% percentile was used. Probe intensities were background-subtracted and scaled to the PhyloChip™ Control Mix. Array fluorescence intensity was collected as integer values ranging from 0 to 65,536 ($2^{16}$). Fluorescence intensities for sets of probes

54 complementing an operational taxonomic unit (OUT) were averaged after discarding the

55 highest and lowest and the mean was $\log_2$ transformed into numbers ranging from 0 to

56 16. For compatibility with some statistical operations, the scores were multiplied by

57 1000 then rounded, allowing a range of integers from 0 to 16,000. These values are

58 referred to as the hybridization score (HybScore). For the complete distribution see

59 Hazen et al, Supplemental Information [1]. The data was reduced to consider the

60 bacterial taxa deemed present as described in Hazen *et al*. [1]. Taxa were filtered to

61 those present in the majority of samples of at least one of the experimental groups and

62 rank-normalized such that taxa in each are represented by their ranked HybScore within

63 that sample only (rank 1 represents the lowest HybScore in that sample).

64 *Sample-to-Sample Distance Function*. All profiles were inter-compared in a pair-wise

65 fashion to determine a dissimilarity score and results were stored as a distance matrix.

66 The Weighted Unifrac distance measure was chosen because it utilizes the

67 phylogenetic distance between OTUs as well as the abundance of those OTUs to

68 compute a community-wide dissimilarity between any pair of profiles [2, 3]. Similar

69 biological samples produce small Weighted Unifrac dissimilarity scores. When

70 comparing the presence or absence of taxa between profiles, the Unweighted Unifrac

71 distance measure was utilized.

72 *Statistical Analysis, Ordination, Clustering, and Classification Methods.* The differences

73 between the microbial communities (the entire number of OTUs detected in any one

74 comparison group versus another) was determined by the Adonis test, which is a

75 permutation test based on a dissimilarity matrix, in this case measured by weighted

76 UniFrac. Because the Adonis test considers the multidimensional structure of the data

77 (e.g., compares entire microbial communities), it does not involve multiple hypotheses

78 testing for each microbial taxon found within those communities.

79 Taxa found increased in their ranked HybScore in one category compared to the

80 alternate categories were identified using the Kruskal-Wallis (KW) test. The aim of a KW

81 filter in the context of this analysis was to reduce the dimensionality of the dataset, and

82 demonstrate that this reduced set of OTUs could still effectively discriminate between

83 samples in terms of their microbial community structures by the ordination and

84 clustering methods listed below.

85 Two-dimensional ordinations and hierarchical clustering maps of the samples in the

86 form of dendrograms were created to graphically summarize the inter-sample

87 relationships. To create dendrograms, the samples from the distance matrix are

88 clustered hierarchically using the average-neighbor (HC-AN) method [4]. Non-Metric

89 Multidimensional Scaling (NMDS) was employed to visualize relationships between

90 samples by two-dimensional ordination plotting [5]. Ordination points are colored by

91 highlighted groupings. Lists of significant taxa whose abundance characterizes each

92 class is performed using Prediction Analysis for Microarrays (PAM), a classifier

93 (supervised machine learning) based method that utilizes a nearest shrunken centroid

94 method [6].

95 *Phylogenetic Tree Visualization*. Bacterial families with OTUs found by the KW test to

96 be differentially abundant between two comparison groups (*e.g.* allergen sensitized WT

97 versus *Il4raF709* mice) were identified, and the one OTU with the greatest difference

98 between the two group means from each family was selected. For those families

99 containing OTUs with both higher and lower abundance scores between the two

100 comparison groups, two OTUs were selected. A phylogenetic tree was constructed

101 using FastTree, which was built using one representative 16S ribosomal DNA (rDNA)

102 gene sequence from each of the OTUs selected from the Greengenes multiple

103 sequence alignment [7, 8]. The Tree was displayed with iTOL software [9].

104 **16S rDNA sequencing methods and data analysis**

105 *Summary of methodology.* The microbial community structure in each stool sample was

106 assessed by 16S amplicon sequencing on the Roche 454 platform. Sequencing data

107 was processed through a bioinformatics pipeline to obtain distributions of OTUs for each

108 sample. We tested differences in overall microbial community structure between stool

109 samples from different groups using the Dirichlet Multinomial model and a likelihood

110 ratio test [10, 11]. We used hierarchical clustering with the Bray-Curtis (BC) dissimilarity

111 measure to visualize the differences between the distributions of OTUs in samples [12].

112 The BC measure quantifies the difference between a pair of ecosystems based on the

113 species or OTU composition of samples. A BC value of zero indicates identical OTU

114 distributions; a BC value of one indicates no overlap in the OTUs present in the pair of

115 samples. We used a bootstrapping procedure to estimate 95% confidence intervals on

116 BC measures, and thus evaluate the reproducibility of sample clusterings. Results were

117 visualized with a dendrogram constructed using the bootstrapped values. We also used

118 the UniFrac measure with Principal Coordinates Analysis (PCoA) to visualize

119 differences between microbial communities in samples; this measure takes into account

120 phylogenetic relationships among sequences and does not require clustering

121 sequences into OTUs [2, 3]. Individual OTUs that discriminate between different groups

122 were determined using a Random Forests supervised machine learning approach [13, 14].

123  *16S rDNA Amplicon sequencing.* DNA pyrosequencing was performed by the Human

124  Genome Sequencing Center at Baylor College of Medicine following protocols

125  benchmarked for the Human Microbiome Project.  The V3-V5 hypervariable regions of

126  the 16S rRNA gene were amplified using primer 357F (5'-CCTACGGGAGGCAGCAG-

127  3') modified with the addition of the 454 FLX-titanium adaptor "B" sequence

128  (5'CCTATCCCCTGTGTGCCTTGGCAGTCTCAG-3') and primer 926R (5'-

129  CCGTCAATTCMTTTRAGT-3') modified with the addition of unique 6-8 nucleotide

130  barcode sequences and the 454 FLX-titanium adaptor "A" sequence (5'-

131  CCATCTCATCCCTGCGTGTCTCCGACTCAG-3').  Barcode and adaptor sequences

132  are                          found                          at

133  http://www.hmpdacc.org/doc/HMP_MDG_454_16S_Protocol_V4_2_102109.pdf.   PCR

134  amplification was performed on 2 uL of DNA template in a total volume of 25 uL

135  containing 1x AccuPrime Buffer II (Invitrogen Corp., Carlsbad, CA), 320 uM of each

136  primer, and 0.03 U/uL AccuPrime High Fidelity *Tag* polymerase.   Reactions were

137  heated at 95$^{o}$C for 2 min followed by 30 cycles of 95$^{o}$C for 20 sec, 50$^{o}$C for 30 sec, and

138  72$^{o}$C for 5 min.  The concentration of amplicons in each reaction was determined in

139  triplicate using the PicoGreen fluorescent assay (Invitrogen Corp.) and amplicons were

140  pooled before being sequenced via a multiplexed 454-FLX-titanium pyrosequencing run

141  according to manufacturer's specifications.

142  *Bioinformatics for 16S data.* Sequences were pre-processed using custom scripts and

143  the packages mothur, CloVR, and QIIME [15-17]. Filtering criteria were: no ambiguous

144  bases, maximum homopolymer length of 8, 1 base difference allowed for barcode

145  matches, and 2 base differences allowed for primer matches. Each sample had

146  approximately 3000 reads after filtering. Sequences were trimmed based on a minimum

147  average quality score of 35 over a window of length 50 nt, and clustered into OTUs with

148  a similarity threshold of 95%.

149  *Testing for differences in OTU distributions between groups.* OTU relative abundances

150  were assumed to follow the Dirichlet Multinomial (DM) distribution[1,2]. To test for

151  differences in overall community structure between two groups, denoted A and B, we

152  used a likelihood ratio test:

153  $S = -2\ln\{\mathrm{P}(X_A, X_B \mid M_{A+B}) / [\mathrm{P}(X_A \mid M_A)\mathrm{P}(X_B \mid M_B)]\}$

154  Here, $X_A$ and $X_B$ represent the set of vectors of OTU counts for groups A and B

155  respectively. $M_{A+B}$ represents the DM model estimated from the combined groups, and

156  $M_A$ and $M_B$ the corresponding DM models estimated from the separate groups. DM

157  parameters were estimated using the Maximum Likelihood method. The *S* statistic

158  asymptotically follows a $\chi^2$ distribution with degrees of freedom equal to the number of

159  OTUs in the samples.

160  *Clustering and visualizing samples.* Bootstrapping was performed to standardize the

161  effects of differing numbers of sequencing reads between samples, and to obtain

162  estimates of the variability of dissimilarity measures between samples. For each pair of

163  samples *i* and *j*, *m* reads were drawn independently and with replacement, and the

164  Bray-Curtis Dissimilarity measure[3] was calculated between the bootstrapped reads. We

165  set *m* equal to the median number of sequencing reads over all samples, and repeated

166  the bootstrapping procedure on each pair of samples 10,000 times. The 95%

167  confidence interval for each sample pair was then estimated from the empirical

168     distribution of values. An average linkage dendrogram was constructed using the 95[th]

169     centile values between nodes.

170     *Finding OTUs that discriminate between groups.* To find OTUs discriminating between

171     groups, we used Random Forests[5] (RF) with a wrapper feature based method as

172     implemented in the Boruta[6] package. Briefly, RF is an ensemble based classification

173     method that uses multiple weak classifier decision trees. An importance measure is

174     calculated for each feature (OTU) based on the loss of accuracy in classification. The

175     statistical significance of the importance measure is determined using a permutation

176     based method.

177     **16S rDNA Pyrosequencing Analysis: Results**

178     *Comparisons between sensitized and sham sensitized Il4raF709 mice.* We assessed

179     the difference in overall microbial community structure among stool samples from

180     *Il4raF709* homozygous mutant mice sensitized with OVA or sham sensitized with PBS.

181     The distributions of OTUs differed significantly between the groups (Dirichlet

182     Multinomial model, *p*-value = < $10^{-20}$). BC dissimilarity dendrograms and UniFrac PCoA

183     plots visualizing differences in overall microbial community structure between groups

184     showed overall separation between the two groups, although two mutant PBS samples

185     were close to outlying mutant OVA samples (**Figure E1A, B**). Several bacterial families

186     and genera were found to discriminate between the groups using a supervised machine

187     learning based method, including OTUs classifying to the genera Clostridium,

188     Bacteroides, Alistipes and Streptococcus (**Table E3**).

189     *Comparisons between unsensitized WT versus Il4raF709 mutant mice.* We assessed

190     the difference in overall microbial community structure between stool samples from

191    unsensitized WT and *Il4raF709* homozygous mutant littermate mice, and found that the

192    distributions of OTUs differed significantly between the two groups (Dirichlet Multinomial

193    model *P* value = $7 \times 10^{-11}$). We visualized differences in overall microbial community

194    structure between the two groups using a dendrogram with the BC measure (**Figure**

195    **E4A**) and a UniFrac PCoA plot (**Figure E4B**). Consistent across both techniques, the

196    samples from the *Il4raF709* homozygous mutant mice overall clustered together,

197    although few WT samples clustered with outlying *Il4raF709* samples. These findings

198    suggest that differences between the two groups were relatively subtle and not well-

199    visualized using a dimensionality reduction method. To explore differences in individual

200    OTUs, we used a supervised machine learning based method, and found differences in

201    several OTUs, including those classifying to the genera Helicobacter, Clostridium,

202    Lactobacillus and Odoribacter (data not shown).

203    *Comparisons between WT and Il4raF709 mutant sensitized mice*. We assessed the

204    difference in overall microbial community structure among stool samples from *Il4raF709*

205    homozygous mutant mice and WT controls sensitized with OVA. The distributions of

206    OTUs differed significantly between each group (Dirichlet Multinomial model, *P* value =

207    $< 10^{-20}$). BC dissimilarity dendrograms and UniFrac PCoA plots visualizing differences in

208    overall microbial community structure between groups showed clear separation

209    between the WT and mutant OVA groups (**Figure E5A, B**). Several bacterial families

210    and genera were found to optimally discriminate between the groups using a supervised

211    machine learning-based method, including Alistipes, Clostridium, Anaeroplasma,

212    Lachnobacterium, and Bacteroides (**Table E9**).

213 *Assessing recolonization of WT GF mice by flora of OVA-sensitized WT versus mutant*

214 *mice.* We assessed the difference in overall microbial community structure between

215 stool samples from the two groups of recipient mice collected 8 weeks after colonization

216 (at the end of the OVA sensitization period). The distributions of OTUs in stool samples

217 from the group of mice receiving donor microbiota from allergen-sensitized mutant mice

218 differed significantly from those of the group receiving donor microbiota from WT mice

219 (Dirichlet Multinomial model[1,2] *P* value < $10^{-20}$). We visualized differences between

220 samples using a dendrogram with the BC measure (**Figure 8A**).

221 The samples from the group of mice that received donor microbiota from WT mice all

222 clustered tightly, and clustered with the respective donor sample. The samples from the

223 group of mice that received donor microbiota from allergen-sensitized *Il4raF709* mutant

224 mice also clustered closely with one another, but were distinct from those of WT flora

225 recipients. The donor sample from allergen-sensitized mutant mice essentially clustered

226 separately, but was closer to its respective recipient samples in aggregate than it was to

227 the other samples.

228 **Supplementary Figure Legends**

229 **Figure E1.** The microbial signature and dysbiosis associated with the allergen

230 sensitization of *Il4raF709* mice is reproduced by 16S rDNA pyrosequencing. **A.**

231 Agglomerative clustering of fecal samples from OVA- and sham PBS sensitized

232 *Il4raF709* mice based on 16S OTUs Samples were clustered based on bootstrapped

233 BC dissimilarity values computed on the relative abundances of OTUs in each sample.

234 BC = 0 indicates identical microbial communities; BC = 1 indicates communities with no

235 overlapping OTUs. Heights of lines on the dendrogram indicate bootstrapped BC values
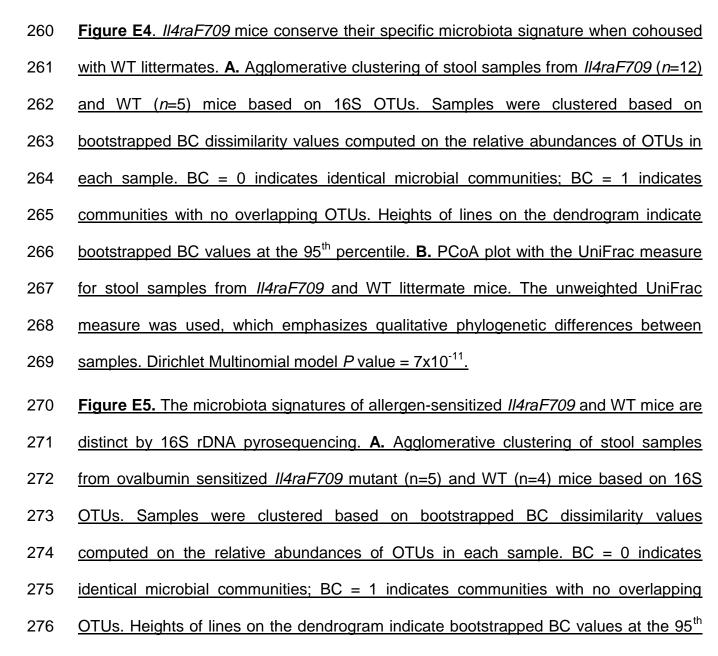
236 at the 95<sup>th</sup> percentile. **B.** Principal Coordinate Analysis (PCoA) plot with the UniFrac

237 measure for stool samples from OVA and sham PBS sensitized *Il4raF709* mice. The

238 unweighted UniFrac measure was used, which emphasizes qualitative phylogenetic

239 differences between samples. *n*=5 for the *Il4raF709* OVA group versus 9 mice for the

240 PBS sham sensitized *Il4raF709* group. Dirichlet Multinomial model, *p*-value = $< 10^{-20}$.

241 **Figure E2.** $T_R$ cell-treatment resets the microbiota of allergen-sensitized *Il4raF709* mice

242 into a new baseline distinct from that of sham sensitized and treated control *Il4raF709*

243 mice. **A**. NMDS based on Weighted Unifrac distance between samples of PBS-

244 sensitized mice (*n*=4) versus those of $T_R$-cell treated (*n*=5), OVA-sensitized mice, based

245 on the 786 taxa whose abundance was significantly different between groups using the

246 Kruskal-Wallis (KW) test. **B**. Hierarchical Clustering based on Weighted Unifrac

247 distance between samples. **C**. Nearest shrunken centroid analysis of OTUs that best

248 characterize the difference between allergen-sensitized versus tolerant groups. **D**.

249 Representation of the abundance of the OTUs identified by the nearest shrunken

250 centroid analysis using the PAM method.

251 **Figure E3**. The microbiota of sham-sensitized *Il4raF709* and WT mice are distinct. **A**.

252 NMDS based on Weighted Unifrac distance between samples of PBS-sensitized

253 *Il4raF709* mice (*n*=4) versus those of WT, PBS-sensitized mice (*n*=4), based on the 813

254 taxa whose abundance was significantly different between groups using the KW test. **B**.

255 Hierarchical Clustering based on Weighted Unifrac distance between samples. **C**.

256 Nearest shrunken centroid analysis of OTUs that best characterize the difference

257 between allergen-sensitized versus tolerant groups. **D**. Representation of the

258 abundance of the OTUs identified by the nearest shrunken centroid analysis using the

259 PAM method.

260 **Figure E4**. *Il4raF709* mice conserve their specific microbiota signature when cohoused

261 with WT littermates. **A.** Agglomerative clustering of stool samples from *Il4raF709* ($n$=12)

262 and WT ($n$=5) mice based on 16S OTUs. Samples were clustered based on

263 bootstrapped BC dissimilarity values computed on the relative abundances of OTUs in

264 each sample. BC = 0 indicates identical microbial communities; BC = 1 indicates

265 communities with no overlapping OTUs. Heights of lines on the dendrogram indicate

266 bootstrapped BC values at the 95[th] percentile. **B.** PCoA plot with the UniFrac measure

267 for stool samples from *Il4raF709* and WT littermate mice. The unweighted UniFrac

268 measure was used, which emphasizes qualitative phylogenetic differences between

269 samples. Dirichlet Multinomial model $P$ value = $7 \times 10^{-11}$.

270 **Figure E5.** The microbiota signatures of allergen-sensitized *Il4raF709* and WT mice are

271 distinct by 16S rDNA pyrosequencing. **A.** Agglomerative clustering of stool samples

272 from ovalbumin sensitized *Il4raF709* mutant (n=5) and WT (n=4) mice based on 16S

273 OTUs. Samples were clustered based on bootstrapped BC dissimilarity values

274 computed on the relative abundances of OTUs in each sample. BC = 0 indicates

275 identical microbial communities; BC = 1 indicates communities with no overlapping

276 OTUs. Heights of lines on the dendrogram indicate bootstrapped BC values at the 95[th]

277 percentile. **B.** PCoA plot with the UniFrac measure for stool samples from ovalbumin

278 sensitized mutant and WT mice. The unweighted UniFrac measure was used, which

279 emphasizes qualitative phylogenetic differences between samples. Dirichlet Multinomial

280 model, *p*-value = $< 10^{-20}$.

281    **Figure E6.** The microbiota signatures of allergen- (OVA +OVA/SEB) sensitized

282    *Il4raF709* and WT mice are distinct. **A**. NMDS based on Weighted Unifrac distance

283    between samples of OVA/SEB-sensitized WT (*n*=6) versus OVA+OVA/SEB sensitized

284    *Il4raF709* mice (*n*=9), based on the 430 taxa whose abundance was significantly

285    different between groups using the KW test. **B**. Hierarchical Clustering based on

286    Weighted Unifrac distance between samples. **C**. Nearest shrunken centroid analysis of

287    OTUs that best characterize the difference between the groups. **D**. Representation of

288    the abundance of the OTUs identified by the nearest shrunken centroid analysis using

289    the PAM method. **E**. Venn diagram showing the abundance levels of different OTUs in

290    relation to the sensitization state of WT and *Il4raF709* mice. The labels define the

291    abundance states of sets of OTUs in relation to specific comparison groups, *e.g.* F709

292    OVA+OVA/SEB<F709 PBS identifies those OTUs that are less abundant in allergen

293    sensitized (with OVA or with OVA/SEB) *Il4raF709* mice as compared to sham (PBS)

294    sensitized mice. The number of OTUs thus identified is indicated in parentheses.

295    Spheres indicate intersections between two sets, while the colored webs show which

296    intersection of sets form the spheres. **F**. Contingency table representation of the results

297    shown in the Venn diagram. P<0.0001 by the $X^2$ test (excluding the 2993 OTUs that did

298    not change upon sensitization in both WT and *Il4raF709* mice).

299    **Table E1**. Annotations of Prediction Analysis for Microarrays (PAM)-selected bacterial

300    taxa that discriminate between sham and allergen (OVA and OVA/SEB-sensitized)

301    *Il4raF709* mice (see **Figure 3C, D**).

302 **Table E2.** Annotations of bacterial taxa showing significantly different abundances

303 between sham and allergen-sensitized *Il4raF709* mice, as shown in the phylogenetic

304 tree in **Figure 4**.

305 **Table E3.** Annotations of bacterial genera that optimally discriminate between sham

306 and allergen-sensitized *Il4raF709* mice, as revealed by 16S rDNA pyrosequencing

307 (**Figure E1)** and determined using the Random Forest machine learning method.

308 **Table E4**. Annotations of PAM-selected bacterial taxa that discriminate between

309 allergen (OVA- and OVA/SEB)-sensitized versus $T_R$-cell treated and OVA-sensitized

310 mice. The taxa selected correspond to those graphically presented in **Figure 5C, D**.

311 **Table E5.** Annotations of bacterial taxa showing significantly different abundances

312 between sham and allergen-sensitized *Il4raF709* mice, as shown in the phylogenetic

313 tree in **Figure 6**.

314 **Table E6**. Annotations of PAM-selected bacterial taxa that discriminate between $T_R$ cell-

315 treated, allergen (OVA)- and sham-sensitized *Il4raF709* mice. The taxa selected

316 correspond to those graphically presented in **Figure E2C, D**.

317 **Table E7**. Annotations of PAM-selected bacterial taxa that discriminate between sham-

318 and allergen (OVA/SEB)-sensitized WT mice. The taxa selected correspond to those

319 graphically presented in **Figure E3C, D**.

320 **Table E8**. Annotations of PAM-selected bacterial taxa that discriminate between

321 allergen (OVA/SEB)-sensitized WT *and Il4raF709* mice. The taxa selected correspond

322 to those graphically presented in **Figure 7C, D**.

323   **Table E9**. Annotations of bacterial genera that optimally discriminate between allergen-

324   sensitized WT versus *Il4raF709* mice, as revealed by 16S rDNA pyrosequencing

325   (**Figure E5)** and determined using the Random Forest machine learning method.

326   **Supplementary References**

327   1.   Hazen TC, Dubinsky EA, DeSantis TZ, Andersen GL, Piceno YM, Singh N, et al.

328        Deep-sea oil plume enriches indigenous oil-degrading bacteria. Science 2010;

329        330:204-8.

330   2.   Lozupone C, Hamady M, Knight R. UniFrac--an online tool for comparing

331        microbial community diversity in a phylogenetic context. BMC Bioinformatics

332        2006; 7:371.

333   3.   Lozupone C, Knight R. UniFrac: a new phylogenetic method for comparing

334        microbial communities. Appl Environ Microbiol 2005; 71:8228-35.

335   4.   Legendre P, Legendre L. Numerical ecology. 2nd English ed. Amsterdam ; New

336        York: Elsevier; 1998.

337   5.   Shepard RN. Multidimensional scaling, tree-fitting, and clustering. Science 1980;

338        210:390-8.

339   6.   Tibshirani R, Hastie T, Narasimhan B, Chu G. Diagnosis of multiple cancer types

340        by shrunken centroids of gene expression. Proc Natl Acad Sci U S A 2002;

341        99:6567-72.

342   7.   Price MN, Dehal PS, Arkin AP. FastTree 2--approximately maximum-likelihood

343        trees for large alignments. PLoS One 2010; 5:e9490.

344    8.      DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, et al.

345            Greengenes, a chimera-checked 16S rRNA gene database and workbench

346            compatible with ARB. Appl Environ Microbiol 2006; 72:5069-72.

347    9.      Letunic I, Bork P. Interactive Tree Of Life (iTOL): an online tool for phylogenetic

348            tree display and annotation. Bioinformatics 2007; 23:127-8.

349    10.    Mosimann JE. On the compound multinomial distribution, the multivariate $\beta$-

350            distribution, and correlations among proportions. Biometrika 1962; 49:65-82.

351    11.    Tvedebrink T. Overdispersion in allelic counts and theta-correction in forensic

352            genetics. Theor Popul Biol 2010; 78:200-10.

353    12.    Bray JR, Curtis JT. An Ordination of the Upland Forest Communities of Southern

354            Wisconsin. Ecological Monographs 1957; 27:325-49.

355    13.    Breiman L. Random Forests. Learning Machines 2001; 45:5-32.

356    14.    Kursa MB, Rudnicki WR. Feature selection with the boruta package. j. Stat. Soft.

357            2010; 36:1-13.

358    15.    Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK,

359            et al. QIIME allows analysis of high-throughput community sequencing data. Nat

360            Methods 2010; 7:335-6.

361    16.    Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al.

362            Introducing mothur: open-source, platform-independent, community-supported

363            software for describing and comparing microbial communities. Appl Environ

364            Microbiol 2009; 75:7537-41.

365    17.    Angiuoli SV, Matalka M, Gussman A, Galens K, Vangala M, Riley DR, et al.

366          CloVR: a virtual machine for automated and portable sequence analysis from the

367          desktop using cloud computing. BMC Bioinformatics 2011; 12:356.

368