

1 **Supplementary tables S1-S9 are available at:**

2 [https://figshare.com/articles/Supplementary\\_Tables\\_-\\_Levade\\_et\\_al\\_2020/12440417](https://figshare.com/articles/Supplementary_Tables_-_Levade_et_al_2020/12440417)

3

#### 4 **Data availability**

5 After removal of human reads (Supplementary Methods), the sequence data has been deposited in  
6 NCBI under BioProject PRJNA608678.

7

#### 8 **Supplementary Methods**

##### 9 **Clinical outcomes and metadata collection**

10 All study participants lived in Dhaka city, Bangladesh. Multiple household contacts could be  
11 enrolled from one household, and only one index case was enrolled per household. Household  
12 contacts were defined as individuals who shared the same cooking pot with an index patient for  
13 three or more days. Rectal swab samples were obtained from contacts beginning the day after  
14 enrollment of the index case. During the follow up period for contacts, two sets of rectal swabs  
15 were collected. The first set were collected during daily home visits, when report of symptoms  
16 was also collected, and was used for *V. cholerae* culture. This data was used for the classification  
17 of clinical outcomes. The second set of rectal swabs collected over the first ten days, and day 30,  
18 also at the home visit, were used for DNA extraction to conduct the microbiome analyses.

19 Infection status was determined using the first set of rectal swab samples, serologic data from the  
20 blood draw, and report of symptoms collected during the observation period. Blood samples  
21 collected during the same days from contacts were drawn at the International Center for Diarrheal  
22 Disease Research, Bangladesh, in Dhaka, Bangladesh. Symptomatic *V. cholerae* infection was  
23 defined by a contact reporting diarrhea within 3 days of a rectal swab positive for *V. cholerae*

24 during the follow-up home visits, or by a four-fold increase in vibriocidal titer with diarrhea  
25 during the follow-up period. Two symptomatic contacts were rectal swab negative for *V.*  
26 *cholerae* but showed a four-fold increase in vibriocidal titer and reported diarrhea, and then were  
27 determined to be infected based on our inclusion criteria. *V. cholerae* infection (rectal swab  
28 positive) was defined as asymptomatic if no diarrhea was reported. All contacts in our study that  
29 developed infection had had the same serotype of *V. cholerae* as the infected household case  
30 (either Inaba or Ogawa).

### 31 **Midani and Expanded cohort description**

32 We used metagenomic reads from 65 of the 76 contacts from a previously published study  
33 by our collaborators, and these contacts are referred to as the Midani 2018 cohort (1). Eleven of  
34 the contacts from this study had insufficient DNA to perform metagenomic sequencing. Each  
35 sample from this cohort was sampled on day 2, the day of the enrollment of household contacts,  
36 one day after the presentation of the household index case at the hospital. To minimize the chance  
37 that the contact acquired a *V. cholerae* infection between index case enrollment and the first  
38 contact sample, the first household contact sample was always taken within a day, and generally  
39 <24h since the enrollment of the index case. To increase power in our study, we added household  
40 contacts samples from additional households, and also added samples from additional time points  
41 of Midani 2018 cohort contacts (see Table S1). This larger group of samples is referred to as the  
42 Expanded cohort, and is comprised of the 65 samples from the Midani 2018 cohort plus 33  
43 additional samples from 16 infected contacts (Figure 1). The 16 infected new contacts in the  
44 expanded cohort have multiple samples from different days that are prior to when infection  
45 occurred, and therefore have multiple “pre-infection” samples (Table S1). Enrollment and sample  
46 collection were identical for the two groups of samples collected. In total, 129 household contacts  
47 were initially enrolled. Eighteen contacts were excluded due to a positive rectal swab at the

48 enrollment evaluation, and 3 were excluded due to recent antibiotic use. Nine contacts reported  
49 ambiguous clinical outcomes based on clinical and history evaluation (i.e. report of diarrhea with  
50 no serologic or culture evidence of *V. cholerae* infection). Six additional enrollees were excluded  
51 due to report of diarrhea during the week prior to enrollment. Among the remaining contacts, 10  
52 resulted with DNA evidence of *V. cholerae* infection despite a rectal swab culture negative for *V.*  
53 *cholerae* at the time of enrollment. Lastly, 13 contacts were excluded due to failure to amplify  
54 DNA from rectal swabs or sequencing failure.

55

### 56 **DNA extraction and sequencing**

57 Fecal samples and rectal swabs were collected during home visit and immediately placed  
58 on ice and stored at -80°C until DNA extraction. For each sample, DNA was extracted as  
59 previously described <sup>4</sup>(1). Briefly, samples were processed using PowerSoil DNA extraction kits  
60 (Qiagen) after pre-heating to 65°C for 10 min and to 95°C for 10 min. A total of 98 samples were  
61 selected for library construction with NEBNext Ultra II DNA library prep kit and sequenced on  
62 the Illumina HiSeq 2500 (paired-end 125 bp) and the Illumina NovaSeq 6000 S4 (paired-end 150  
63 bp) at the Genome Québec sequencing platform (McGill University).

64

### 65 **Sequence preprocessing**

66 Sequencing fastq files were quality checked with FastQC  
67 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). For removal of human and  
68 technical contaminant DNA, metagenomic shotgun sequences were aligned to the PhiX genome  
69 and the human genome (hg19) with Bowtie2 (2). Host-decontaminated fastq files were then  
70 quality filtered using Trimmomatic, a computational tool for removing low quality data (3),  
71 discarding all reads shorter than 75 nucleotides and with quality Phred score less than 20.

72

### 73 **Taxonomic and functional profiling**

74           We performed taxonomic profiling of archaea, bacteria and microbial eukaryotes using  
75 MetaPhlAn2 (version 2.9) (4), with default parameters. Briefly, MetaPhlAn2 maps shotgun  
76 metagenomic sequencing reads against a precomputed database of clade-specific markers to  
77 produce a robust estimate of the taxonomic clades present in the microbiome sample and estimate  
78 their relative abundance. In addition to species-level relative abundance, we used MetaPhlAn2 to  
79 characterize the presence/absence of strain-specific markers for each sample. Species abundances  
80 are real numbers in the range [0,1] while markers assume binary values.

81           Species-resolved functional profiling was completed using HMP Unified Metabolic  
82 Analysis Network 2 (HUMAN2, version 2.8.1) (5), which maps reads to a sample-specific  
83 reference database from the pangenomes of the subset of species detected in the sample by  
84 MetaPhlAn2, quantifying gene presence and abundance on a per-species basis. A translated  
85 search is then performed against a UniRef-based protein sequence catalog for reads that fail to  
86 map to one of the detected species. The results are abundance profiles of gene families  
87 (UniRef90s) stratified by species contributing to those genes, which can be regrouped in higher  
88 level grouping categories (Pfam domains in this data set). Finally, gene families were further  
89 analyzed to reconstruct metabolic pathways abundance according to the MetaCyc pathway  
90 catalog (5).

91

### 92 **Statistical analyses**

93           Univariate analyses was performed to identify discriminatory biomarkers (species, protein  
94 families and pathways) using linear discriminant analysis (LDA) effect sizes (LEfSe) (6). LEfSe  
95 first uses the non-parametric factorial Kruskal-Wallis test to detect features with significant

96 differential abundance for each class of interest, then uses LDA to estimate the effect size of each  
97 differential abundant feature (6). In our case, the significance threshold was set at  $p=0.05$  and an  
98 LDA effect size of 2 was used for discriminative features.

99

## 100 **Random-forest based machine learning approach**

101 MetAML (7) was applied to four types of quantitative profiles: taxonomic species-level  
102 relative abundances, strain-specific markers presence/absence patterns, MetaCyc pathways and  
103 gene-families grouped by Pfam domains relative abundances. In each case, we used Random  
104 Forest (RF) as classifier and set the following parameters as previously described by (8): the  
105 number of estimator trees was equal to 1000 and the quality of split at each tree node was  
106 evaluated using Shannon entropy. The minimum number of samples per leaf and the number of  
107 features per tree were respectively set as 1 and 30%, except for the strain-specific markers  
108 presence/absence profile, where the number of features equal the square root of the total number  
109 of features, due to the higher number of features. One of the advantages of RF models is that they  
110 can intrinsically integrate binary and multi-class classification problems, and implicitly provide a  
111 list of the most informative features to the predictive model. Another advantage of RFs is the  
112 built-in feature selection step during the model generation phase, which allows a selection of a  
113 reduced subset of the most important features for discriminating between classes. Adding a  
114 feature selection step to a model is a useful way to remove irrelevant features, especially in  
115 datasets with high dimensionality. Feature selection can also reduce overall training time and the  
116 risk of overfitting (7).

117 The feature relative importance values were used to perform an embedded feature  
118 selection strategy, implemented as described in (7), in addition to the RF model, for two  
119 (Uninfected vs Infected contacts) and three classes (Asymptomatic infected vs Symptomatic

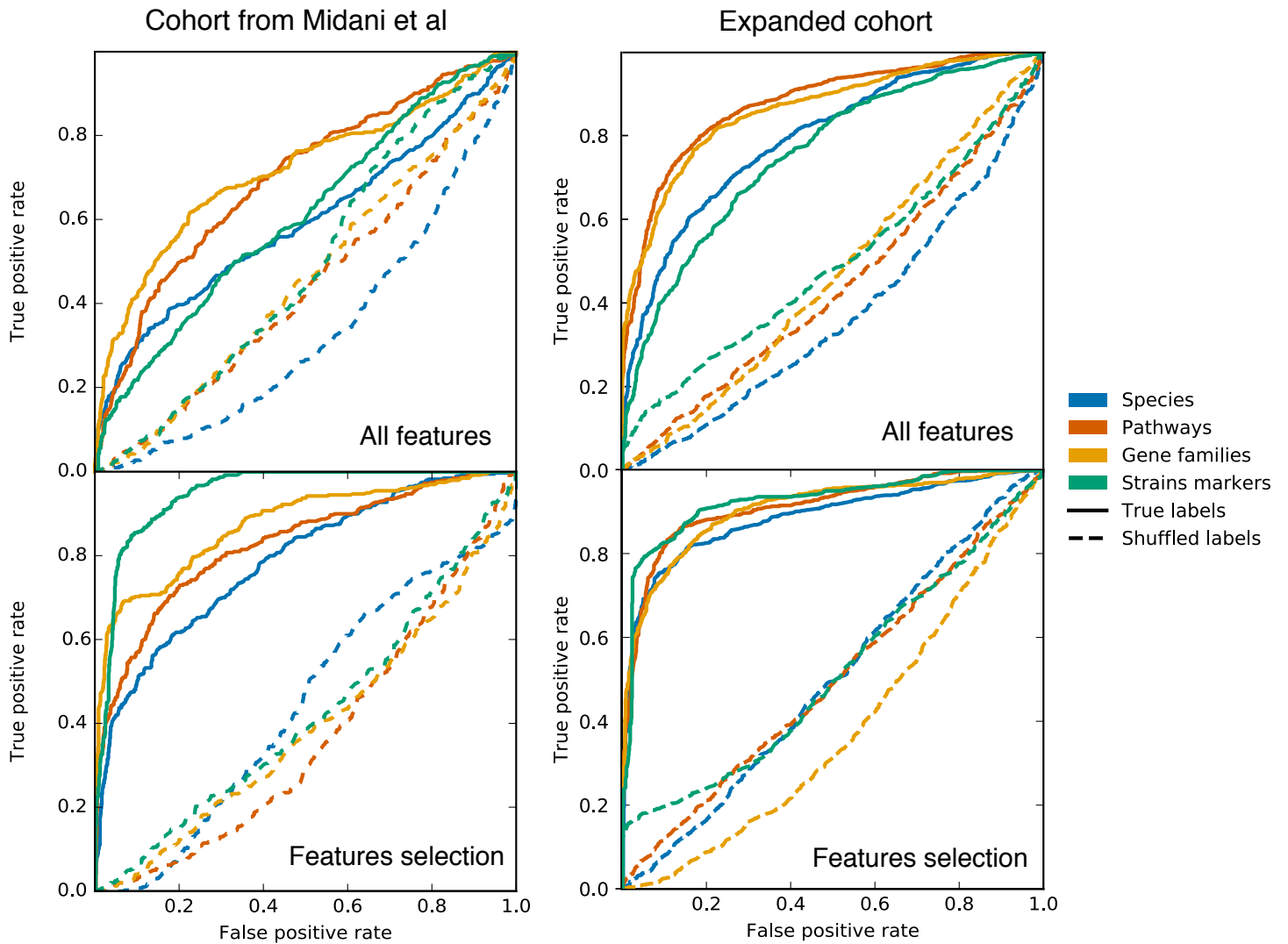
120 infected vs Uninfected contacts). Feature selection benefits include removal of irrelevant features,  
121 especially in datasets with high dimensionality, to decrease the overall training time and reduce  
122 the risk of model overfitting (9). Each cohort (the Midani\_2018 and the expanded cohort) was  
123 analyzed independently for the RF model, using a stratified cross-validation approach. In the  
124 feature selection model, we chose the minimum number of features that maximized the accuracy  
125 in order to generate a final model on this limited number of features (Figure S2). The same set of  
126 selected features determined by the Midani\_2018 dataset was used for the two cohorts in this  
127 approach. In all cases, the sensitivity and accuracy of the models generated were tested by  
128 performing stratified 3-fold cross validations, repeated and averaged on 20 independent  
129 runs. Finally, the results obtained for the original classification problem were compared with  
130 those obtained by a random classifier (denoted in the paper as Shuffled). For this purpose, we  
131 shuffled randomly the labels of all the samples, and used these same settings. Difference of  
132 performance between classifiers and statistical significance were calculated as described in  
133 Pasolli *et al* (7).

134         The following performance metrics were reported: 1) the overall accuracy (i.e., the  
135 proportion of outcomes correctly predicted), 2) the precision (i.e., the number of correct positive  
136 samples divided by the number of samples predicted as positive), 3) the recall (i.e. the true  
137 positive rate, or power), and 4) the F1 score, which is the harmonic mean of precision and recall.  
138 For binary classification problems, class posterior probabilities were used to plot the Receiver  
139 Operating Characteristic (ROC) curve, which represents the true positive rate (i.e., the recall)  
140 against the false positive rate (i.e., the number of wrong positive samples divided by the total  
141 number of non-positive samples). From the ROC curve, the program computes the area under the  
142 curve (AUC), which can be interpreted as the probability that a randomly selected positive

143 sample will have a higher classification result than a randomly selected negative one. The AUC  
144 ranges in  $[0, 1]$ , where 0.5 corresponds to random prediction.

145

146

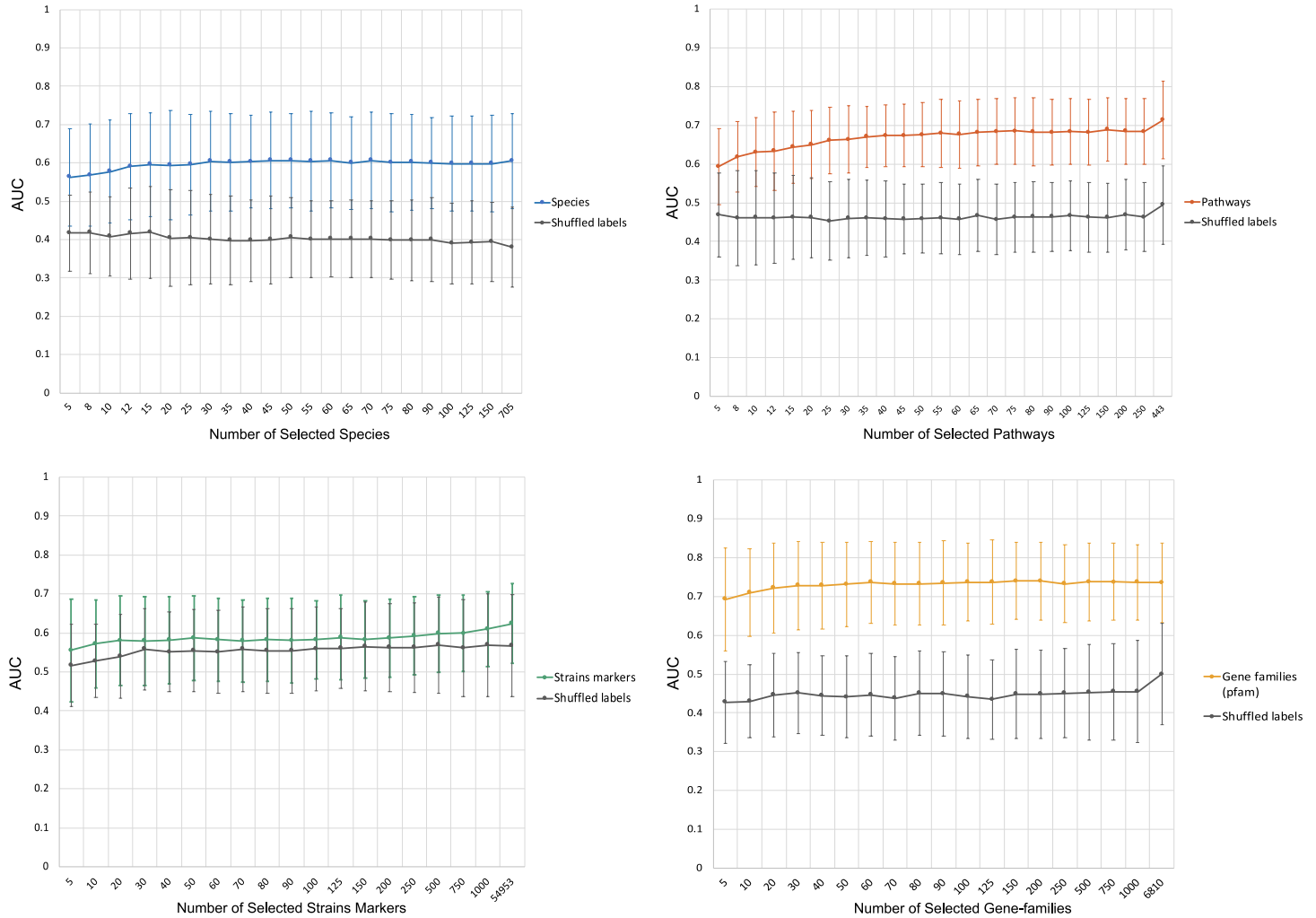


148  
 149 **Figure S1. Classification result from species abundance, strain-level markers**  
 150 **presence/absence, pathways and gene family abundances.** Four different classifications of  
 151 microbiome information were used to predict disease susceptibility (Uninfected vs Infected)  
 152 using a random forest algorithm and stratified 3-fold cross validation. Plots show Average ROC  
 153 curves by using the four type of features for two different datasets (the Midani et al data and the



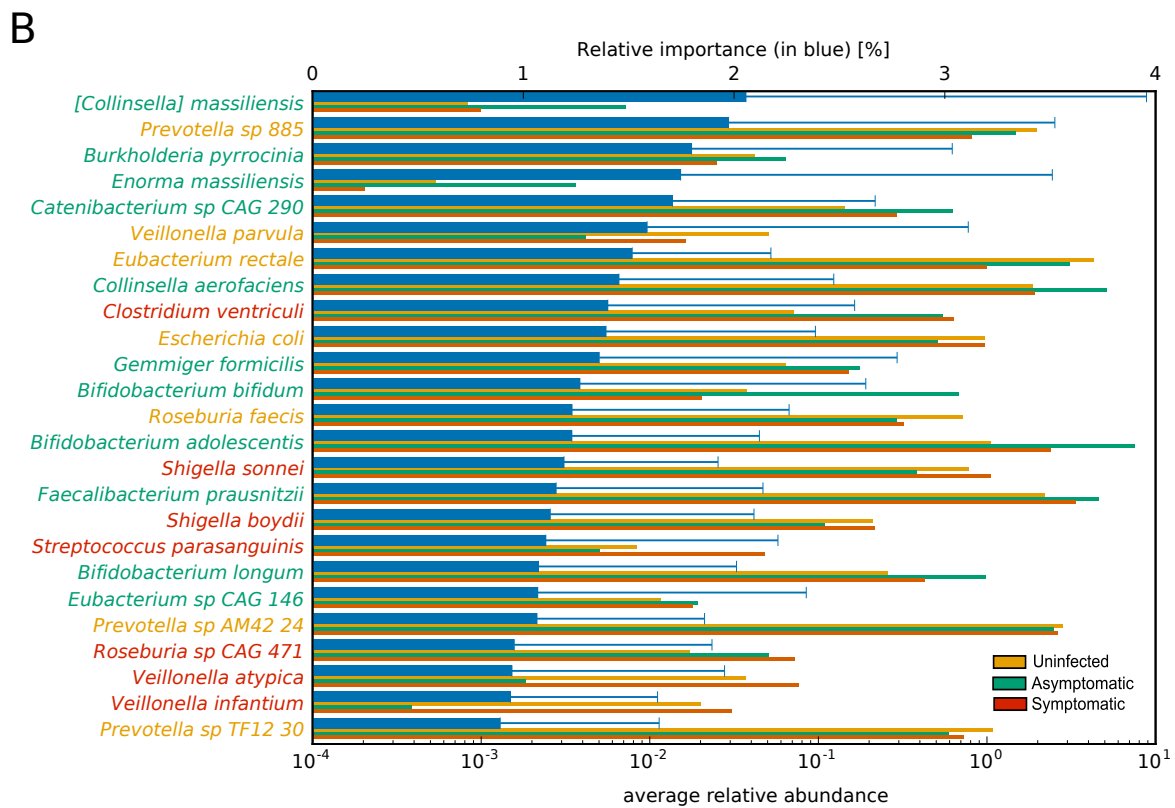
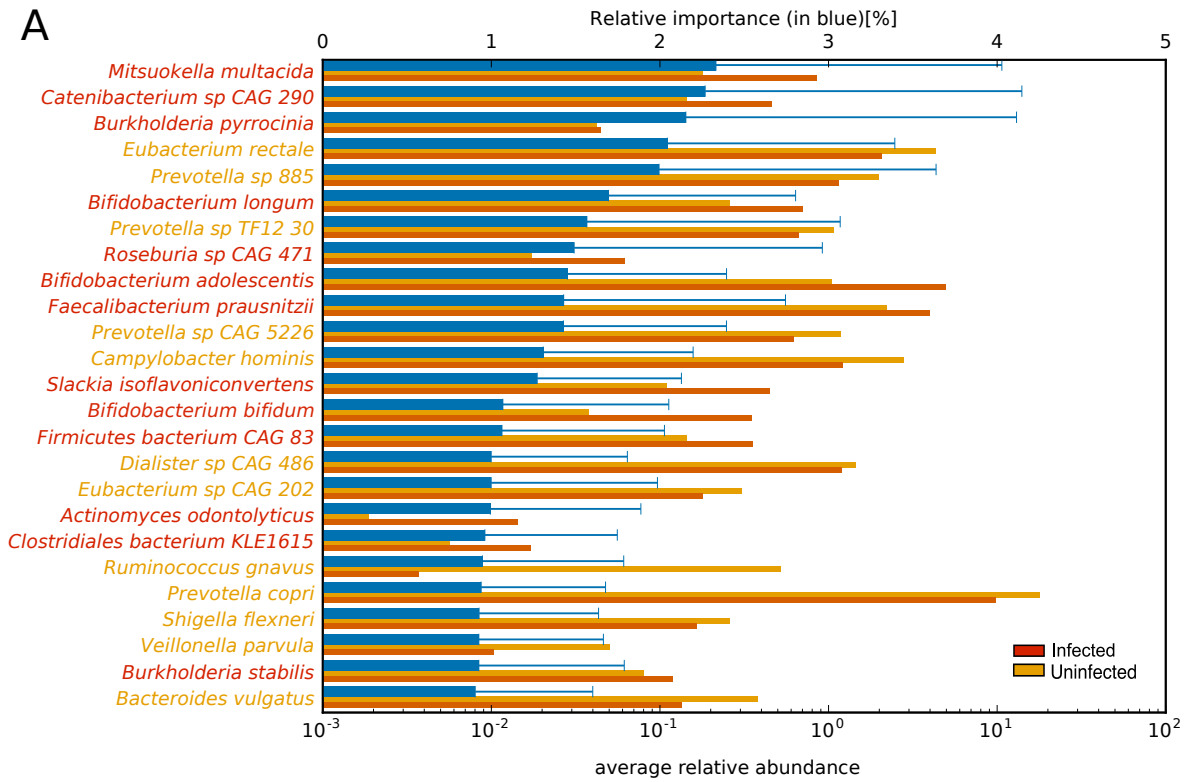
154 expanded dataset). Random Forest was applied on all features (top row) and on a set of selected  
155 features (bottom row).

156

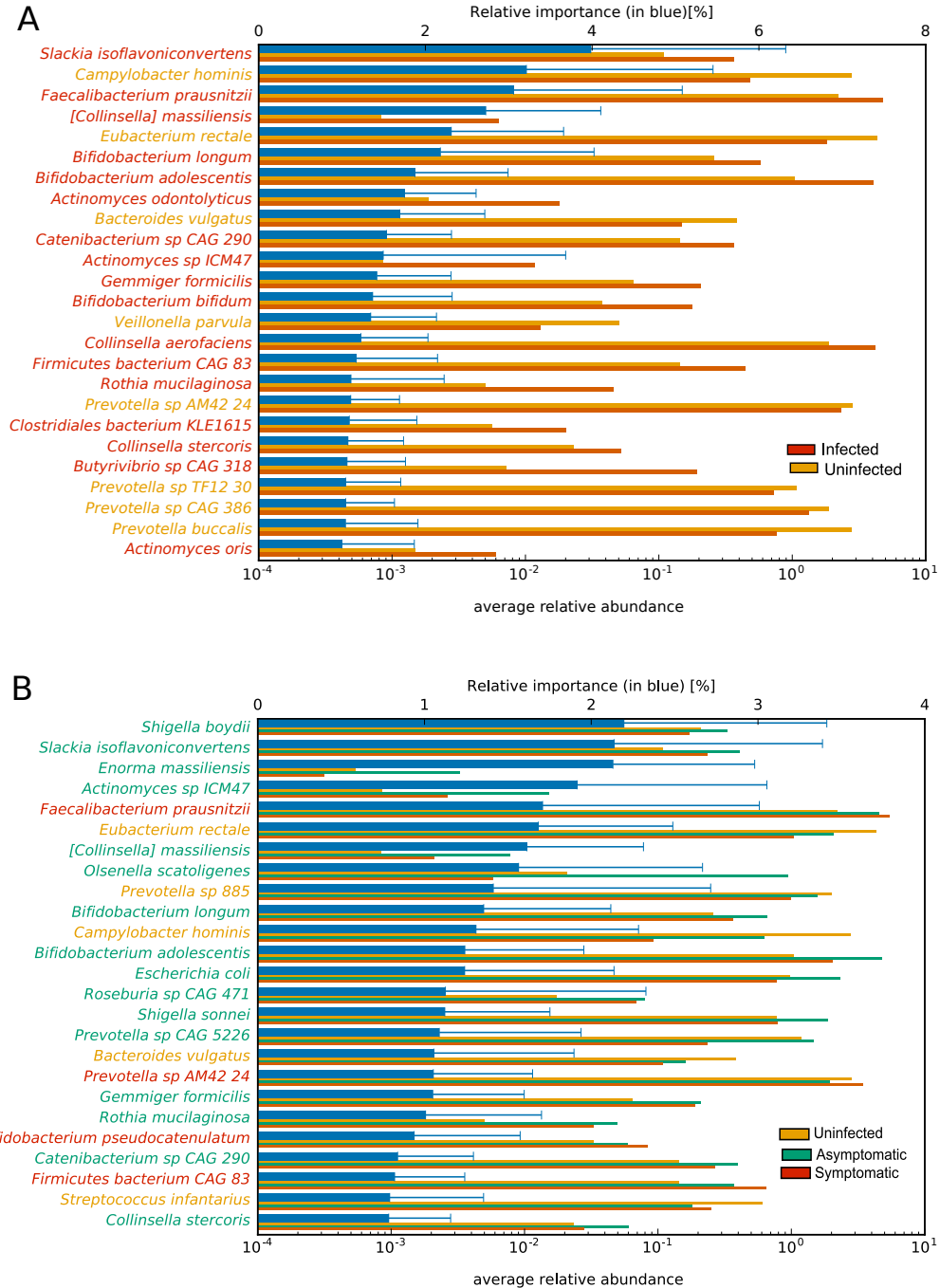


157 **Fig. S2. Identification of a minimal number of features for each type of biomarker.**

158 Prediction performances (AUC values) at increasing number of microbial species, pathways,  
159 strains markers and genes families, obtained by retraining the random forest model on the top-  
160 ranking features identified with a first random forest model training with a 3-fold cross-validation  
161 approach.



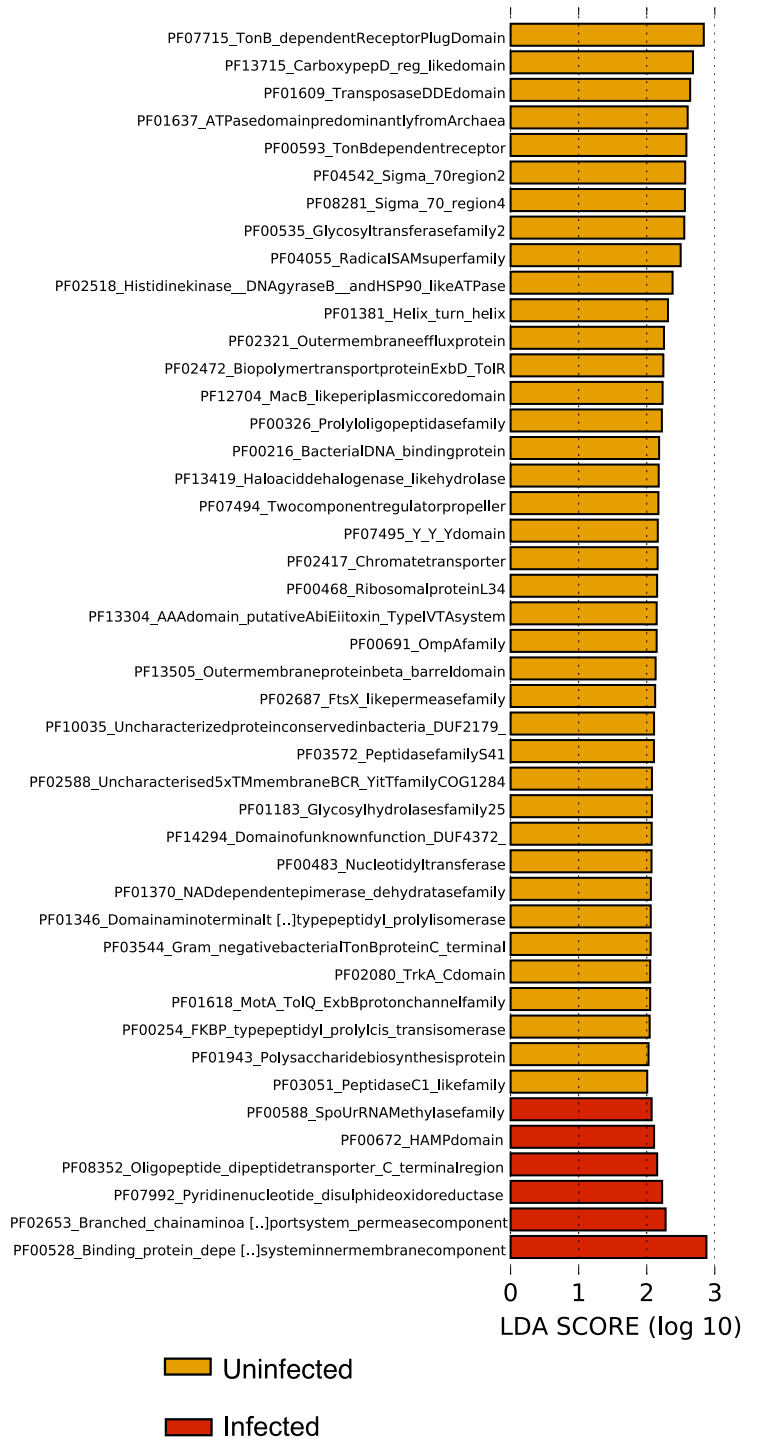
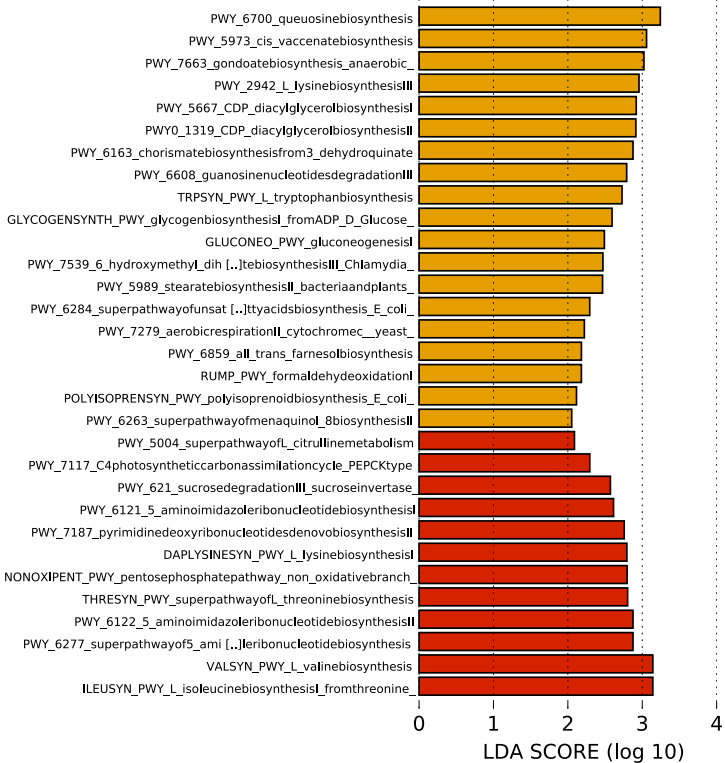
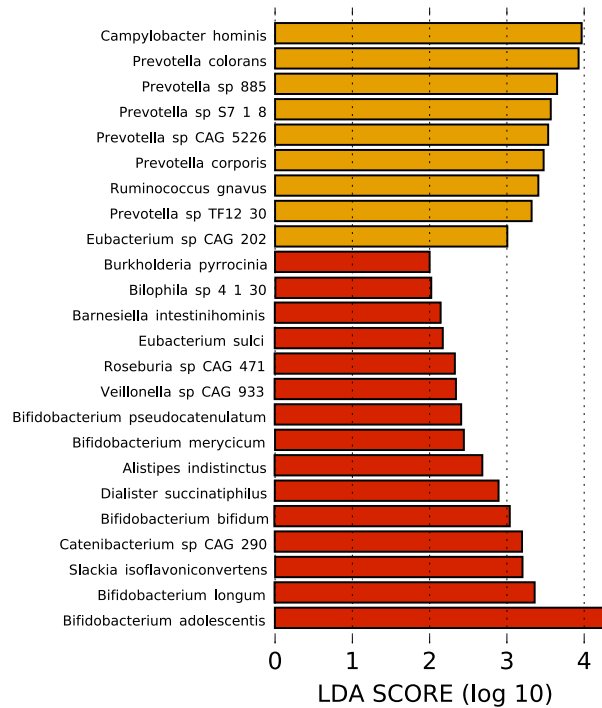
163 **Fig. S3. Most important discriminating species of the gut microbiome at the time of**  
164 **exposure to *V. cholerae* identified in the Midani 2018 dataset, classified by clinical outcome.**  
165 (A) Species associated with contacts that became infected (or remained uninfected) during  
166 follow-up. (B) Species associated with contacts that became  
167 uninfected/asymptomatic/symptomatic during follow-up. For each species reported on the  
168 vertical axis, the top bar (blue) corresponds to the species relative importance (with standard  
169 deviation) and the other bars refer to the average relative abundance. The top 25 most important  
170 features are shown here; See Table S6 for the full list. Feature relative importance was computed  
171 using the Mean Decrease in Impurity strategy, as described in the Methods.  
172



173 **Fig. S4. Most important discriminating species identified with the random forest algorithm**  
 174 **on the expanded dataset for (A) contacts that became infected or remained uninfected and**  
 175 **(B) contacts that remained uninfected, or became infected and were asymptomatic vs**  
 176 **symptomatic.** For each species reported on the vertical axis, the top bar (in blue) corresponds to  
 177 the feature relative importance (with standard deviation) and the other bars refer to the average

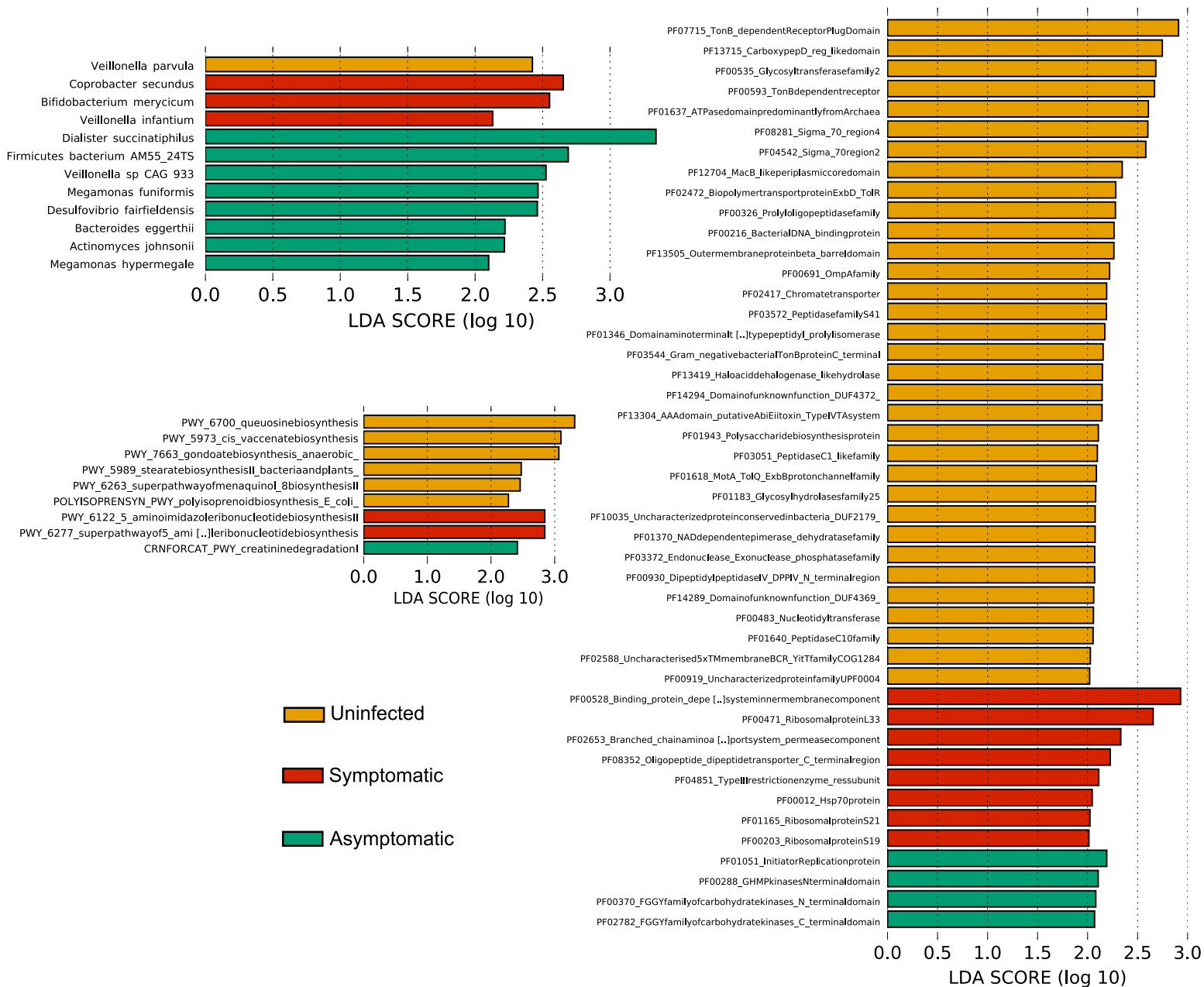
178 relative abundance. Feature relative importance (blue) is computed using the mean decrease

179 impurity strategy.



180

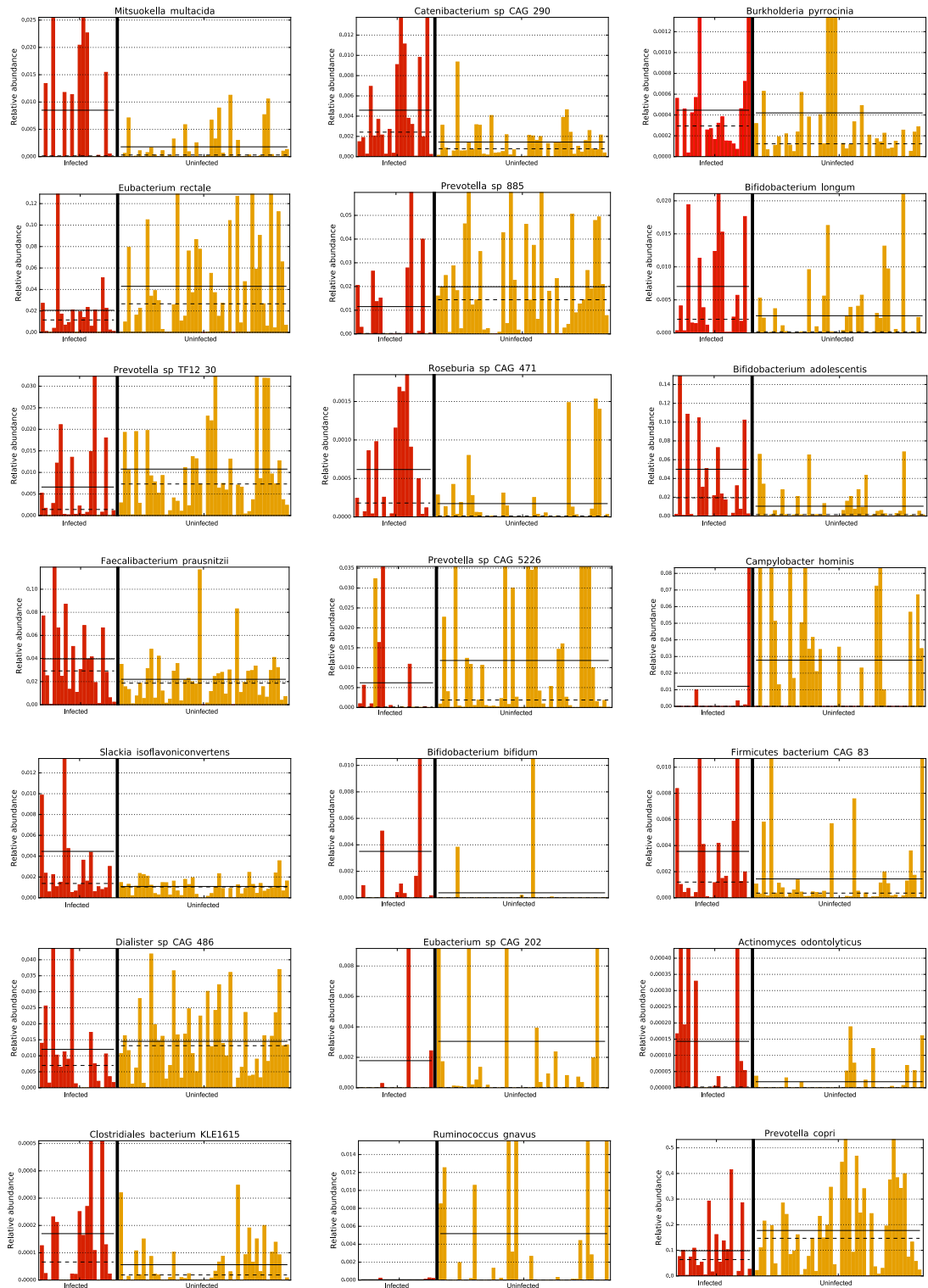
181 **Fig. S5. Linear discriminant analysis (LDA) scores computed for differentially abundant**  
182 **species, pathways and genes families in the fecal microbiomes of samples from the**  
183 **Midani\_2018 cohort for 2 categories.** The cohort is divided in two categories: controls who  
184 remained uninfected (yellow) and controls who became infected (red). Length of the bar indicates  
185 effect size associated with a feature.  $p = 0.05$  for the Kruskal-Wallis H statistic; features with  
186 LDA score  $> 2$  are shown.



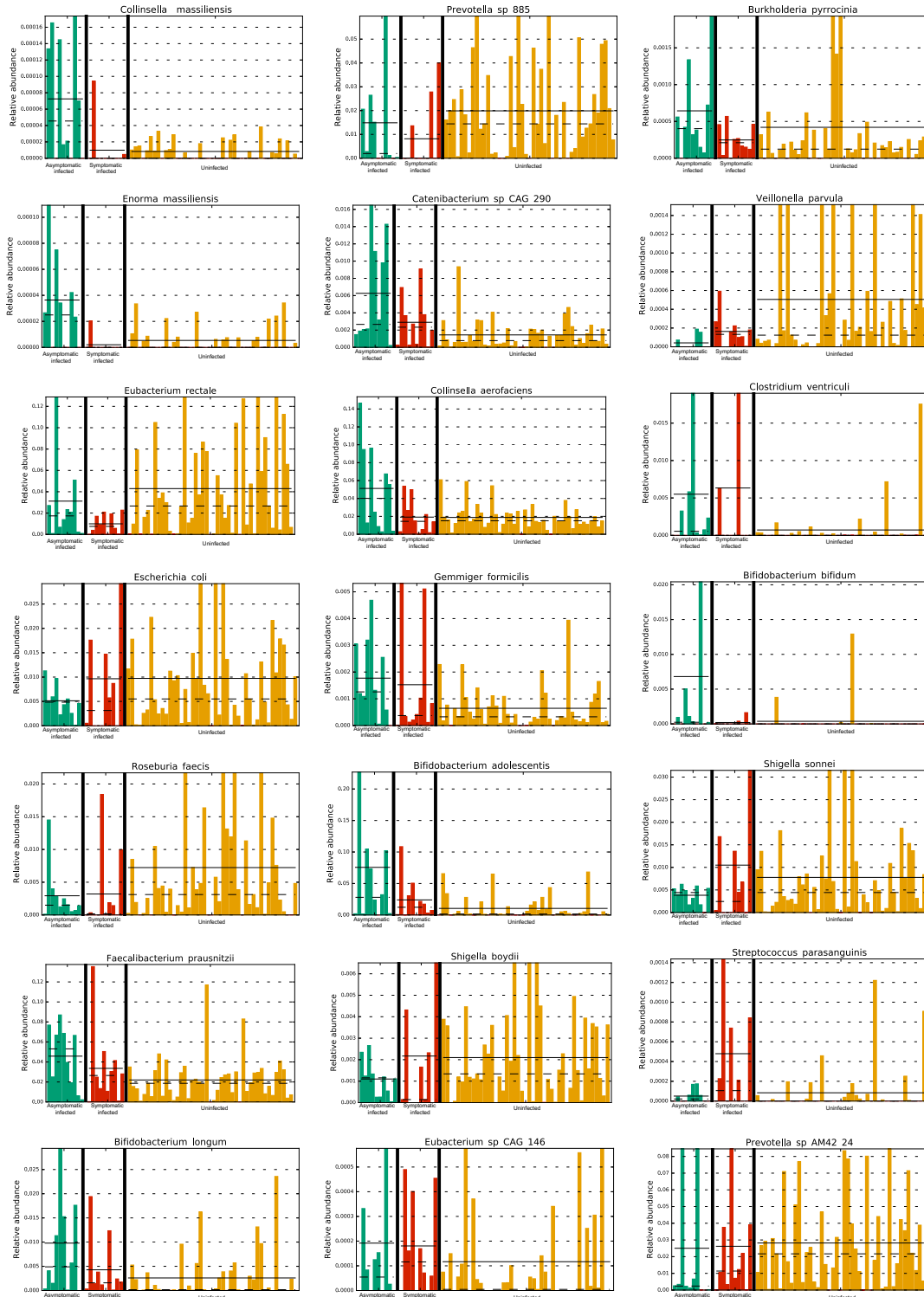
189 **Fig. S6. Linear discriminant analysis (LDA) scores computed for differentially abundant**  
 190 **species, pathways and genes families in the fecal microbiomes of samples from the**  
 191 **Midani\_2018 cohort for 3 categories.** The cohort is divides in three categories: controls who  
 192 remained uninfected (yellow), controls who became infected and symptomatic (red), and controls

193 who became infected but stayed asymptomatic (green). Length indicates effect size associated  
194 with a feature.  $p = 0.05$  for the Kruskal-Wallis H statistic; features with LDA score  $> 2$  are  
195 shown.

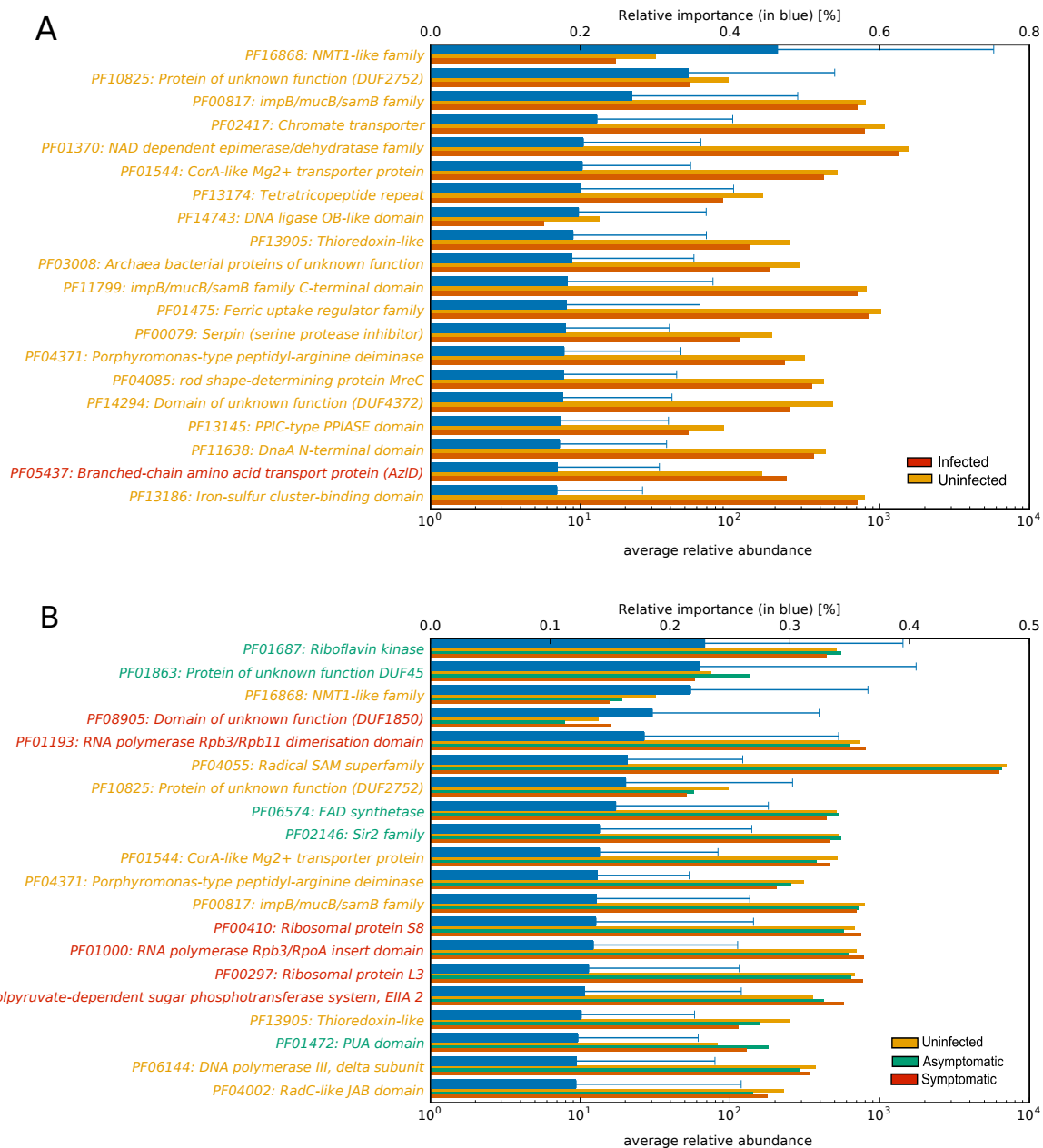




196  
 197 **Fig. S7. Relative abundance of the top 21 most important discriminating species of the gut**  
 198 **microbiome of contacts at the time of exposure to *V. cholerae* identified in the Midani 2018**  
 199 **dataset for two classes (Uninfected vs Infected). The straight line indicates the group means,**  
 200 **and the dotted line indicates the group medians.**

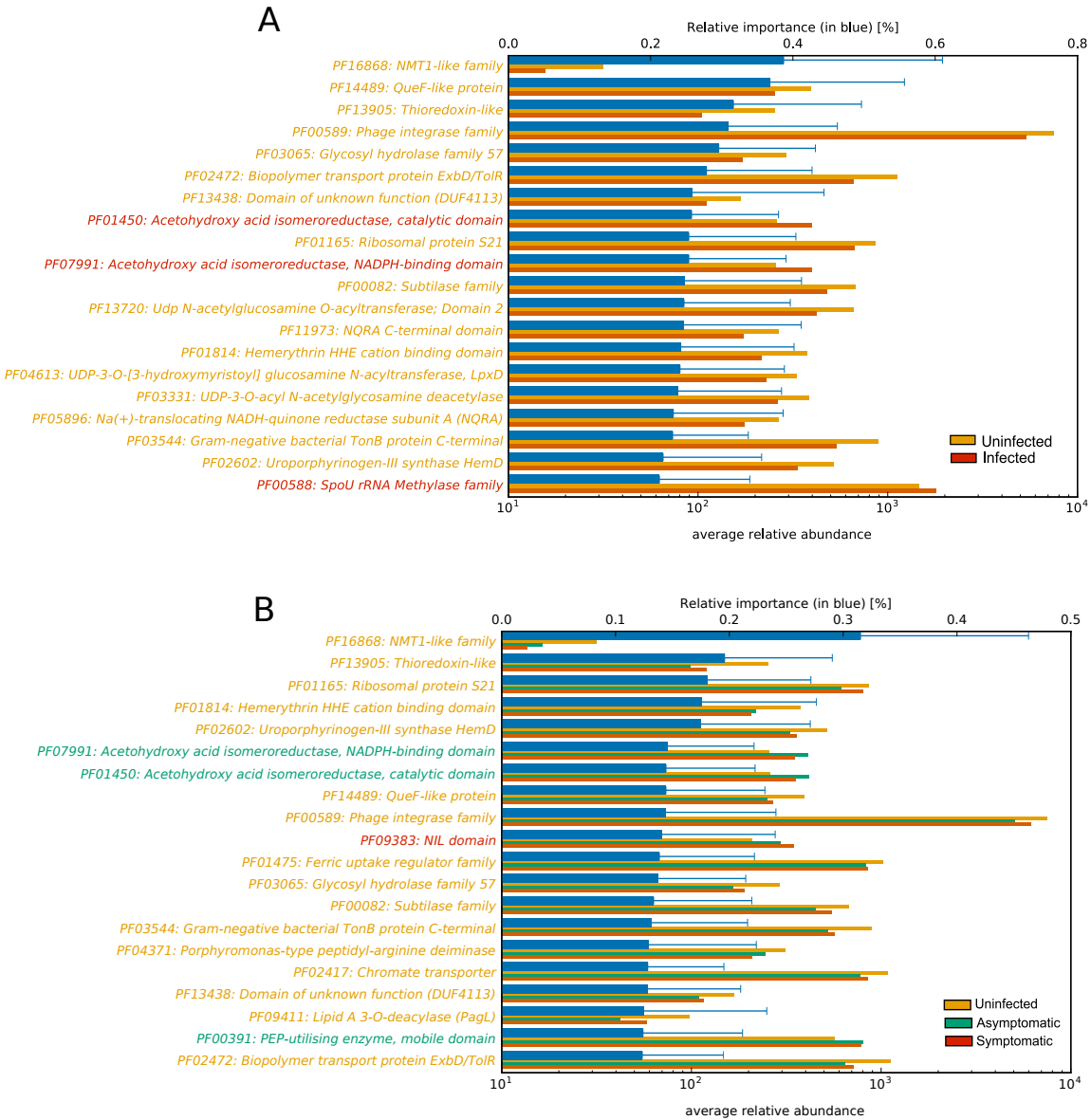


201  
 202 **Fig. S8. Relative abundance of the top 21 most important discriminating species in the gut**  
 203 **microbiome of contacts at the time of exposure to *V. cholerae* identified in the Midani 2018**  
 204 **dataset for three classes (Uninfected vs Asymptomatic Infected and Symptomatic Infected).**  
 205 The straight line indicates the group means, and the dotted line indicates the group medians



206  
 207 **Figure S9. Most important discriminating gene families of the gut microbiome at the time of**  
 208 **exposure to *V. cholerae* identified in the Midani 2018 dataset, classified by clinical outcome.**  
 209 (A) Genes associated with contacts that became uninfected/infected during follow-up. (B) Genes  
 210 associated with contacts that became uninfected/asymptomatic/symptomatic during follow-up.  
 211 For each gene-family reported on the vertical axis, the top bar (in blue) corresponds to the feature

212 relative importance (with standard deviation) and the other bars refer to the average relative  
213 abundance (copies per million). The top 25 most important features are shown here; See Table S8  
214 for the full list. Feature relative importance was computed using the mean decrease in impurity  
215 strategy, as described in the Methods.



216

217

218 **Fig. S10. Most important discriminating gene-families (grouped by Pfam domain) identified**

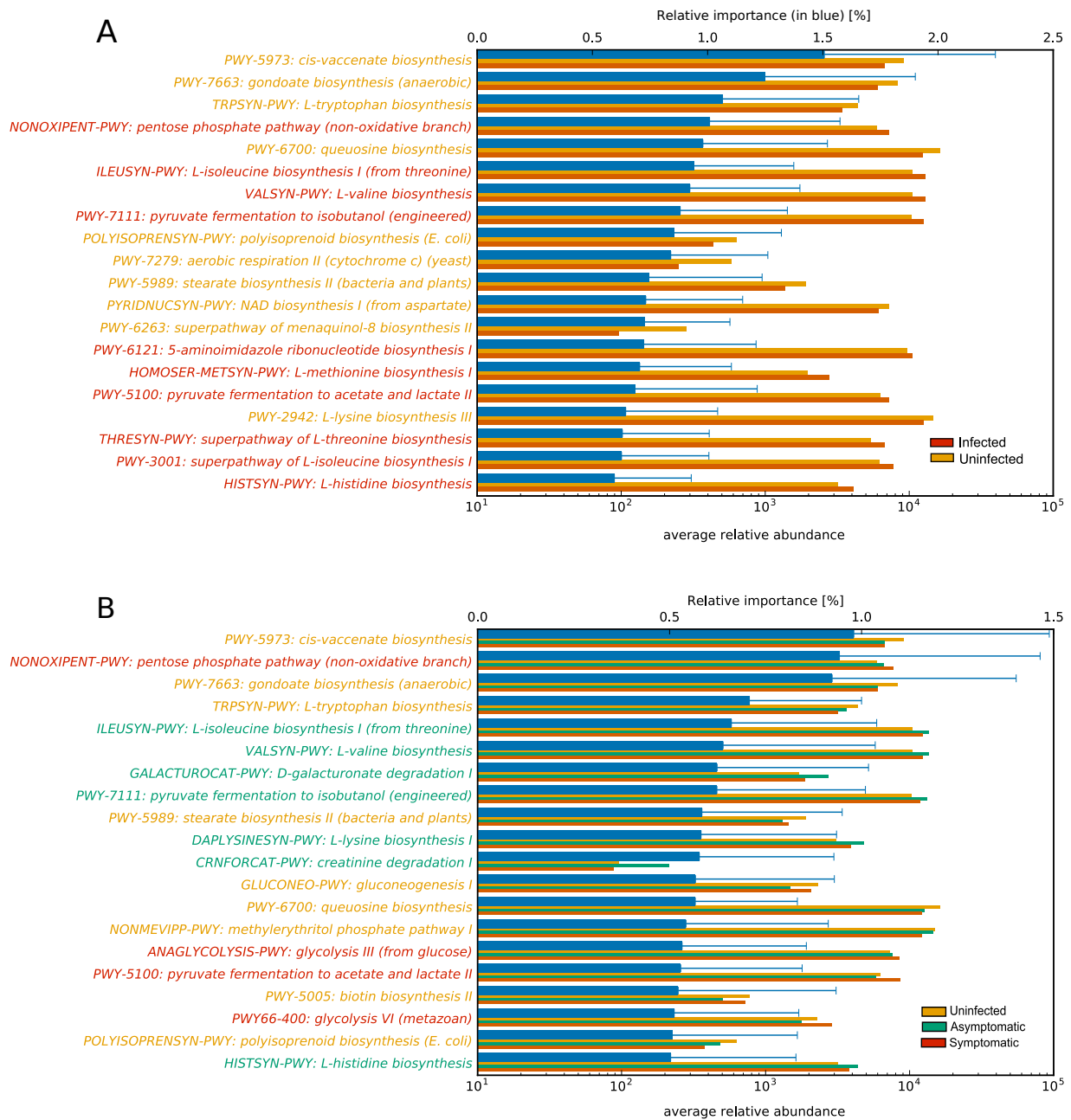
219 **with the random forest algorithm on the Expanded dataset for (A) contacts that became**

220 **infected or remained uninfected and (B) contacts that remained uninfected, or became**

221 **infected and were asymptomatic vs symptomatic.** For each gene-families reported on the

222 vertical axis, the top bar (in blue) corresponds to the feature relative importance (with standard

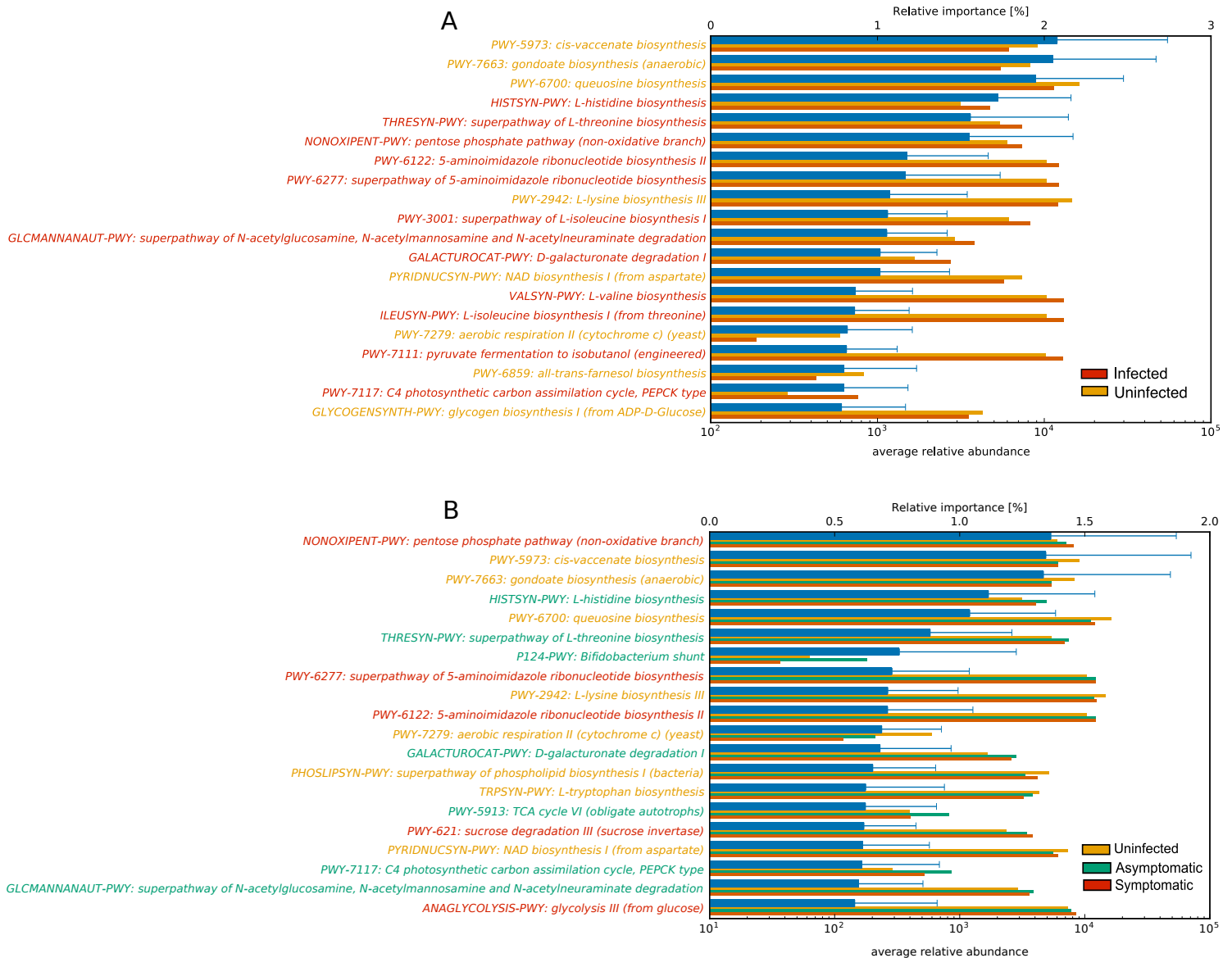
223 deviation) and the other bars refer to the average relative abundance (copies per million). Feature  
224 relative importance (blue) is computed using the mean decrease impurity strategy.



226

227 **Fig. S11. Most important discriminating pathways identified with the random forest**  
 228 **algorithm on the Midani\_2018 dataset for (A) contacts that became infected or remained**  
 229 **uninfected and (B) contacts that remained uninfected or became infected and were**  
 230 **asymptomatic vs symptomatic.** For each pathway reported on the vertical axis, the top bar

231 (blue) corresponds to the feature relative importance (with standard deviation) and other bars  
 232 refer to the average relative abundance (copies per million). Feature relative importance (blue) is  
 233 computed using the mean decrease impurity strategy.



234

235 **Fig. S12. Most important discriminating pathways identified with the random forest**

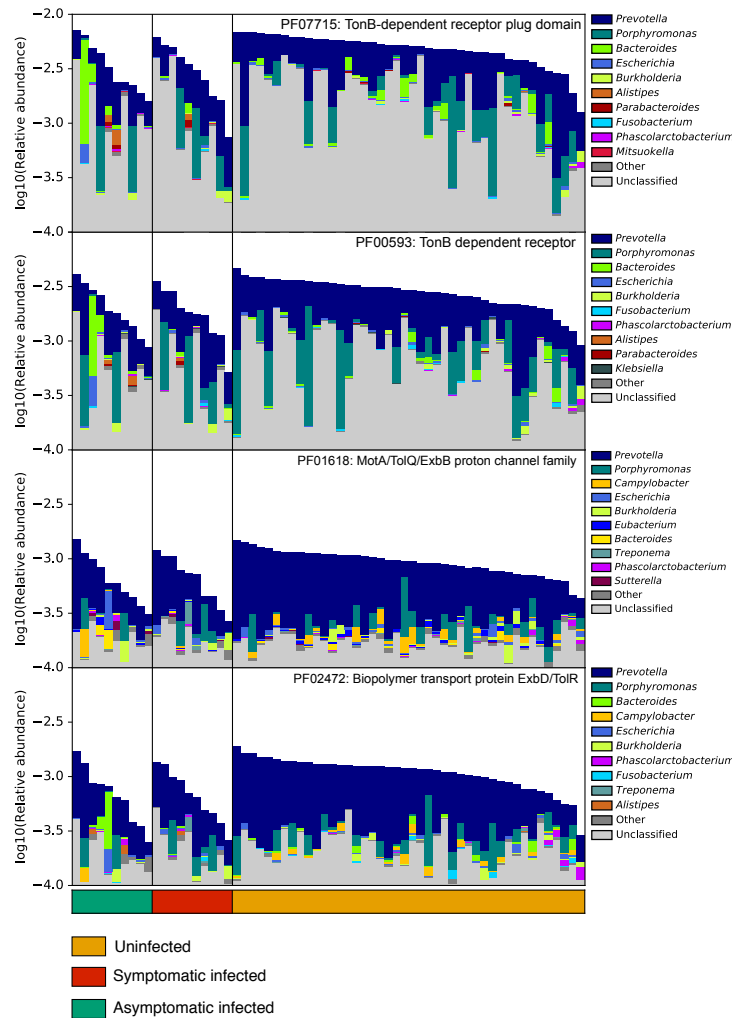
236 **algorithm on the expanded dataset for (A) contacts that became infected or remained**

237 **uninfected and (B) contacts that remained uninfected, or became infected and were**

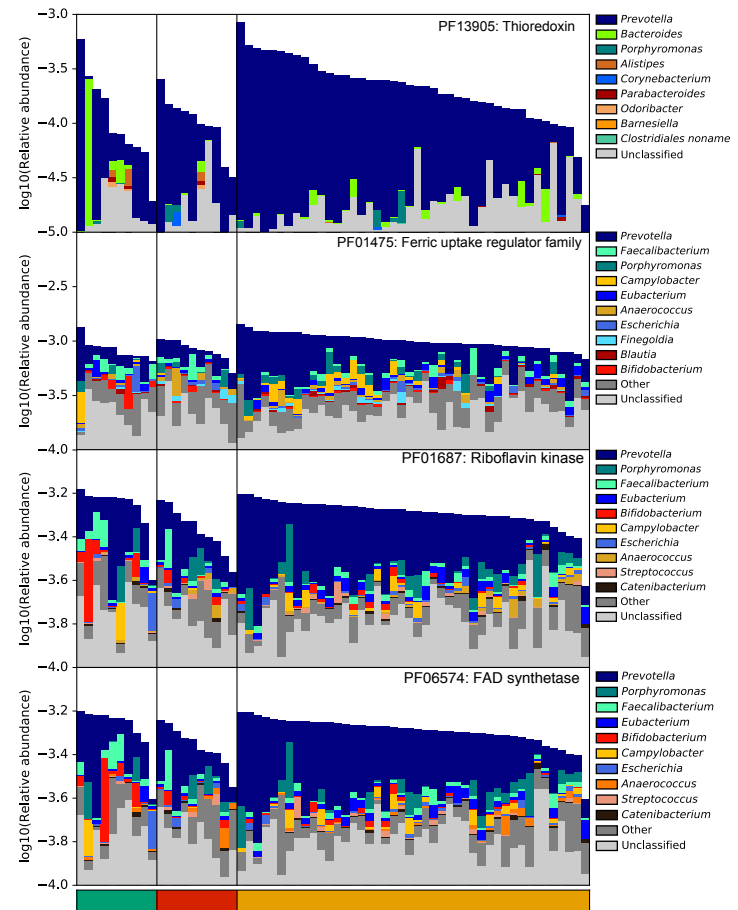


238 **asymptomatic vs symptomatic.** For each pathway reported on the vertical axis, the top bar (in  
239 blue) corresponds to the feature relative importance (with standard deviation) and the other bars  
240 refer to the average relative abundance (copies per million). Feature relative importance (blue) is  
241 computed using the mean decrease impurity strategy.

Microbial gene families enriched in contacts who remained **uninfected**

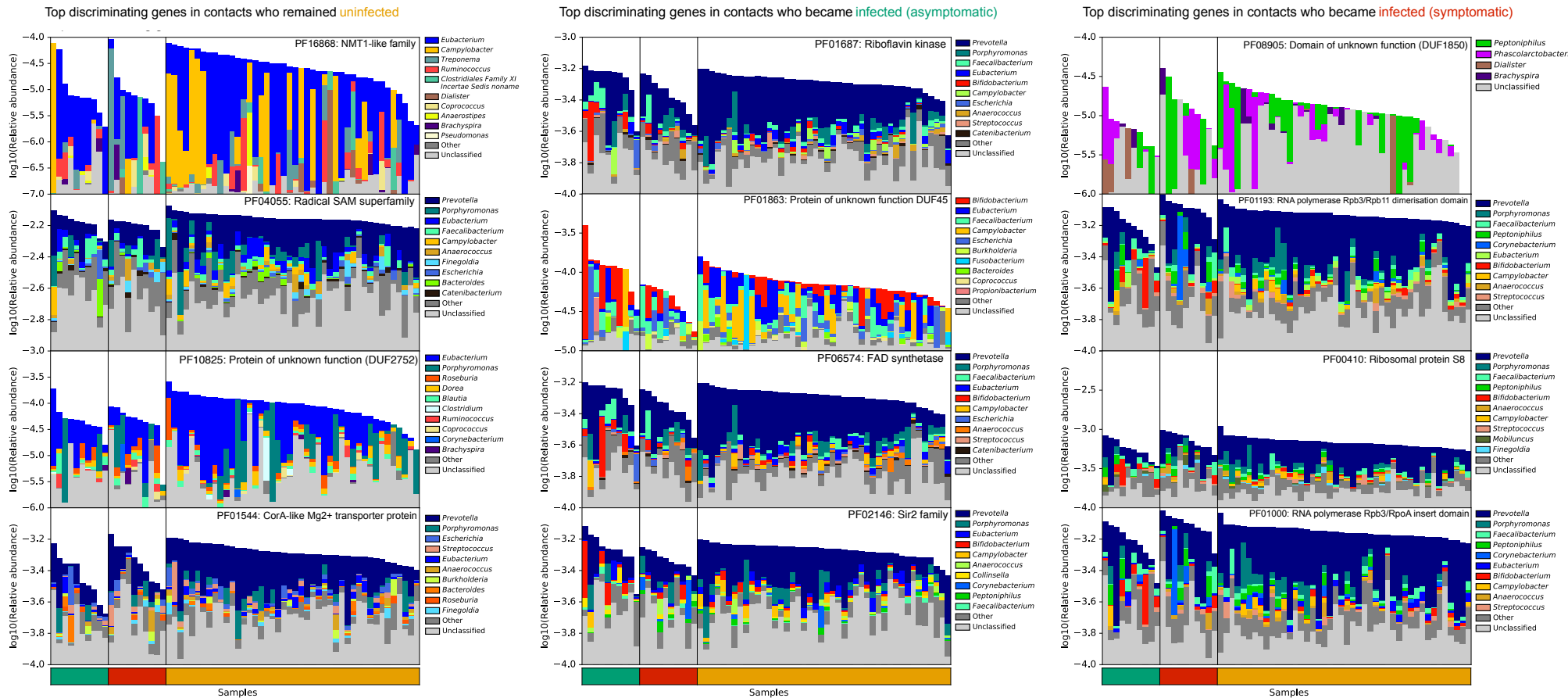


Top discriminating gene families in contacts who remained **uninfected** and in contacts who became **infected (asymptomatic)**



242 **Fig. S13. Enriched and top predictive gene families (Pfam domain grouping) involved in iron metabolism, annotated by their**  
 243 **taxonomic contributors. Gene families involved with iron metabolism that were differentially abundant in contacts who remained**

244 uninfected were identified with LEfSe are represented on the left (shown in Supplementary Fig 2). On the right are the top predictive  
 245 gene families involved with iron metabolism identified with MetAML. Total bar height reflects log<sub>10</sub>-scaled community abundance.  
 246 Genera contributions are linearly scaled within total.  
 247



248 **Fig. S14. Top predictive gene families (Pfam domain grouping) for each class (uninfected, asymptomatic and symptomatic**  
249 **infected contacts), annotated by their taxonomic contributors.** Total bar height reflects  $\log_{10}$ -scaled community abundance. Genera  
250 contributions are linearly scaled within total.  
251

252 **Supplementary References**

- 253
- 254 1. Midani FS, Weil AA, Chowdhury F, Begum YA, Khan AI, Debela MD, et al. Human Gut  
255 Microbiota Predicts Susceptibility to *Vibrio cholerae* Infection. *J Infect Dis.* **2018**.
- 256 2. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature Methods.*  
257 **2012**;9: 357–359.
- 258 3. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence  
259 data. *Bioinformatics.* **2014**;30: 2114–2120.
- 260 4. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, et al. MetaPhlAn2  
261 for enhanced metagenomic taxonomic profiling. *Nature Methods.* **2015**;12: 902–903.
- 262 5. Franzosa EA, McIver LJ, Rahnnavard G, Thompson LR, Schirmer M, Weingart G, et al.  
263 Species-level functional profiling of metagenomes and metatranscriptomes. *Nature*  
264 *Methods.* **2018**;15: 962–968.
- 265 6. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, et al. Metagenomic  
266 biomarker discovery and explanation. *Genome Biol.* **2011**;12: 1–18.
- 267 7. Pasolli E, Truong DT, Malik F, Waldron L, Segata N. Machine Learning Meta-analysis of  
268 Large Metagenomic Datasets: Tools and Biological Insights. *PLoS Comput Biol.* **2016**;12:  
269 e1004977.
- 270 8. Thomas AM, Manghi P, Asnicar F, Pasolli E, Armanini F, Zolfo M, et al. Metagenomic  
271 analysis of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures  
272 and a link with choline degradation. *Nat Med.* **2019**;25: 667–678.
- 273 9. Zhou Y-H, Gallins P. A Review and Tutorial of Machine Learning Methods for Microbiome  
274 Host Trait Prediction. *Front Genet.* **2019**;10.