

1 **Title: Integrating Gut Microbiota and Host Immune Markers for Highly**
2 **Accurate Diagnosis of *Clostridioides difficile* Infection**

3

4 **Authors:** Shanlin Ke^{1,2}, Nira R. Pollock^{3,4}, Xu-Wen Wang¹, Xinhua Chen⁵, Kaitlyn Daugherty⁵,
5 Qianyun Lin⁵, Hua Xu⁵, Kevin W. Garey⁶, Anne J. Gonzales-Luna⁶, Ciarán P. Kelly^{5*}, Yang-Yu
6 Liu^{1,7*}

7 **Affiliations:**

8 ¹Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical
9 School, Boston, Massachusetts 02115, USA.

10 ²State Key Laboratory of Pig Genetic Improvement and Production Technology, Jiangxi
11 Agricultural University 330045, China.

12 ³Division of Infectious Diseases, Department of Medicine, Beth Israel Deaconess Medical
13 Center, Boston, Massachusetts 02115, USA.

14 ⁴Department of Laboratory Medicine, Boston Children's Hospital, Boston, Massachusetts 02115,
15 USA.

16 ⁵Division of Gastroenterology, Department of Medicine, Beth Israel Deaconess Medical Center,
17 Boston, Massachusetts 02115, USA.

18 ⁶Department of Pharmacy Practice and Translation Research, University of Houston College of
19 Pharmacy, Houston, Texas 77204, USA.

20 ⁷Center for Cancer Systems Biology, Dana-Farber Cancer Institute, Boston, Massachusetts
21 02115, USA.

22 *To whom correspondence should be addressed: Y.-Y.L. (yyl@channing.harvard.edu) or C.P.K.
23 (ckelly2@bidmc.harvard.edu).

24

25 **One Sentence Summary:** Incorporating both gut microbiome and host immune marker data into
26 classification models can better distinguish CDI from other groups than can either type of data
27 alone.

28

29

30

31

32

33

34

35

36

37 **Abstract:** Exposure to *Clostridioides difficile* can result in asymptomatic carriage or infection
38 with symptoms ranging from mild diarrhea to fulminant pseudomembranous colitis. A reliable
39 diagnostic approach for *C. difficile* infection (CDI) remains controversial. Accurate diagnosis is
40 paramount not only for patient management but also for epidemiology and disease research.
41 Here, we characterized gut microbial compositions and a broad panel of innate and adaptive
42 immunological markers in 243 well-characterized human subjects, who were divided into four
43 phenotype groups: CDI, Asymptomatic Carriage, Non-CDI Diarrhea, and Control. We found that
44 CDI is associated with alteration of many different aspects of the gut microbiota, including
45 overall microbial diversity and microbial association networks. We demonstrated that
46 incorporating both gut microbiome and host immune marker data into classification models can
47 better distinguish CDI from other groups than can either type of data alone. Our classification
48 models display robust diagnostic performance to differentiate CDI from Asymptomatic carriage
49 (AUC~0.916), Non-CDI Diarrhea (AUC~0.917), or Non-CDI that combines all other three
50 groups (AUC~0.929). Finally, we performed symbolic classification using selected features to
51 derive simple mathematic formulas for highly accurate CDI diagnosis. Overall, this study
52 provides evidence supporting important roles of gut microbiota and host immune markers in CDI
53 diagnosis, which may also inform the design of future therapeutic strategies.

54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72

73 INTRODUCTION

74 *Clostridioides difficile* infection (CDI) is the most common cause of healthcare-associated
75 infection and an important cause of morbidity and mortality among hospitalized patients¹⁻³.
76 Current treatment strategies for CDI, including vancomycin, metronidazole and fidaxomicin,
77 have inconsistent cure rates and treatment failure or CDI recurrence may occur in approximately
78 one third of cases^{4,5}. Antibiotic exposure is considered the most important factor predisposing
79 patients to CDI^{6,7}. In fact, treatments with antibiotics have a tremendous impact on the
80 composition and functionality of the gut microbiota, and accordingly are associated with reduced
81 colonization resistance against pathogens such as *C. difficile*⁸⁻¹⁰. A distinct microbial community
82 structure has been reported to be associated with CDI in human cohorts and animal models^{11,12}.
83 Characterization of the microbial features in individuals with different *C. difficile*
84 infection/colonization status is an essential step in understanding the role of the gut microbiome
85 in the development of CDI.

86 The pathophysiology of CDI is mainly associated with the production of two exotoxins,
87 toxin A (TcdA) and toxin B (TcdB)¹³. TcdA and TcdB act on intestinal epithelial cells, inducing
88 pro-inflammatory cytokines, loss of tight junctions, cell detachment and an impaired mucosal
89 barrier¹⁴⁻¹⁶. The innate and adaptive immune responses to CDI play crucial roles in disease onset,
90 expression, severity, progression, and overall prognosis^{17,18}. The innate immune defense
91 mechanisms against *C. difficile* and its toxins include the commensal intestinal flora, mucosal
92 barrier, intestinal epithelial cells, and mucosal immune system^{19,20}. TcdA and TcdB have
93 multiple effects on the innate immune system, including inducing expression of numerous pro-
94 inflammatory mediators (e.g., cytokines, chemokines and neuroimmune peptides) and the
95 recruitment and activation of a variety of innate immune cells^{21,22}. Adaptive immunity is also
96 sufficient to provide some protection from CDI, likely via antibody-mediated neutralization of
97 TcdA and TcdB²³⁻²⁶. These immune markers also have the potential to act as clinically useful
98 diagnostic markers of CDI.

99 Exposure to toxinogenic *C. difficile* can lead to a range of clinical outcomes ranging from
100 asymptomatic colonization to mild diarrhea and more severe disease syndromes such as
101 pseudomembranous colitis, toxic megacolon, bowel perforation, sepsis, and death^{27,28}.
102 Asymptomatic *C. difficile* carriage is characterized by *C. difficile* colonization in the absence of
103 symptoms of infection. Previous studies suggest that *C. difficile* asymptomatic carriers have the
104 potential to contribute to *C. difficile* transmission and hospital-onset CDI in inpatient facilities, as
105 carriers can shed spores into the hospital environment^{29,30}.

106 The diagnosis of CDI is based on clinical signs and symptoms in combination with
107 laboratory testing. Several diagnostic laboratory tests are available including enzyme
108 immunoassays (EIA) for TcdA and TcdB, nucleic acid amplification tests (NAAT), selective
109 toxinogenic culture, cell cytotoxicity neutralization assay, and glutamate dehydrogenase EIA³¹⁻³³.
110 However, currently available approaches do not accurately differentiate CDI from diarrhea with
111 another cause in a patient colonized with toxinogenic *C. difficile*. Over-diagnosis of disease
112 could result in overtreatment of CDI, delayed recognition of other causes of illness, and
113 unnecessary antibiotic exposures³⁴.

114 Machine learning has a great impact in many areas of medical research, as it offers a
115 principled approach for developing sophisticated, automatic, and objective algorithms for
116 analysis of complex data. Indeed, previous studies indicate that supervised learning can be
117 successfully employed for clinical disease assessment for diverse disorders including Parkinson's
118 disease³⁵, diabetes³⁶, inflammatory bowel disease³⁷ and glaucoma³⁸. In our previous work, we
119 found that specific immune markers, particularly G-CSF, can be used to distinguish adults with
120 CDI from other groups including asymptomatic carriers and NAAT-negative patients with and
121 without diarrhea³⁹. Here, we integrate the host immune marker data and newly obtained gut
122 microbiome data from subjects of the same cohort to build classification models to optimally
123 distinguish CDI from other groups. Our aim is to identify consistent biological signatures for
124 highly accurate diagnosis of CDI.

125

126 RESULTS

127 Baseline demographic and clinical characteristics of participants

128 Our clinical cohort consists of 243 well-characterized recruited participants, who were divided
129 into four groups (see Materials and Methods)³⁹: (1) Control: subjects without diarrhea and with
130 NAAT-negative stool (n=47); (2) Non-CDI Diarrhea: subjects with diarrhea but NAAT-negative
131 stool (n=44); (3) Asymptomatic Carriage: subjects without diarrhea but with NAAT-positive
132 stool (n=40); (4) CDI: subjects with diarrhea and NAAT-positive stool (n=112). The first three
133 groups can be combined as the Non-CDI group. The entire clinical cohort had a mean \pm SD age
134 of 63.66 ± 14.85 year and was 48.15% female. Demographic data of the cohort are summarized
135 in Table 1. In total, 187 participants (76.95%) had both gut microbiome and immune marker data
136 available (see Table S1).

137

138 Microbial community structure

139 To compare the overall microbial community structure of the four groups, we first calculated the
140 alpha diversity (i.e., the within-sample taxonomic diversity) of each sample at the genus level
141 using four different measures: *taxa richness* (the observed number of different taxa present in the
142 sample), *Chao1* (abundance-based estimator of taxa richness), *Evenness* (the uniformity of the
143 population size of each taxa present in the sample), and *Shannon diversity index* (estimator of
144 taxa richness and evenness: more weight on richness). (See Materials and Methods for detailed
145 definitions of those alpha diversity measures.) As shown in Fig.1, We found that richness indices
146 (taxa richness and Chao1) did not differ significantly among these groups. The gut microbiota of
147 Non-CDI Diarrhea subjects showed lower evenness than that of the Control group. Shannon
148 diversity was significantly lower in the Non-CDI Diarrhea and CDI groups than in the Control
149 group.

150 To determine whether the gut microbial compositions of participants are affected by *C.*
151 *difficile* infection/colonization status, we performed Principal Coordinates Analysis (PCoA) at
152 the genus level using Bray-Curtis dissimilarity (which is a beta diversity measure to quantify the
153 between-sample compositional dissimilarity). We found no distinct clusters corresponding to the

154 four different phenotype groups, implying that the gut microbial compositions of participants
155 from the four groups are not significantly different (Fig. 2A). Interestingly, by directly
156 comparing the beta diversity of each group, we did find that the CDI group displays higher beta
157 diversity than other groups (Fig. 2B), indicating that the microbial compositions of participants
158 within the CDI group vary more prominently than other groups. Permutational multivariate
159 analysis of variance (PERMANOVA) showed that the overall bacterial composition differed
160 significantly among different groups based on the CDI status ($P < 0.001$; Table S2), whereas
161 other host factors such as age, sex, race and ethnicity had no significant effect on the microbiome
162 composition.

163 To identify microbiome markers (i.e., certain taxa with very high discriminatory ability)
164 to differentiate those different phenotype groups, we performed differential abundance analysis.
165 In particular, we used ANCOM⁴⁰ (analysis of composition of microbiomes) with a Benjamini-
166 Hochberg correction, and adjusted for age and sex. We found that the abundances of 15 genera
167 were significantly different between CDI and Asymptomatic Carriage groups (Fig. 3A and Table
168 S3). Among the 15 genera, 4 of them (*Veillonella*, *Enterobacter*, *Granulicatella* and *Dialister*) of
169 these genera were enriched in the CDI group, while the other 11 genera (*Lactococcus*, *Dorea*,
170 *Moryella*, [*Ruminococcus*] *gauvreauii* group, *Stenotrophomonas*, *Agathobacter*, *Blautia*,
171 *Sellimonas*, *Eggerthella*, *Faecalitalea* and *Lachnospiraceae* UCG-008) were enriched in the
172 Asymptomatic Carriage group. We also found 16 differentially abundant genera between the
173 Non-CDI Diarrhea group and the CDI group (Fig. 3B and Table S4). Of these, 10 genera
174 (*Clostridioides*, *Enterobacter*, *Epulopiscium*, *Escherichia-Shigella*, *Eisenbergiella*, *Dialister*,
175 *Ruminiclostridium*, *Fusobacterium*, *Klebsiella* and *Veillonella*) were enriched in the CDI group,
176 and the other 6 genera ([*Eubacterium*] *hallii* group, *Collinsella*, *Agathobacter*, *Dorea*,
177 *Stenotrophomonas* and *Streptococcus*) were enriched in the Non-CDI Diarrhea group. ANCOM
178 analysis also enabled us to identify 40 genera (including *Clostridioides* and *Veillonella*) that have
179 significant differential abundances between the CDI group and the whole Non-CDI group (Fig.
180 3C and Table S5). Note that a total of 6 differentially abundant genera were identified from all
181 the three comparisons: CDI vs. Asymptomatic Carriage; CDI vs. Non-CDI Diarrhea; CDI vs.
182 Non-CDI. Among them, *Veillonella*, *Enterobacter* and *Dialister* were enriched in the CDI group,
183 while *Dorea*, *Stenotrophomonas* and *Agathobacter* were depleted in the CDI group.

184

185 **Microbial correlation networks**

186 To compare the microbial communities of the four groups at the network-level, we constructed
187 the genus-level microbial correlation network for each group using SparCC⁴¹ (sparse correlations
188 for compositional data). We found that the microbial correlation network of the CDI group has
189 quite different structure compared to other groups (Fig. 4). In order to quantify the difference of
190 the network structure, we calculated the number of nodes, number of edges, average degree (the
191 average number of connections per node), graph density (measure of how close the network is to
192 a complete graph), clustering coefficient (measure of how complete the neighborhood of a node
193 is) and modularity (measure of how well a network decomposes into modular communities)

194 (Table S6). In general, compared with the networks of other groups, the network of the CDI
195 group has fewer nodes and edges, lower average degree, but higher modularity. These indicate
196 that the overall microbial correlations in the CDI group are much weaker than those in other
197 groups.

198 To analyze these patterns in more detail, we used NetShift⁴² to identify potentially
199 important “driver” taxa responsible for the change of microbial correlations. This analysis
200 revealed 24 potential driver taxa linked with the change of microbial correlations between CDI
201 and Asymptomatic Carriage groups (Fig. S1). The top driver taxa were *Alistipes*, *Clostridioides*,
202 *Desulfovibrio*, *Eggerthella*, *Erysipelatoclostridium*, *Klebsiella*, *Odoribacter* *Proteus*,
203 *[Ruminococcus]_torques_group*, *Streptococcus*, *Vagococcus* and *Veillonella*. We then identified
204 24 genera as potential driver taxa underlying the change of microbial correlations between CDI
205 and Non-CDI Diarrhea groups (Fig. S2). The top driver taxa were *Alistipes*, *Buttiauxella*,
206 *Citrobacter*, *Clostridium_sensu_stricto_13*, *Desulfovibrio*, *Klebsiella*, *Oscillibacter*,
207 *Phascolarctobacterium*, *Streptococcus* and *Veillonella*. Finally, Netshift analysis revealed 38
208 potential driver taxa underlying the change of microbial correlations between CDI and Non-CDI
209 groups. The top driver taxa were *Bifidobacterium*, *Clostridioides*, *Klebsiella*, *Oscillibacter*,
210 *Streptococcus* and *Veillonella* (Fig. S3). Together, these results suggested that certain bacterial
211 taxa (e.g., *Clostridioides*, *Klebsiella*, *Streptococcus* and *Veillonella*) could play an important role
212 in driving the changes of microbial correlations in subjects with different *C. difficile*
213 infection/colonization status.

214

215 **Host immune markers and CDI**

216 To determine the systemic levels of proinflammatory cytokines in CDI, we measured the
217 circulating levels of granulocyte-colony stimulating factor (G-CSF), interleukin-1 β (IL-1 β), IL-2,
218 IL-4, IL-6, IL-8, IL-10, IL-13, IL-15, monocyte chemoattractant protein-1 (MCP-1), vascular
219 endothelial growth factor-A (VEGF-A), and tumor necrosis factor-alpha (TNF- α) as previously
220 reported³⁹. Serum concentrations of immunoglobulin A (IgA), IgG, and IgM antibodies against
221 *C. difficile* toxin A and toxin B were measured by semi-quantitative enzyme-linked
222 immunosorbent assay (ELISA). We previously demonstrated specific markers of these innate
223 and adaptive immunity that can distinguish CDI from each of the other three groups³⁹. In the
224 current study, we are particularly interested in comparing the CDI group and the combined Non-
225 CDI group. Based on the Mann-Whitney U test, we identified in total 11 immune markers that
226 displayed significantly different concentrations in these two groups, including G-CSF, IL-4, IL-
227 6, IL-8, IL-10, IL-15, TNF- α , MCP1, IgA anti-toxin A and B, and IgG anti-toxin A in blood
228 (Table S7). All of these immune markers had higher concentrations in the CDI group than in the
229 Non-CDI group. Host immune marker variations between samples were evaluated using the
230 Principal Component Analysis (PCA) (Fig. 2C). PCA plot showed no clear clustering of those
231 subjects based on immune marker concentrations. However, boxplot of Euclidean distance of
232 immune marker profiles from CDI patients showed higher within-group variation than that in all
233 the other three groups (Fig. 2D). PERMANOVA analysis indicated that the immune homeostasis

234 was significantly different among different groups based on the CDI status ($P = 0.016$; Table
235 S2). But age, gender, race and ethnicity did not have significant effects on the host immune
236 marker levels.

237

238 **Interplay between gut microbiome and host immune markers**

239 To reveal the interplay between the gut microbiome and the host immune system, we calculated
240 the correlations between microbial compositions (at the genus level) and the circulating levels of
241 host immune markers for each of the four groups (Spearman correlation with Benjamini-
242 Hochberg correction). The results are shown in Fig. 5 and Fig. S4. For the Control group, the
243 most significant associations were identified as *Chiristensenellaceae R-7 group* negatively
244 associated with $\text{TNF}\alpha$, *Bifidobacterium* positively associated with VEGFA and IL-13, *Rothia*
245 positively associated with IL-15, and *Veillonella* positively related with IL-4 (Fig.5A and Fig.
246 S4). For the Non-CDI Diarrhea group, *Ruminococcaceae UCG-011* was negatively correlated
247 with IL-8 and IL-6, *Defluviitaleaceae UCG-011* was positively correlated with IL-1b, and
248 *Blautia* was negatively correlated with MCP1 levels (Fig. 5B). For the Asymptomatic Carriage
249 group, we found that *Lactobacillus* was negatively associated with VEGFA, *Akkermansia* was
250 positively associated with IL-6, and *Enterococcus* was positively related to $\text{TNF}\alpha$ (Fig. 5C). For
251 the CDI group, negative associations involved *Akkermansia* and IL-10, *Lactococcus* and G-CSF,
252 while positive associations involved *Lactobacillus* and IgG and IgA anti-toxin B (Fig. 5D).
253 Interestingly, none of these most significant associations was universally present across different
254 groups. This indicated that the interactions between gut microbiota and host immunological
255 markers can be very sensitive to the status of *C. difficile* colonization and infection. More
256 importantly, this result implies that the integration of gut microbiota and host immune markers
257 might be quite useful for highly accurate diagnosis of CDI.

258

259 **Diagnostic accuracy for CDI classification based on host immune markers and gut 260 microbiota**

261 To determine whether host immune markers or gut microbiota could serve as biomarkers to
262 classify subjects into different groups, we constructed a multi-class classifier based on random
263 forests (RF). One of the most popular performance metrics of a classifier is the Area Under the
264 receiver operating characteristic Curve (AUC). The performance of a multi-class classifier is
265 measured by both micro-average and macro-average AUCs. (For micro-average AUC, we
266 calculated the AUC from the individual true positive rates and false positive rates of the multi-
267 class model. For the macro-average AUC, we calculated the AUC independently for each class
268 and then took the average.) We considered three different feature types: (1) host immune maker
269 concentrations alone; (2) gut microbial compositions alone; and (3) the integration of (1) and (2)
270 in our classification analysis. To eliminate confounding effects, we excluded the genus
271 *Clostridioides* from our classification analysis. The immune marker-based classifier achieved
272 macro-average AUC ~ 0.827 and micro-average AUC ~ 0.828 (Fig. S5A), which are quite
273 comparable to the performance of microbiota-based classifier (Fig. S5B). Interestingly,

274 integrating immune marker with gut microbiota showed much better classification performance
275 (macro-average AUC \sim 0.926 and micro-average AUC \sim 0.869) (Fig. S5C).

276 We further performed binary classifications to distinguish CDI subjects from
277 Asymptomatic Carriage, Non-CDI Diarrhea, and Non-CDI subjects, using different feature types
278 (Fig. 6). The goal of this analysis was to assess whether any single taxon or immune marker
279 could reliably differentiate CDI status. The importance of each feature was quantified by the
280 Mean Decrease in Accuracy (MDA) of the classifier due to the exclusion (or permutation) of this
281 feature. The more the accuracy of the classifier decreases due to the exclusion (or permutation)
282 of a single feature, the more important that feature is deemed for classification of the data.

283 In the classification of CDI vs. Asymptomatic Carriage, we found that G-CSF and
284 *Moryella* were the most important immune and microbial features, respectively (Fig. S6:A-B).
285 But the classification based on G-CSF (or *Moryella*) alone did not yield very high performance:
286 mean AUC \sim 0.817 (or 0.701), respectively (Fig. 6:A1-A2). When we used all the immune
287 markers (or all the genera) as features, we achieved mean AUC \sim 0.867 (or 0.805), respectively
288 (Fig. 6:A3-A4). Interestingly, when we integrated all the host immune markers and gut microbial
289 composition data together, we achieved a much higher performance with mean AUC \sim 0.900
290 (Fig. 6:A5). In order to select a subset of features that is as discriminatory as the whole set of
291 features, we followed the “1-SE” rule (i.e., one chooses the model with fewest features such that
292 its classification performance is less than one standard error away from that of the model with all
293 the features), and selected the following 4 features: 2 bacterial genera (*Moryella* and *Veillonella*)
294 and 2 immune markers (G-CSF and IL-6) in classifying CDI and Asymptomatic Carriage groups
295 (Fig. S6:G-J). The RF classifier with those selected features displayed an outstanding
296 classification performance, with mean AUC \sim 0.916 (Fig. 6:A6). Note that a significant negative
297 correlation between *Moryella* and G-CSF was found in the Asymptomatic Carriage group (Fig.
298 5C), which might contribute to the outstanding performance of the RF classifier with *Moryella*
299 and G-CSF as selected features.

300 In the classification of CDI vs. Non-CDI Diarrhea groups, we found that G-CSF and
301 *[Eubacterium]_hallii_group* are the top immune and microbial features, respectively (Fig. S6:C-
302 D). But the classification based on G-CSF (or *[Eubacterium]_hallii_group*) alone did not
303 perform very well: mean AUC \sim 0.747 (or \sim 0.630), respectively (Fig. 6:B1-B2). When we used
304 all the immune marker (or all the microbial genera) as features, we achieved mean AUC \sim 0.851
305 (or \sim 0.884), respectively (Fig. 6:B3-B4). By integrating all features from both host immune
306 marker and gut microbial genera, we further improved the classification performance to mean
307 AUC \sim 0.918 (Fig. 6:B5). Following the “1-SE” rule, we selected the following 5 features: 3
308 genera: *Enterococcus*, *Epulopiscium* and *[Eubacterium]_hallii_group*; and 2 immune markers:
309 G-CSF and IgA anti-toxin A (Fig. S6:H-K). The RF classifier with those selected features
310 achieved mean AUC \sim 0.917 (Fig. 6:B6), which is quite comparable to that of using all the
311 features. Note that *Enterococcus* was found to be significantly associated with G-CSF in the
312 Non-CDI Diarrhea group (Fig. 5B). This might partially explain the outstanding performance of
313 the RF classifier with *Enterococcus* and G-CSF as selected features.

314 In the classification of CDI vs. Non-CDI groups, we found that G-CSF and *Curvibacter*
 315 are the top immune and microbial features, respectively (Fig. S6:E-F). Classification based on G-
 316 CSF (or *Curvibacter*) alone achieved mean AUC ~ 0.802 (or ~ 0.683), respectively (Fig. 6:C1-
 317 C2). When we used all the immune marker (or all the microbial genera) as features, we achieved
 318 mean AUC ~ 0.878 (or ~ 0.903), respectively (Fig. 6:C3-C4). Integrating all features from both
 319 host immune marker and gut microbial genera, we further improved the classification
 320 performance to mean AUC ~ 0.941 (Fig. 6:C5). Following the “1-SE” rule, we selected the
 321 following 10 features: 6 genera: *Stenotrophomonas*, *Curvibacter*, *Enterobacter*, *Anaerobacillus*,
 322 *Fusobacterium* and *Veillonella*; and 4 immune markers: G-CSF, IL-6, TNF- α and IgA anti-toxin
 323 B (Fig. S6:I-L). Classification with those well selected features achieved mean AUC ~ 0.929
 324 (Fig. 6:C6).

325

326 Using symbolic classification to derive diagnostic scores.

327 The outstanding classification results based on well-selected features prompt us to derive simple
 328 mathematical models for CDI diagnosis. To achieve that, we leveraged symbolic classification
 329 (SC)^{43,44}, a genetic programming technique that automatically searches the space of
 330 mathematical expressions to find the model that best fits a given dataset. The fitness function in
 331 SC is a maximization function, and the number of generations is chosen based on the saturation
 332 of the fitness score (Fig. S7). Using the same set of selected features and trained with the entire
 333 dataset, the SC model outperformed logistic regression (LR) in differentiating CDI from
 334 Asymptomatic Carriage (or Non-CDI Diarrhea, or Non-CDI), based on various performance
 335 metrics: Accuracy, Precision, Recall and F1-score (see Table 2).

336 Indeed, as shown in Table 2, we derived a simple SC model with selected features,
 337 reaching a very high accuracy (0.896) in distinguishing CDI subjects from Asymptomatic
 338 Carriage. Basically, for each subject i , we calculate the diagnostic score $f(i)$ that will be used
 339 for CDI diagnosis: the class of subject i is CDI if $f(i) > 0$; Asymptomatic Carriage, if $f(i) \leq 0$.
 340 Here,

$$341 \quad f(i) = x_{\text{GCSF}} * x_{\text{Veillonella}} (x_{\text{GCSF}}^3 - 0.2 * x_{\text{Moryella}} + 0.4) + 1.1 * x_{\text{GCSF}} - 0.1 * x_{\text{IL6}} -$$

$$342 \quad 18.25, \tag{1}$$

343 with x_a representing the abundance or concentration of feature- a in subject- i . Similarly, we
 344 derived a SC model with accuracy of 0.900 in distinguishing CDI (if $f(i) > 0$) from Non-CDI
 345 Diarrhea (if $f(i) \leq 0$) with the diagnostic score

$$346 \quad f(i) = x_{\text{Enterococcus}} * x_{\text{IgA_toxA}} (0.5 * x_{\text{Epulopiscium}} - 1) + x_{\text{[Eubacterium]_hallii_group}} (0.02 *$$

$$347 \quad -x_{\text{GCSF}}) + x_{\text{IgA_toxA}} \left(1 - \frac{10}{x_{\text{GCSF}}} \right) - \frac{0.003}{x_{\text{Enterococcus}}}. \tag{2}$$

348 Finally, we derived a SC model with accuracy of 0.882 in distinguishing CDI (if $f(i) > 0$) from
 349 Non-CDI (if $f(i) \leq 0$) with the diagnostic score

$$350 \quad f(i) = x_{\text{GCSF}} * x_{\text{IgA_toxB}} (0.2 * x_{\text{Anaerobacillus}} * x_{\text{GCSF}} * x_{\text{Stenotrophomonas}} * x_{\text{TNF}\alpha} + 0.04 *$$

$$351 \quad x_{\text{Curvibacter}} * x_{\text{GCSF}} + 0.3 * x_{\text{Enterobacter}}^4 * x_{\text{GCSF}} * x_{\text{Veillonella}}) + x_{\text{Fusobacterium}} *$$

$$x_{\text{GCSF}}(0.5 * x_{\text{Curvibacter}} + x_{\text{GCSF}} * x_{\text{Stenotrophomonas}}) + x_{\text{Curvibacter}}(0.1 * x_{\text{IL6}} - x_{\text{Anaerobacillus}}) + x_{\text{Stenotrophomonas}}(x_{\text{Stenotrophomonas}} - 2). \quad (3)$$

To ensure the SC models learned from the entire dataset are not overfitting, we performed cross-validation by randomly splitting the dataset to form a training set (80% of the data) and a held-out test set (20% of the data) in 10 different ways. Each time, for each classification task, we learned the SC model from the training dataset and evaluated it on the test dataset. Due to the different training sets, SC will derive different mathematical formulas (i.e., diagnostic scores). However, those SC models learned from different training datasets demonstrated quite robust performance in terms of Accuracy, Precision, Recall and F1-score (see Table S8). More importantly, even trained with less data, the SC models still outperformed LR models learned from the entire dataset.

These SC models consisted of explicit mathematical equations, which are more transparent than black-box classifiers such as RF. At the same time, the SC models are also more accurate than traditional classifiers (such as LR). The transparency and high accuracy highlight the importance of SC models in the clinical diagnosis of CDI.

DISCUSSION

Current methods for CDI diagnosis are unable to combine high sensitivity and high clinical specificity, which can result in either underdiagnosis or overdiagnosis of CDI⁴⁵. A more accurate diagnostic approach for CDI could optimize therapeutic decision-making and reduce transmission. Here, we employed 16S rRNA gene sequencing to profile the gut microbial compositions and combined the gut microbiome data with data from a broad panel of innate and adaptive host immune response markers to investigate the potential roles of these markers in the diagnosis of CDI. We demonstrated that the combination of host immune markers and gut microbial data can provide a potential route to optimize CDI diagnosis. Importantly, this work derived specific diagnostic models (in terms of mathematic equations) that yielded robust accuracy in differentiating CDI subjects from Asymptomatic Carriage, Non-CDI Diarrhea and Non-CDI groups.

Taxonomic diversity is a fundamental property of ecological systems. It is generally believed to be an important determinant of the structure and functioning of ecological communities^{46,47}. Diversity indices have been routinely calculated in the study of human microbiome⁴⁸. Consistent with previous studies⁴⁹⁻⁵², we found that the gut microbiomes of CDI patients were characterized by lower Shannon diversity than that of the Control group. Interestingly, we observed an increased variation of both immune markers and gut microbial compositions in the CDI group with respect to other studied groups. This suggests that CDI is characterized by a significantly less stable microbiome and immune homeostasis. Our findings are in line with the Anna Karenina principle, which suggests that CDI linked changes in the microbiome and immune homeostasis are likely stochastic, leading to community instability⁵³⁻⁵⁵.

We were able to identify several candidate driver taxa (e.g., *Desulfovibrio*, *Klebsiella*, *Streptococcus* and *Veillonella*) that played a key role in driving the changes of microbial

392 correlation networks between CDI and Asymptomatic Carriage (or Non-CDI Diarrhea, Non-
393 CDI) groups. Among those driver taxa, *Desulfovibrio* has previously been shown to have a
394 pathogenic role in ulcerative colitis due to its ability to generate sulfides⁵⁶. *Streptococcus* has
395 previously been shown to produce lactate thus impacting *C. difficile* TcdA production and *tcdA*
396 expression to alleviate CDI⁵⁷. *Klebsiella* is a Gram-negative bacterium that cause different types
397 of healthcare-associated infections including pneumonia, bloodstream infections, and
398 meningitis⁵⁸. *Klebsiella* bacteria have been increasingly shown to develop antimicrobial
399 resistance, most recently to the class of antibiotics known as carbapenems^{59,60}. It is thus possible
400 that the CDI pathogenesis is further enforced by the enrichment of antagonistic bacteria present
401 in the gut microbiome of CDI subjects. In addition, our analysis demonstrated that the
402 associations between host immunological markers and gut microbial compositions in the CDI
403 group were dramatically different from those in other groups. However, further investigations are
404 needed to determine whether these alterations are integral to the CDI pathogenesis.

405 The diagnosis of CDI remains challenging, especially the ability to distinguish CDI and
406 *C. difficile* colonization⁶¹⁻⁶³. To address this issue, we developed classification models aimed at
407 differentiating CDI status based on host immune markers and gut microbiome data. We excluded
408 the genus *Clostridioides* in further classification analysis to eliminate confounding effects.
409 Evaluating the classification performance of host immune markers or/and microbiome data in
410 multi-class models, it appears that a combination of host immune markers and gut microbiome
411 data can further improve the accuracy of classification. More specifically, we were able to
412 identify specific immune and microbial features that could accurately distinguish CDI subjects
413 from Asymptomatic Carriage, Non-CDI Diarrhea, and Non-CDI subjects. In addition, most of
414 the selected features identified by feature selection were also differentially abundant genera and
415 differentially expressed immune markers.

416 From the classification of CDI and Asymptomatic Carriage, we were able to select a few
417 features with outstanding discriminability, including *Veillonella* and *Moryella*. Interestingly, a
418 positive relationship between *Veillonella* and CDI has been identified in recent studies⁶⁴⁻⁶⁷. An
419 important role for *Veillonella* in CDI is supported by the fact that *Veillonella* species were
420 associated with low coprostanol levels that correlated strongly with CDI⁶⁴. A similar negative
421 relationship between *Moryella* species and CDI has previously been observed⁶⁸. *Enterococcus*, a
422 feature selected from the classification of CDI vs. Non-CDI Diarrhea, has been reported to be
423 associated with CDI due to vancomycin resistance⁶⁹. Consistent with the findings from previous
424 reports^{70,71}, *Epulopiscium* was significantly enriched in the CDI group and played an important
425 role in differentiating this comparison. Among those features selected from the classification of
426 CDI and Non-CDI groups, *Enterobacter* and *Fusobacterium* have been considered as
427 opportunistic pathogens involved in multiple diseases^{72,73}.

428 Machine learning method has the potential to identify biomarkers and aid in the diagnosis
429 of many diseases. However, the learnt relationships between predictors and outcome are
430 typically non-transparent, especially non-linear methods (i.e., decision tree learning). Previous
431 study has shown that an interpretable trees framework can extract, measure, prune, select, and

432 summarize rules from a tree ensemble, and calculates frequent variable interactions⁷⁴. However,
433 these rules from tree ensembles are still too complicated to be clinically meaningful. Classical
434 logistic regression, is one of the most common machine learning models in medicine⁷⁵. The main
435 drawback of LR is its failure to solve non-linear problems and it underperforms where there are
436 multiple or non-linear decision boundaries⁷⁶. Furthermore, the log odds scale in LR is hard to
437 interpret⁷⁷. Symbolic classification based on genetic programming is an automated technique to
438 derive formulas from features for classification purpose⁷⁸. Using the selected integrated features
439 from the random forests model, we demonstrated that the mathematical formulas automatically
440 derived from symbolic classification have robust diagnostic accuracy to differentiate CDI
441 patients from Asymptomatic Carriage (or Non-CDI Diarrhea, and Non-CDI groups).
442 Specifically, symbolic classification provides explicit mathematic formulas as its output, which
443 significantly improves the transparency of the learned relationship between predictors and
444 outcomes. These results hold translational promise in clinical diagnosis of CDI. Further external
445 validation of the derived formulas will require a different cohort with the same inclusion criteria
446 as ours. This is beyond the scope of the current work.

447 We previously demonstrated the potential clinical utility of a specific immunological
448 biomarker (G-CSF) for CDI diagnosis³⁹. This study leverages the newly obtained gut
449 microbiome data from the same unique and well-characterized study cohort, allowing us to study
450 integrated host immune marker and gut microbiome signatures. The fundamental differences
451 between this study and our previous one are the clinical utilization of integrated immune and
452 microbiome signatures to distinguish CDI patients from Asymptomatic Carriage, Non-CDI
453 Diarrhea, and Non-CDI groups, and to derive diagnostic scores for CDI diagnosis. We believe
454 that this study sets the stage to explore the potential role of an immune and microbiome-based
455 test for CDI diagnosis. Of course, observed associations and selected features do not offer any
456 causal relationships. Prospective studies are needed to validate the mechanism underlying the
457 relationship between these selected features/biomarkers and the CDI infection/colonization
458 status. The 16S rRNA sequencing may not have captured additional insights associated with the
459 disease status available at the species or strain level. Further studies are needed to validate the
460 clinical utility of the proposed biomarkers by metagenomics sequencing as well as
461 metatranscriptomics, metaproteomics and metabolomics.

462 In summary, leveraging a well-characterized clinical cohort, we provided strong evidence
463 that integrating gut microbiome and host immune signatures can significantly improve the CDI
464 diagnosis. In particular, these results demonstrate that knowledge of gut microbial compositions
465 in combination with host immune markers is beneficial in generating clinically relevant machine
466 learning models for disease diagnosis. Indeed, the machine learning models show high diagnostic
467 accuracy in differentiating true CDI from asymptomatic carriage of *C. difficile* and from non-*C.*
468 *difficile* diarrhea, which are areas where current laboratory testing for CDI lacks adequate
469 clinical specificity.

470
471

472 **MATERIALS AND METHODS**

473 **Study cohort**

474 The background and design of this cohort has been described in details previously⁶². Concisely,
475 we have four groups associated with different *C. difficile* infection/colonization statuses: (1)
476 Control: subjects without diarrhea who had screened as eligible for the asymptomatic carriage
477 group (see below) but were NAAT-negative on research stool testing; (2) Non-CDI Diarrhea:
478 subjects with diarrhea but NAAT-negative stool on clinical testing; (3) Asymptomatic Carriage:
479 subjects were admitted for at least 72 hours, had received at least one dose of an antibiotic within
480 the past 7 days, did not have diarrhea in the 48 hours prior to stool sample collection, had
481 positive NAAT results on research stool testing and were not treated for CDI; (4) CDI: inpatients
482 with positive clinical stool NAAT result, diarrhea, and a decision to treat for CDI. All subjects
483 were adults (age \geq 18 years old). Clinical serum samples were collected as discards within 24
484 hours of stool sample collection. In our previous study³⁹, the four groups were named as (1) “no
485 Diarrhea NAAT-Negative” = Control; (2) diarrhea NAAT-negative = Non-CDI Diarrhea; (3)
486 Carrier-NAAT = Asymptomatic Carriage and (4) CDI-NAAT = CDI. In this work, for simplicity
487 we used the simpler and more clearly descriptive titles.

488

489 **Serum immune marker measurement**

490 The measurement of host serum cytokines concentrations of IL-2, IL-4, IL-6, IL-8, IL-10, IL-13,
491 IL-15, IL-1 β , G-CSF, IL-1 β , MCP-1, VEGF-A, and TNF- α was performed using a Milliplex
492 magnetic bead kit and Luminex analyzer (MAGPIX) (Millipore Sigma, Inc., Burlington, MA) as
493 per the manufacturer’s instructions. Purified toxin A and B were separately prepared from *C.*
494 *difficile* strain VPI 10463 (American Type Culture Collection 43255-FZ, Manassas, VA). Serum
495 antibody (IgA, IgG, and IgM) levels against *C. difficile* toxins A and B were measured by semi-
496 quantitative enzyme-linked immunosorbent assay (ELISA). All the experimental details have
497 been reported previously^{39,62}.

498

499 **Fecal DNA extraction and bacterial 16S rRNA sequencing data analysis**

500 Stool DNA was extracted using the DNeasy PowerSoil Pro Kit (Qiagen, cat# 12888-100) in a
501 QiaCube automated DNA extraction system (Qiagen) according to instructions. Briefly, 250mg
502 stool was transferred into a PowerBead Pro Tube provided with the kit and 200 μ g RNaseA and
503 800 μ l of CD1 solution were added. Tubes were vortexed briefly, transferred into an adapter, and
504 then vortexed at maximum speed for 10 min. Tubes were centrifuged at 15,000 x g for 1 min and
505 about 500–600 μ l supernatant was used for DNA extraction according to instructions. DNA were
506 eluted in 70 μ l elution solution C6 and stored at -80°C until use. 16S rRNA microbiome
507 characterization was performed by sequencing the V4 region of the 16S rRNA gene using the
508 Illumina MiSeq.⁷⁹ Each sample was amplified using a barcoded primer, which yielded a unique
509 sequence identifier tagged onto each individual sample library. Illumina-based sequencing
510 yielded greater than 15,000 reads per sample. CLC Genomics Workbench version 12 (Qiagen)
511 was used for OTU clustering and generation of abundance tables. Analyses were performed

512 using the tutorial “OTU Clustering Step by Step” updated September 2, 2019 and available on
513 the Qiagen website:

514 https://resources.qiagenbioinformatics.com/tutorials/OTU_Clustering_Steps.pdf

515

516 **Microbial diversity and differential abundance analysis**

517 Both alpha and beta diversity measures were calculated at the genus level using the vegan:
518 Community Ecology Package in R (<https://CRAN.R-project.org/package=vegan>). Measures of
519 alpha diversity included: the richness S (the number of taxa present in the community/sample),
520 Chao 1 index $S_{chao1} = S + \frac{F_1(F_1-1)}{2(F_2+1)}$, Shannon index $H = -\sum_{i=1}^S p_i \log p_i$, and evenness $J =$
521 $H/\log S$. Here, F_1 and F_2 are the count of singletons and doubletons, respectively, and p_i is the
522 relative abundance of taxon- i in the community. For beta diversity, we used the Bray-Curtis
523 dissimilarity measure, which was also used in the Principal Coordinates Analysis (PCoA). We
524 applied principal component analysis (PCA) on the expression levels of all immune markers
525 based on the Euclidean distance.

526 Difference in microbiome compositions and immune expression levels by CDI status
527 (i.e., different groups) and other covariates (i.e., age, sex, race and ethnicity) were tested by the
528 permutational multivariate analysis of variance (PERMANOVA) using the “adonis” function in
529 the vegan R package. All PERMANOVA tests were performed with the default 999
530 permutations based on the Bray-Curtis dissimilarity and Euclidean distance for microbial
531 composition and immune marker data, respectively. Note that in the PERMANOVA tests, we
532 only included subjects with known information of age, sex, race and ethnicity.

533 For differential abundance analysis, we used ANCOM⁴⁰ (analysis of composition of
534 microbiomes), with a Benjamini–Hochberg correction at 5% level of significance, and adjusted
535 for age and sex. The Mann–Whitney U test was used to compare the difference of immune
536 marker levels between different groups.

537

538 **Microbial correlation network analysis**

539 The microbial correlation networks were constructed using SparCC⁴¹ (sparse correlations for
540 compositional data, <https://github.com/luispedro/sparcc>). Significant interactions were
541 determined by the bootstrapped results ($N = 100$) using the script PseudoPvals in SparCC.
542 Significant correlations with absolute sparse correlations ≥ 0.3 were visualized using Gephi
543 (<https://gephi.org/>). We also used the NetShift⁴² (<https://web.rniapps.net/netshift>) to identify
544 potential “driver” taxa underlying the differences of microbial correlation networks associated
545 with CDI and Asymptomatic Carriage (or Non-CDI Diarrhea, and Non-CDI). The key driver
546 taxa were identified based on the neighbor shift (NESH) score, Jaccard Index and delta
547 betweenness (ΔB)⁴².

548

549 **Microbiome-Immune marker association analysis**

550 Associations between the gut microbiota and host immune markers were quantified by Spearman
551 correlation coefficients in combination with Benjamini-Hochberg FDR correction to account for

552 multiple hypothesis testing (significance threshold $\alpha \leq 0.05$). All included genera were required
553 to be detected in $\geq 50\%$ of all samples in each group.

554

555 **Classification with Random Forests model**

556 To build a classification model capable of testing the overall contribution of immunological or
557 microbial data in distinguishing the CDI status, we developed a multi-class random forests (RF)
558 classifier. The data is split into a training set and a test set, with 70% of the data forming the
559 training data and the remaining 30% forming the test set. The performance of the multi-class
560 model was measured by micro-average and macro-average AUC. A macro-average score
561 computed the metric independently for each group and then was averaged across all levels
562 regardless of the number of samples in each group, whereas a micro-average will aggregate the
563 contributions of all groups to compute the average metric.

564 To determine whether more specific host immune markers or gut microbial taxa could
565 differentiate CDI subjects from Asymptomatic Carriage, Non-CDI Diarrhea and Non-CDI
566 groups, we constructed the binary classifiers based on RF models with integrated immune
567 markers and microbiome data. The performance of the classifiers were evaluated by a 5-fold
568 cross validation. In order to reduce computation complexity and feature redundancy, a feature
569 selection procedure was performed as follows. We first ranked all the features based on their
570 mean decrease accuracy (MDA). Then we followed the “1-SE strategy” to select the minimum
571 set of top features whose mean AUC is within one standard error of the mean AUC from the
572 model with all of the features.

573

574 **Symbolic classification with genetic programming**

575 Genetic programming (GP) is a genetic algorithm that searches the space of mathematical
576 equations without any constraints on their forms⁸⁰. GP involves reproduction, random mutation,
577 crossover, a fitness function, and multiple generations of evolution of a population of computer
578 programs to resolve a given task. GP is commonly used to investigate a functional relationship
579 (i.e., a mathematical formula) between features in data (symbolic regression: SR) or to group
580 data into categories (symbolic classification: SC). We employed Karoo GP⁸¹, a genetic
581 programming application suite written in Python that support both SR and SC analysis, to derive
582 simple formulas for CDI diagnosis. We performed a random data-split to create a training set
583 (80% of the data) and a held-out test set (20% of the data) for ten times, which were used to
584 evaluate the SC performance. Due to the different training sets, SC will derive different
585 formulas, but their classification performances (in terms of Accuracy, Precision, Recall, F1-
586 score) are quite comparable (Table S8). The formulas shown in Table 2 were derived based on
587 the whole dataset. The Karoo GP was used with the following settings: (1) the fitness function
588 (Kernel) is c (representing “classification”); (2) the type of tree is r (ramped half/half); (3) the
589 maximum tree depth for the initial population is 6; (4) the number of trees per generation is 100;
590 (5) the maximum number of generations is 190 (based on the converging results shown in Fig.
591 S7); (6) constants include 0.1, 0.2, 0.3, 0.4 and 0.5; and (7) all other parameters are set as default

592 values. The fitness function in SC is a maximization function, which will seek the highest fitness
593 score among the trees in each generation. The sign of the final formula $f(i)$ will be used for CDI
594 diagnosis: the class of subject i is CDI if $f(i) > 0$; or Asymptomatic Carriage (or Non-CDI
595 Diarrhea, Non-CDI) if $f(i) \leq 0$.

596 To demonstrate the advantage of SC, for each classification task (i.e., CDI vs.
597 Asymptomatic Carriage, CDI vs. Non-CDI Diarrhea, and CDI vs. Non-CDI), we also performed
598 logistic regression (LR) using the same set of selected features as used in SC (Table 2). The LR
599 models were constructed using the `glm()` function in R. The class of subject i is CDI if
600 $p(i) \geq 0.5$; or Asymptomatic Carriage (or Non-CDI Diarrhea, Non-CDI) if $p(i) < 0.5$.

601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627

628 REFERENCES AND NOTES

- 629 1 Lessa, F. C. *et al.* Burden of Clostridium difficile infection in the United States. *N Engl J*
630 *Med* **372**, 825-834, doi:10.1056/NEJMoa1408913 (2015).
- 631 2 Depestel, D. D. & Aronoff, D. M. Epidemiology of Clostridium difficile infection. *J*
632 *Pharm Pract* **26**, 464-475, doi:10.1177/0897190013499521 (2013).
- 633 3 McDonald, L. C. *et al.* Clinical Practice Guidelines for Clostridium difficile Infection in
634 Adults and Children: 2017 Update by the Infectious Diseases Society of America (IDSA)
635 and Society for Healthcare Epidemiology of America (SHEA). *Clin Infect Dis* **66**, e1-
636 e48, doi:10.1093/cid/cix1085 (2018).
- 637 4 Bagdasarian, N., Rao, K. & Malani, P. N. Diagnosis and treatment of Clostridium
638 difficile in adults: a systematic review. *JAMA* **313**, 398-408,
639 doi:10.1001/jama.2014.17103 (2015).
- 640 5 Rineh, A., Kelso, M. J., Vatansever, F., Tegos, G. P. & Hamblin, M. R. Clostridium
641 difficile infection: molecular pathogenesis and novel therapeutics. *Expert Rev Anti Infect*
642 *Ther* **12**, 131-150, doi:10.1586/14787210.2014.866515 (2014).
- 643 6 Stevens, V., Dumyati, G., Fine, L. S., Fisher, S. G. & van Wijngaarden, E. Cumulative
644 antibiotic exposures over time and the risk of Clostridium difficile infection. *Clin Infect*
645 *Dis* **53**, 42-48, doi:10.1093/cid/cir301 (2011).
- 646 7 Slimings, C. & Riley, T. V. Antibiotics and hospital-acquired Clostridium difficile
647 infection: update of systematic review and meta-analysis. *J Antimicrob Chemother* **69**,
648 881-891, doi:10.1093/jac/dkt477 (2014).
- 649 8 Lewis, B. B. *et al.* Loss of Microbiota-Mediated Colonization Resistance to Clostridium
650 difficile Infection With Oral Vancomycin Compared With Metronidazole. *J Infect Dis*
651 **212**, 1656-1665, doi:10.1093/infdis/jiv256 (2015).
- 652 9 Becattini, S., Taur, Y. & Pamer, E. G. Antibiotic-Induced Changes in the Intestinal
653 Microbiota and Disease. *Trends Mol Med* **22**, 458-478,
654 doi:10.1016/j.molmed.2016.04.003 (2016).
- 655 10 Buffie, C. G. *et al.* Profound alterations of intestinal microbiota following a single dose
656 of clindamycin results in sustained susceptibility to Clostridium difficile-induced colitis.
657 *Infect Immun* **80**, 62-73, doi:10.1128/IAI.05496-11 (2012).
- 658 11 Perez-Cobas, A. E. *et al.* Structural and functional changes in the gut microbiota
659 associated to Clostridium difficile infection. *Front Microbiol* **5**, 335,
660 doi:10.3389/fmicb.2014.00335 (2014).
- 661 12 Theriot, C. M. *et al.* Antibiotic-induced shifts in the mouse gut microbiome and
662 metabolome increase susceptibility to Clostridium difficile infection. *Nat Commun* **5**,
663 3114, doi:10.1038/ncomms4114 (2014).
- 664 13 Leffler, D. A. & Lamont, J. T. Clostridium difficile Infection. *N Engl J Med* **373**, 287-
665 288, doi:10.1056/NEJMc1506004 (2015).
- 666 14 Genth, H., Dreger, S. C., Huelsenbeck, J. & Just, I. Clostridium difficile toxins: more
667 than mere inhibitors of Rho proteins. *Int J Biochem Cell Biol* **40**, 592-597,
668 doi:10.1016/j.biocel.2007.12.014 (2008).
- 669 15 Sun, X., He, X., Tzipori, S., Gerhard, R. & Feng, H. Essential role of the
670 glucosyltransferase activity in Clostridium difficile toxin-induced secretion of TNF-alpha
671 by macrophages. *Microb Pathog* **46**, 298-305, doi:10.1016/j.micpath.2009.03.002 (2009).

- 672 16 Riegler, M. *et al.* Clostridium difficile toxin B is more potent than toxin A in damaging
673 human colonic epithelium in vitro. *J Clin Invest* **95**, 2004-2011, doi:10.1172/JCI117885
674 (1995).
- 675 17 Sun, X. & Hirota, S. A. The roles of host and pathogen factors and the innate immune
676 response in the pathogenesis of Clostridium difficile infection. *Mol Immunol* **63**, 193-202,
677 doi:10.1016/j.molimm.2014.09.005 (2015).
- 678 18 Kelly, C. P. & Kyne, L. The host immune response to Clostridium difficile. *J Med*
679 *Microbiol* **60**, 1070-1079, doi:10.1099/jmm.0.030015-0 (2011).
- 680 19 Bibbo, S. *et al.* Role of microbiota and innate immunity in recurrent Clostridium difficile
681 infection. *J Immunol Res* **2014**, 462740, doi:10.1155/2014/462740 (2014).
- 682 20 Iacob, S., Iacob, D. G. & Luminos, L. M. Intestinal Microbiota as a Host Defense
683 Mechanism to Infectious Threats. *Front Microbiol* **9**, 3328,
684 doi:10.3389/fmicb.2018.03328 (2018).
- 685 21 Madan, R. & Petri, W. A., Jr. Immune responses to Clostridium difficile infection.
686 *Trends Mol Med* **18**, 658-666, doi:10.1016/j.molmed.2012.09.005 (2012).
- 687 22 Sun, X., Savidge, T. & Feng, H. The enterotoxicity of Clostridium difficile toxins. *Toxins*
688 (*Basel*) **2**, 1848-1880, doi:10.3390/toxins2071848 (2010).
- 689 23 Kyne, L., Warny, M., Qamar, A. & Kelly, C. P. Association between antibody response
690 to toxin A and protection against recurrent Clostridium difficile diarrhoea. *Lancet* **357**,
691 189-193, doi:10.1016/S0140-6736(00)03592-3 (2001).
- 692 24 Wilcox, M. H. *et al.* Bezlotoxumab for Prevention of Recurrent Clostridium difficile
693 Infection. *N Engl J Med* **376**, 305-317, doi:10.1056/NEJMoa1602615 (2017).
- 694 25 Giannasca, P. J. *et al.* Serum antitoxin antibodies mediate systemic and mucosal
695 protection from Clostridium difficile disease in hamsters. *Infect Immun* **67**, 527-538
696 (1999).
- 697 26 Johnston, P. F., Gerding, D. N. & Knight, K. L. Protection from Clostridium difficile
698 infection in CD4 T Cell- and polymeric immunoglobulin receptor-deficient mice. *Infect*
699 *Immun* **82**, 522-531, doi:10.1128/IAI.01273-13 (2014).
- 700 27 Rupnik, M., Wilcox, M. H. & Gerding, D. N. Clostridium difficile infection: new
701 developments in epidemiology and pathogenesis. *Nat Rev Microbiol* **7**, 526-536,
702 doi:10.1038/nrmicro2164 (2009).
- 703 28 Schaffler, H. & Breitruck, A. Clostridium difficile - From Colonization to Infection.
704 *Front Microbiol* **9**, 646, doi:10.3389/fmicb.2018.00646 (2018).
- 705 29 Blixt, T. *et al.* Asymptomatic Carriers Contribute to Nosocomial Clostridium difficile
706 Infection: A Cohort Study of 4508 Patients. *Gastroenterology* **152**, 1031-1041 e1032,
707 doi:10.1053/j.gastro.2016.12.035 (2017).
- 708 30 Guerrero, D. M. *et al.* Asymptomatic carriage of toxigenic Clostridium difficile by
709 hospitalized patients. *J Hosp Infect* **85**, 155-158, doi:10.1016/j.jhin.2013.07.002 (2013).
- 710 31 Tenover, F. C., Baron, E. J., Peterson, L. R. & Persing, D. H. Laboratory diagnosis of
711 Clostridium difficile infection can molecular amplification methods move us out of
712 uncertainty? *J Mol Diagn* **13**, 573-582, doi:10.1016/j.jmoldx.2011.06.001 (2011).
- 713 32 Burnham, C. A. & Carroll, K. C. Diagnosis of Clostridium difficile infection: an ongoing
714 conundrum for clinicians and for clinical laboratories. *Clin Microbiol Rev* **26**, 604-630,
715 doi:10.1128/CMR.00016-13 (2013).

- 716 33 Musher, D. M. *et al.* Detection of *Clostridium difficile* toxin: comparison of enzyme
717 immunoassay results with results obtained by cytotoxicity assay. *J Clin Microbiol* **45**,
718 2737-2739, doi:10.1128/JCM.00686-07 (2007).
- 719 34 Kociolek, L. K. *et al.* Impact of a Healthcare Provider Educational Intervention on
720 Frequency of *Clostridium difficile* Polymerase Chain Reaction Testing in Children: A
721 Segmented Regression Analysis. *J Pediatric Infect Dis Soc* **6**, 142-148,
722 doi:10.1093/jpids/piw027 (2017).
- 723 35 Abos, A. *et al.* Discriminating cognitive status in Parkinson's disease through functional
724 connectomics and machine learning. *Sci Rep* **7**, 45347, doi:10.1038/srep45347 (2017).
- 725 36 Dagliati, A. *et al.* Machine Learning Methods to Predict Diabetes Complications. *J*
726 *Diabetes Sci Technol* **12**, 295-302, doi:10.1177/1932296817706375 (2018).
- 727 37 Mossotto, E. *et al.* Classification of Paediatric Inflammatory Bowel Disease using
728 Machine Learning. *Sci Rep* **7**, 2427, doi:10.1038/s41598-017-02606-2 (2017).
- 729 38 Kim, S. J., Cho, K. J. & Oh, S. Development of machine learning models for diagnosis of
730 glaucoma. *PLoS One* **12**, e0177726, doi:10.1371/journal.pone.0177726 (2017).
- 731 39 Kelly, C. P. *et al.* Host Immune Markers Distinguish *Clostridioides difficile* Infection
732 From Asymptomatic Carriage and Non-*C. difficile* Diarrhea. *Clin Infect Dis*,
733 doi:10.1093/cid/ciz330 (2019).
- 734 40 Mandal, S. *et al.* Analysis of composition of microbiomes: a novel method for studying
735 microbial composition. *Microb Ecol Health Dis* **26**, 27663, doi:10.3402/mehd.v26.27663
736 (2015).
- 737 41 Friedman, J. & Alm, E. J. Inferring correlation networks from genomic survey data. *PLoS*
738 *Comput Biol* **8**, e1002687, doi:10.1371/journal.pcbi.1002687 (2012).
- 739 42 Kuntal, B. K., Chandrakar, P., Sadhu, S. & Mande, S. S. 'NetShift': a methodology for
740 understanding 'driver microbes' from healthy and disease microbiome datasets. *ISME J*
741 **13**, 442-454, doi:10.1038/s41396-018-0291-x (2019).
- 742 43 Bannister, C. A., Halcox, J. P., Currie, C. J., Preece, A. & Spasic, I. A genetic
743 programming approach to development of clinical prediction models: A case study in
744 symptomatic cardiovascular disease. *PLoS One* **13**, e0202685,
745 doi:10.1371/journal.pone.0202685 (2018).
- 746 44 Schmidt, M. & Lipson, H. Distilling free-form natural laws from experimental data.
747 *Science* **324**, 81-85, doi:10.1126/science.1165893 (2009).
- 748 45 Gateau, C., Couturier, J., Coia, J. & Barbut, F. How to: diagnose infection caused by
749 *Clostridium difficile*. *Clin Microbiol Infect* **24**, 463-468, doi:10.1016/j.cmi.2017.12.005
750 (2018).
- 751 46 Loreau, M. *et al.* Biodiversity and ecosystem functioning: current knowledge and future
752 challenges. *Science* **294**, 804-808, doi:10.1126/science.1064088 (2001).
- 753 47 Ives, A. R. & Carpenter, S. R. Stability and diversity of ecosystems. *Science* **317**, 58-62,
754 doi:10.1126/science.1133258 (2007).
- 755 48 Ma, Z. S., Li, L. & Gotelli, N. J. Diversity-disease relationships and shared species
756 analyses for human microbiome-associated diseases. *ISME J* **13**, 1911-1919,
757 doi:10.1038/s41396-019-0395-y (2019).
- 758 49 Song, Y. *et al.* Microbiota dynamics in patients treated with fecal microbiota
759 transplantation for recurrent *Clostridium difficile* infection. *PLoS One* **8**, e81330,
760 doi:10.1371/journal.pone.0081330 (2013).

- 761 50 Milani, C. *et al.* Gut microbiota composition and *Clostridium difficile* infection in
762 hospitalized elderly individuals: a metagenomic study. *Sci Rep* **6**, 25945,
763 doi:10.1038/srep25945 (2016).
- 764 51 Jiang, Z. D. *et al.* Randomised clinical trial: faecal microbiota transplantation for
765 recurrent *Clostridium difficile* infection - fresh, or frozen, or lyophilised microbiota from
766 a small pool of healthy donors delivered by colonoscopy. *Aliment Pharmacol Ther* **45**,
767 899-908, doi:10.1111/apt.13969 (2017).
- 768 52 Shankar, V. *et al.* Species and genus level resolution analysis of gut microbiota in
769 *Clostridium difficile* patients following fecal microbiota transplantation. *Microbiome* **2**,
770 13, doi:10.1186/2049-2618-2-13 (2014).
- 771 53 Zaneveld, J. R., McMinds, R. & Vega Thurber, R. Stress and stability: applying the Anna
772 Karenina principle to animal microbiomes. *Nat Microbiol* **2**, 17121,
773 doi:10.1038/nmicrobiol.2017.121 (2017).
- 774 54 Giongo, A. *et al.* Toward defining the autoimmune microbiome for type 1 diabetes. *ISME*
775 *J* **5**, 82-91, doi:10.1038/ismej.2010.92 (2011).
- 776 55 Caussy, C. *et al.* A gut microbiome signature for cirrhosis due to nonalcoholic fatty liver
777 disease. *Nat Commun* **10**, 1406, doi:10.1038/s41467-019-09455-9 (2019).
- 778 56 Rowan, F. *et al.* *Desulfovibrio* bacterial species are increased in ulcerative colitis. *Dis*
779 *Colon Rectum* **53**, 1530-1536, doi:10.1007/DCR.0b013e3181f1e620 (2010).
- 780 57 Kolling, G. L. *et al.* Lactic acid production by *Streptococcus thermophilus* alters
781 *Clostridium difficile* infection and in vitro Toxin A production. *Gut Microbes* **3**, 523-529,
782 doi:10.4161/gmic.21757 (2012).
- 783 58 van de Beek, D. *et al.* Clinical features and prognostic factors in adults with bacterial
784 meningitis. *N Engl J Med* **351**, 1849-1859, doi:10.1056/NEJMoa040845 (2004).
- 785 59 Arnold, R. S. *et al.* Emergence of *Klebsiella pneumoniae* carbapenemase-producing
786 bacteria. *South Med J* **104**, 40-45, doi:10.1097/SMJ.0b013e3181fd7d5a (2011).
- 787 60 Navon-Venezia, S., Kondratyeva, K. & Carattoli, A. *Klebsiella pneumoniae*: a major
788 worldwide source and shuttle for antibiotic resistance. *FEMS Microbiol Rev* **41**, 252-275,
789 doi:10.1093/femsre/fux013 (2017).
- 790 61 Cohen, S. H. *et al.* Clinical practice guidelines for *Clostridium difficile* infection in
791 adults: 2010 update by the society for healthcare epidemiology of America (SHEA) and
792 the infectious diseases society of America (IDSA). *Infect Control Hosp Epidemiol* **31**,
793 431-455, doi:10.1086/651706 (2010).
- 794 62 Pollock, N. R. *et al.* Comparison of *Clostridioides difficile* Stool Toxin Concentrations in
795 Adults With Symptomatic Infection and Asymptomatic Carriage Using an Ultrasensitive
796 Quantitative Immunoassay. *Clin Infect Dis* **68**, 78-86, doi:10.1093/cid/ciy415 (2019).
- 797 63 Crobach, M. J. T. *et al.* Understanding *Clostridium difficile* Colonization. *Clin Microbiol*
798 *Rev* **31**, doi:10.1128/CMR.00021-17 (2018).
- 799 64 Antharam, V. C. *et al.* An Integrated Metabolomic and Microbiome Analysis Identified
800 Specific Gut Microbiota Associated with Fecal Cholesterol and Coprostanol in
801 *Clostridium difficile* Infection. *PLoS One* **11**, e0148824,
802 doi:10.1371/journal.pone.0148824 (2016).
- 803 65 Khanna, S. *et al.* Gut microbiome predictors of treatment response and recurrence in
804 primary *Clostridium difficile* infection. *Aliment Pharmacol Ther* **44**, 715-727,
805 doi:10.1111/apt.13750 (2016).

- 806 66 Han, S. H., Yi, J., Kim, J. H., Lee, S. & Moon, H. W. Composition of gut microbiota in
807 patients with toxigenic Clostridioides (Clostridium) difficile: Comparison between
808 subgroups according to clinical criteria and toxin gene load. *PLoS One* **14**, e0212626,
809 doi:10.1371/journal.pone.0212626 (2019).
- 810 67 Daquigan, N., Seekatz, A. M., Greathouse, K. L., Young, V. B. & White, J. R. High-
811 resolution profiling of the gut microbiome reveals the extent of Clostridium difficile
812 burden. *NPJ Biofilms Microbiomes* **3**, 35, doi:10.1038/s41522-017-0043-0 (2017).
- 813 68 Hudson, L. E., Anderson, S. E., Corbett, A. H. & Lamb, T. J. Gleaning Insights from
814 Fecal Microbiota Transplantation and Probiotic Studies for the Rational Design of
815 Combination Microbial Therapies. *Clin Microbiol Rev* **30**, 191-231,
816 doi:10.1128/CMR.00049-16 (2017).
- 817 69 Fujitani, S., George, W. L., Morgan, M. A., Nichols, S. & Murthy, A. R. Implications for
818 vancomycin-resistant Enterococcus colonization associated with Clostridium difficile
819 infections. *Am J Infect Control* **39**, 188-193, doi:10.1016/j.ajic.2010.10.024 (2011).
- 820 70 Antharam, V. C. *et al.* Intestinal dysbiosis and depletion of butyrogenic bacteria in
821 Clostridium difficile infection and nosocomial diarrhea. *J Clin Microbiol* **51**, 2884-2892,
822 doi:10.1128/JCM.00845-13 (2013).
- 823 71 Sokol, H. *et al.* Specificities of the intestinal microbiota in patients with inflammatory
824 bowel disease and Clostridium difficile infection. *Gut Microbes* **9**, 55-60,
825 doi:10.1080/19490976.2017.1361092 (2018).
- 826 72 Mezzatesta, M. L., Gona, F. & Stefani, S. Enterobacter cloacae complex: clinical impact
827 and emerging antibiotic resistance. *Future Microbiol* **7**, 887-902, doi:10.2217/fmb.12.61
828 (2012).
- 829 73 Umana, A. *et al.* Utilizing Whole Fusobacterium Genomes To Identify, Correct, and
830 Characterize Potential Virulence Protein Families. *J Bacteriol* **201**,
831 doi:10.1128/JB.00273-19 (2019).
- 832 74 Deng, H. Interpreting tree ensembles with inTrees. *International Journal of Data Science
833 and Analytics* **7**, 277-287, doi:10.1007/s41060-018-0144-8 (2019).
- 834 75 Dreiseitl, S. & Ohno-Machado, L. Logistic regression and artificial neural network
835 classification models: a methodology review. *J Biomed Inform* **35**, 352-359,
836 doi:10.1016/s1532-0464(03)00034-0 (2002).
- 837 76 Tollenaar, N. & van der Heijden, P. G. M. Optimizing predictive performance of criminal
838 recidivism models using registration data with binary and survival outcomes. *PLoS One*
839 **14**, e0213245, doi:10.1371/journal.pone.0213245 (2019).
- 840 77 Norton, E. C. & Dowd, B. E. Log Odds and the Interpretation of Logit Models. *Health
841 Serv Res* **53**, 859-878, doi:10.1111/1475-6773.12712 (2018).
- 842 78 Liu, K. H. & Xu, C. G. A genetic programming-based approach to the classification of
843 multiclass microarray datasets. *Bioinformatics* **25**, 331-337,
844 doi:10.1093/bioinformatics/btn644 (2009).
- 845 79 Fadrosch, D. W. *et al.* An improved dual-indexing approach for multiplexed 16S rRNA
846 gene sequencing on the Illumina MiSeq platform. *Microbiome* **2**, 6, doi:10.1186/2049-
847 2618-2-6 (2014).
- 848 80 Biesheuvel, C. J., Siccama, I., Grobbee, D. E. & Moons, K. G. Genetic programming
849 outperformed multivariable logistic regression in diagnosing pulmonary embolism. *J Clin
850 Epidemiol* **57**, 551-560, doi:10.1016/j.jclinepi.2003.10.011 (2004).

851 81 Staats, K., Pantridge, E., Cavaglia, M., Milovanov, I. & Aniyan, A. TensorFlow Enabled
852 Genetic Programming. *arXiv e-prints*, arXiv:1708.03157 (2017).

853
854

855 **Acknowledgements:** The authors thank all patients who participated in this study, as well as Carolyn
856 Alonso, Javier Villafuerte Gálvez, and the technologists in the Beth Israel Deaconess Medical Center
857 Clinical Microbiology Laboratory for their help with sample collection. The authors thank Zheng Sun for
858 valuable discussion on the microbiome data analysis.

859

860 **Funding:** Y.-Y.L. acknowledged grants from National Institutes of Health (R01AI141529,
861 R01HD093761, UH3OD023268, U19AI095219 and U01HL089856). N.R.P. and C.P.K. acknowledged
862 grants from National Institutes of Health (R01AI116596) and Institut Mérieux. S.K. was supported by the
863 China Scholarship Council.

864

865 **Author contributions:** Y.-Y.L., N.R.P., X.C., and C.P.K. conceived and designed the project. C.P.K.,
866 N.R.P., X.C. and K.D. performed the clinical study. X.C., H.X., and Q.L. contributed to the serum
867 immune marker measurement. K.W.G. and A.J.G. performed fecal DNA extraction and bacterial 16S
868 rRNA sequencing. S.K., X.-W.W., and Y.-Y.L. performed all the data analysis and wrote the manuscript.
869 N.R.P., K.W.G., C.P.K., and K.D. edited the manuscript.

870

871 **Competing interests:** C.P.K. has acted as a paid consultant to Artugen, Facile Therapeutics, First Light
872 Biosciences, Finch, Matrifax, Merck, Seres Health, and Vedanta and has received grant support from
873 Merck. X. C. has acted as a paid consultant to Artugen. All other authors report no potential conflicts of
874 interest.

875

876 **Data and materials availability:** Data will be available from corresponding authors upon reasonable
877 request.

878

879

880

881

882

883

884

885

886

887

888

889

890

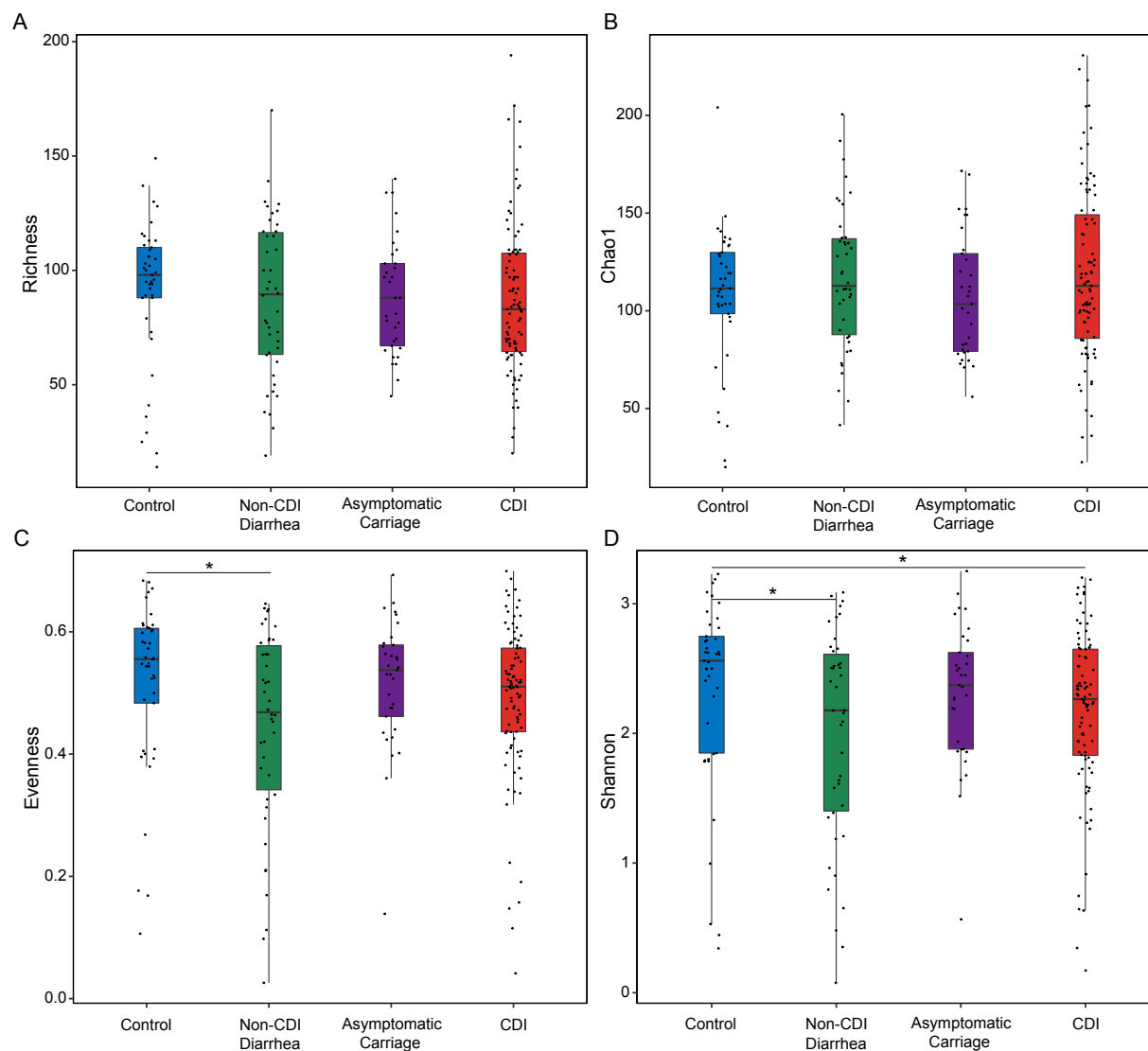
891

892

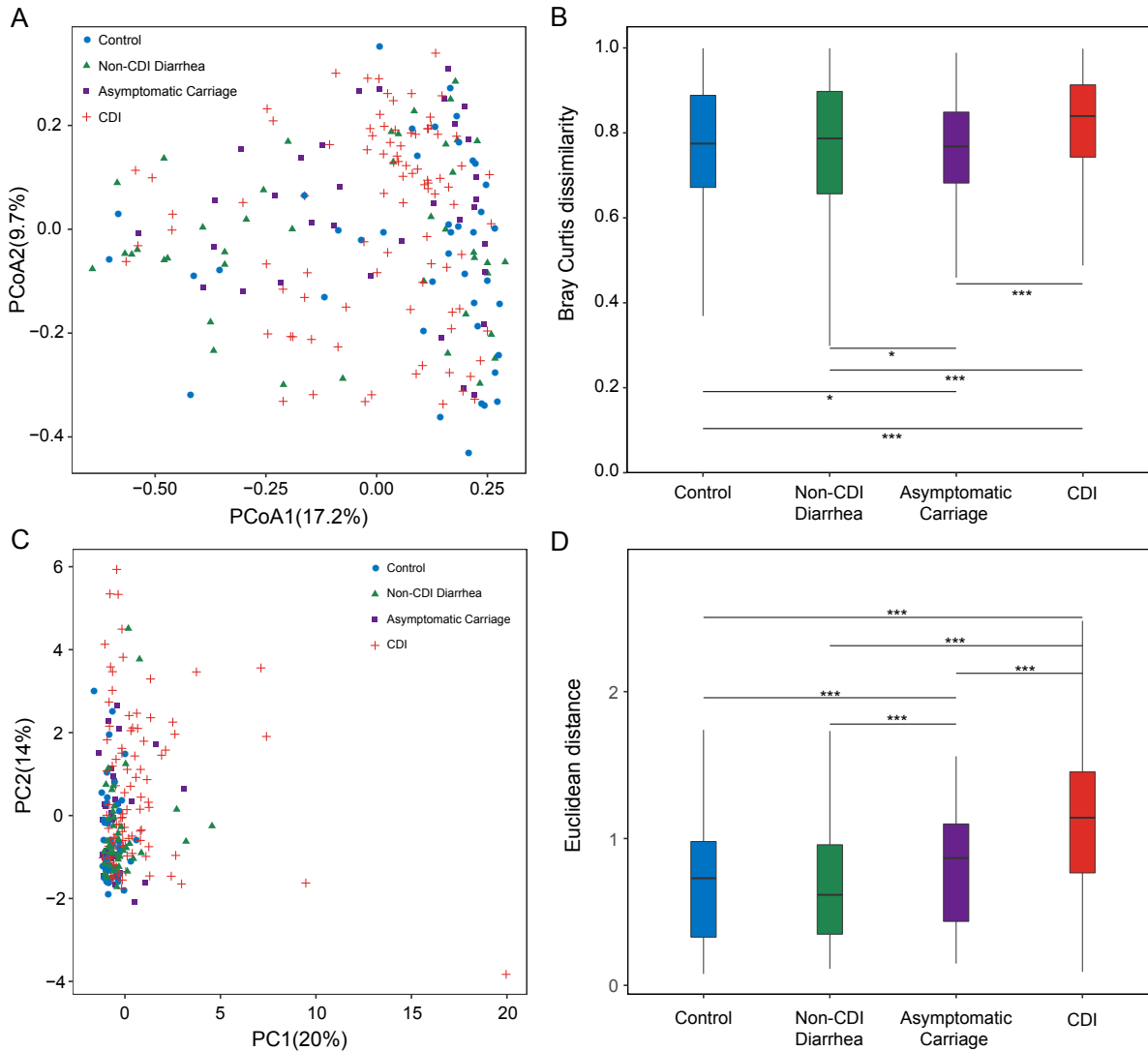
893

894

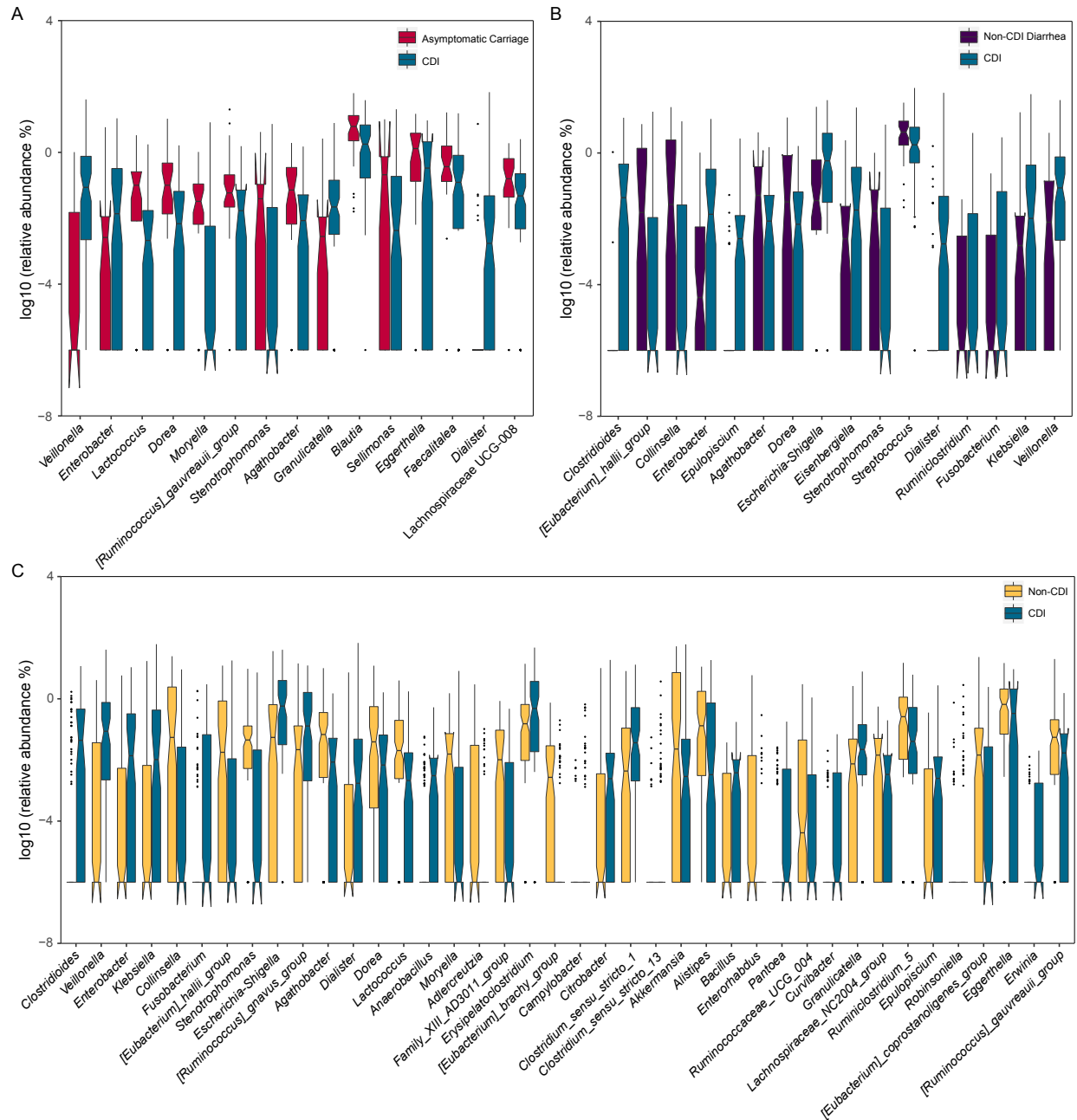
895 **Figures**



896
897 **Fig. 1. Comparing the alpha diversity of the gut microbiota of subjects with different *C.***
898 ***difficile* infection/colonization statuses (Control, Non-CDI Diarrhea, Asymptomatic**
899 **Carriage, and CDI) using different alpha diversity measures. (A) Taxa richness. (B) Chao1.**
900 **(C) Evenness. (D) Shannon index. Each dot represents the alpha diversity value of a particular**
901 **subject's gut microbiota. Statistical significance was determined by Mann-Whitney test,**
902 *** $P < 0.05$.**
903

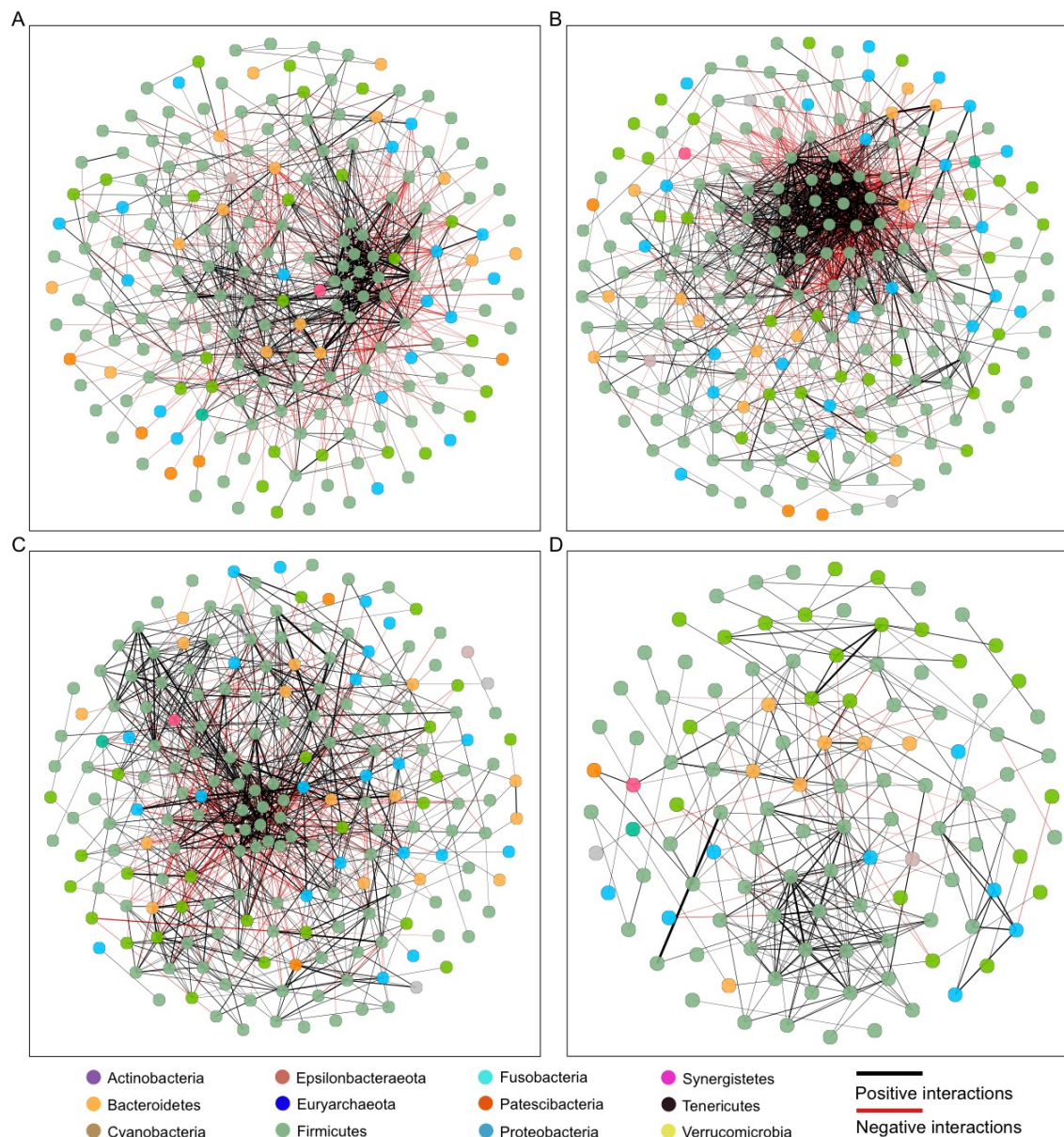


904
905 **Fig. 2. Ordination analysis and beta diversity comparison of the gut microbiota (and host**
906 **immune markers) for subjects with different *C. difficile* infection/colonization statuses**
907 **(Control, Non-CDI Diarrhea, Asymptomatic Carriage, and CDI). (A) Principal Coordinates**
908 **Analysis (PCoA) plot based on Bray–Curtis dissimilarities of microbial compositions. (B)**
909 **Boxplot of the gut microbiome Bray–Curtis dissimilarity between subjects within each group.**
910 **(C) Principle component analysis (PCA) plot of host immune marker concentrations. (D)**
911 **Boxplot of the Euclidean distance for the host immune markers of subjects within each group.**
912 **Statistical significance was determined by Mann–Whitney test, * $P < 0.05$, ** $P < 0.01$,**
913 ***** $P < 0.001$.**
914

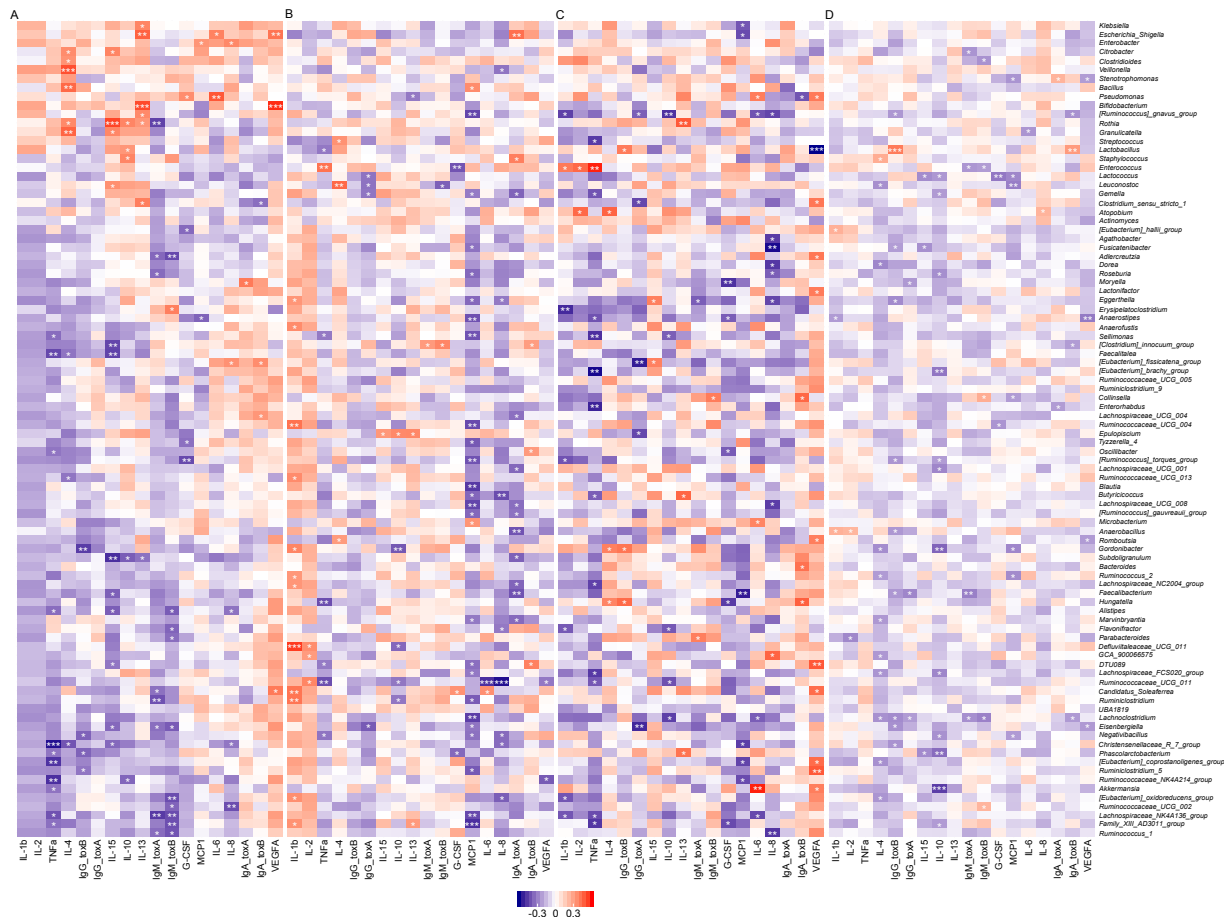


915
 916 **Fig. 3. Relative abundances of differentially abundant genera identified by ANCOM in**
 917 **comparing different groups. (A) CDI vs. Asymptomatic Carriage. (B) CDI vs. Non-CDI**
 918 **Diarrhea. (C) CDI vs. Non-CDI. The top differentially abundant taxa were ranked based on their**
 919 **W statistics (from left to right). The relative abundance (%) are plotted on log₁₀ scale. The**
 920 **notches in the boxplots show the 95% confidence interval around the median.**

921
 922

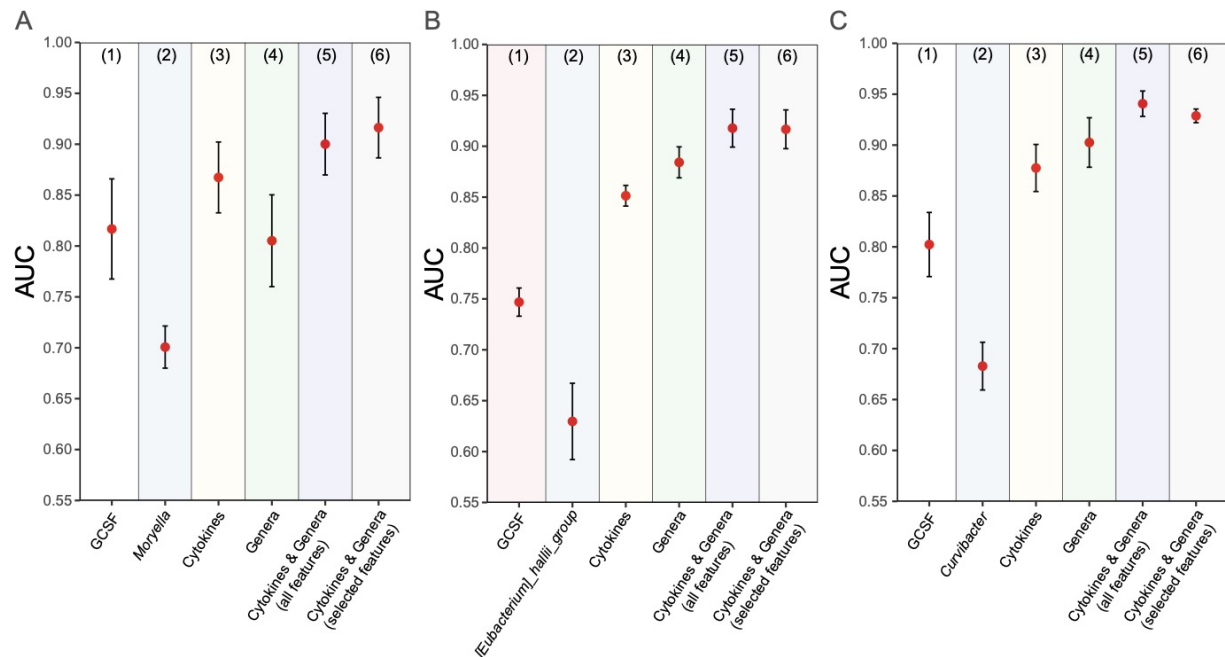


923
 924 **Fig. 4. Microbial correlation networks of different groups. (A) Control. (B) Non-CDI**
 925 **Diarrhea. (C) Asymptomatic Carriage. (D) CDI.** Nodes represent genera and are colored based
 926 on their phylum. Edges represent microbial correlations: green/red means positive/negative
 927 correlations, respectively. Edge thickness indicates correlation strength, and only the high-
 928 confidence interactions (p -value < 0.05) with high absolute correlation coefficients (> 0.3) were
 929 presented. For each group, we further identified the top-three most connected genera/nodes.
 930 They are *Ruminococcus_1*, *Roseburia* and *Lachnospiraceae_UCG-008* for the Control group,
 931 *[Ruminococcus]_torques_group*, *[Eubacterium]_hallii_group* and *Blautia* for the Non-CDI
 932 Diarrhea group, *Ruminiclostridium_5*, *Enterococcus* and *Lachnospiraceae_UCG_008* for the
 933 Asymptomatic Carriage group, and *Alistipes*, *Ruminiclostridium_5* and *Lachnoclostridium* for
 934 the CDI group.



935
 936 **Fig. 5. Correlations between gut microbial abundances and host immune markers in**
 937 **different groups, quantified by Spearman correlation with Benjamini-Hochberg correction.**
 938 (A) Control. (B) Non-CDI Diarrhea. (C) Asymptomatic Carriage. (D) CDI. Rows represent
 939 genera; columns represent immune markers. The layout of the heatmap is followed the
 940 hierarchical clustering results of Control cohort (see Fig.S4). Red/blue represents
 941 positive/negative correlation, respectively. The intensity of the colors denotes the strength of the
 942 correlation. * $\alpha < 0.05$, ** $\alpha < 0.01$, *** $\alpha < 0.001$.

943
 944



945
 946 **Fig. 6. The performance of RF-based classification models based on various types of**
 947 **features in differentiating CDI from other groups. (A) CDI vs. Asymptomatic Carriage. (B)**
 948 **CDI vs. Non-CDI Diarrhea. (C) CDI vs. Non-CDI. For each classification task, we used different**
 949 **types of features: (1) the top-1 immune marker feature (based on mean decrease accuracy); (2)**
 950 **the top-1 genus feature; (3) all immune markers; (4) all genera; (5) integration of all immune**
 951 **markers and genera; (6) selected features from the set of all immune markers and genera. Error**
 952 **bars represent the standard errors of the means (SEM).**

953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

972 **Table 1. Demographic characteristics of the enrolled subjects.**

973

Characteristics	NAAT negative		NAAT positive	
	Control (n=47)	Non-CDI Diarrhea (n=44)	Asymptomatic Carriage (n=40)	CDI (n=112)
Sex				
Female	14 (29.79%)	22 (50.00%)	20 (50.00%)	61 (54.46%)
Male	33 (70.21%)	22 (50.00%)	20 (50.00%)	51 (45.54%)
Age, Avg ± SD	62.40 ± 12.33	63.07 ± 13.15	62.15 ± 17.25	64.99 ± 15.62
Ethnicity				
Hispanic	1 (2.13%)	3 (6.82%)	1 (2.50%)	6 (5.36%)
Non-Hispanic	38 (80.85%)	37 (84.09%)	31 (77.50%)	96 (85.71%)
Unknown	8 (17.02%)	4 (9.09%)	8 (20.00%)	10 (8.93%)
Race				
White	33 (70.21%)	28 (63.64%)	28 (70.00%)	89 (79.46%)
Other	4 (8.51%)	10 (22.73%)	3 (7.50%)	23 (20.54%)
Unknown	10 (21.28%)	6 (13.64%)	9 (22.50%)	0 (0.00%)

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997 **Table 2. Diagnostic scores derived from symbolic classification (SC) and logistic regression**
 998 **(LR).** For each subject i , we calculate his/her diagnostic score $f(i)$ (or $p(i)$) based on one of the
 999 following formulas derived from SC (or LR), respectively. For SC, the class of subject i is CDI
 1000 if $f(i) > 0$; or Asymptomatic Carriage (or Non-CDI Diarrhea, Non-CDI) if $f(i) \leq 0$. For LR,
 1001 the class of subject i is CDI if $p(i) \geq 0.5$; or Asymptomatic Carriage (or Non-CDI Diarrhea,
 1002 Non-CDI) if $p(i) < 0.5$. Here, both $f(i)$ and $p(i)$ were learned from the entire dataset. Features
 1003 used here include: x_1 : GCSF; x_2 : IgA_toxA; x_3 : IgA_toxB; x_4 : IL6; x_5 : TNF α ; x_6 :
 1004 *Anaerobacillus*; x_7 : *Curvibacter*; x_8 : *Enterobacter*; x_9 : *Enterococcus*; x_{10} : *Epulopiscium*; x_{11} :
 1005 [*Eubacterium*] *haillii_group*; x_{12} : *Fusobacterium*; x_{13} : *Moryella*; x_{14} : *Stenotrophomonas*; x_{15} :
 1006 *Veillonella*. In particular, for each classification task (regardless of using SC or LR), the
 1007 following selected features were: (1) CDI vs. Asymptomatic Carriage: x_1, x_4, x_{13} and x_{15} ; (2)
 1008 CDI vs. Non-CDI Diarrhea: x_1, x_2, x_9, x_{10} , and x_{11} ; (3) CDI vs. Non-CDI: $x_1, x_3, x_4, x_5, x_6, x_7,$
 1009 x_8, x_{12}, x_{14} and x_{15} . Note that in the calculation of precision, recall and F1-score, we can treat
 1010 either CDI (or Asymptomatic Carriage, Non-CDI Diarrhea, Non-CDI) as the true positive.
 1011 Results shown in the parenthesis represent the latter case.
 1012

Model	Diagnostic	Formula	Accuracy	Precision	Recall	F1-score
SC	CDI vs. Asymptomatic Carriage	$f(i) = x_1 * x_{15}(x_1^3 - 0.2 * x_{13} + 0.4) + 1.1 * x_1 - 0.1 * x_4 - 18.25$	0.896	0.914 (0.840)	0.949 (0.75)	0.931 (0.792)
	CDI vs. Non-CDI Diarrhea	$f(i) = x_9 * x_2(0.5 * x_{10} - 1) + x_{11}(0.02 * x_{11} - x_1) + x_2 \left(1 - \frac{10}{x_1}\right) - \frac{0.003}{x_9}$	0.900	0.946 (0.826)	0.897 (0.905)	0.921 (0.864)
	CDI vs. Non-CDI	$f(i) = x_1 * x_3(0.2 * x_1 * x_5 * x_6 * x_{14} + 0.04 * x_1 * x_7 + 0.3 * x_1 * x_{15} * x_8^4 + x_1 * x_{12}(0.5 * x_7 + x_1 * x_{14}) + x_7(0.1 * x_4 - x_6) + x_{14}(x_{14} - 2)$	0.882	0.889 (0.878)	0.821 (0.927)	0.853 (0.902)
LR	CDI vs. Asymptomatic Carriage	$\log\left(\frac{p(i)}{1-p(i)}\right) = 0.66725 - 0.04442 * x_1 + 0.01022 * x_4 + 7.51484 * x_{13} - 85.00213 * x_{15}$	0.830	0.895 (0.667)	0.872 (0.714)	0.883 (0.690)
	CDI vs. Non-CDI Diarrhea	$\log\left(\frac{p(i)}{1-p(i)}\right) = 0.01974 - 0.002084 * x_1 - 0.02391 * x_2 + 1.895 * x_9 - 12740 * x_{10} + 163.9 * x_{11}$	0.800	0.814 (0.765)	0.897 (0.619)	0.854 (0.684)
	CDI vs. Non-CDI	$\log\left(\frac{p(i)}{1-p(i)}\right) = 2.122 - 0.01002 * x_1 + 0.01833 * x_3 - 0.006334 * x_4 - 0.009566 * x_5 - 4609 * x_6 - 8576 * x_7 - 40.75 * x_8 - 101.1 * x_{12} + 32.84 * x_{14} - 43.4 * x_{15}$	0.813	0.841 (0.798)	0.679 (0.908)	0.752 (0.850)

1013

1014

1015

1016 **Supplementary Materials**

1017

1018 Fig. S1. The “driver” taxa responsible for the change of microbial correlations between CDI and
1019 Asymptomatic Carriage

1020 Fig. S2. The “driver” taxa responsible for the change of microbial correlations between CDI and
1021 Non-CDI Diarrhea.

1022 Fig. S3. The “driver” taxa responsible for the change of microbial correlations between CDI and
1023 Non-CDI.

1024 Fig. S4. Significant correlations between gut microbial abundances and host immune markers in
1025 the Control group.

1026 Fig. S5. Gut microbiota and host immune markers can accurately differentiate different groups in
1027 multi-class classification models.

1028 Fig. S6. Using the mean decrease accuracy (MDA) ranking and the 1-SE rule to select features to
1029 distinguish CDI from other groups.

1030 Fig. S7. The fitness evolution during the genetic programming.

1031 Table S1. Sample sizes of different data types in different groups.

1032 Table S2. Permutational multivariate analysis of variance (PERMANOVA) in microbial
1033 compositions and immune markers.

1034 Table S3. Differentially abundant genera between CDI and Asymptomatic Carriage groups
1035 detected by ANCOM, adjusted for age and sex.

1036 Table S4. Differentially abundant genera between CDI and Non-CDI Diarrhea groups detected
1037 by ANCOM, adjusted for age and sex.

1038 Table S5. Differentially abundant genera between CDI and Non-CDI groups detected by
1039 ANCOM, adjusted for age and sex.

1040 Table S6. Characteristics of microbial correlation networks associated with different groups.

1041 Table S7. Comparison of host immune markers in different groups.

1042 Table S8. Accuracy, Precision, Recall and F1-score of symbolic classification in CDI diagnosis.

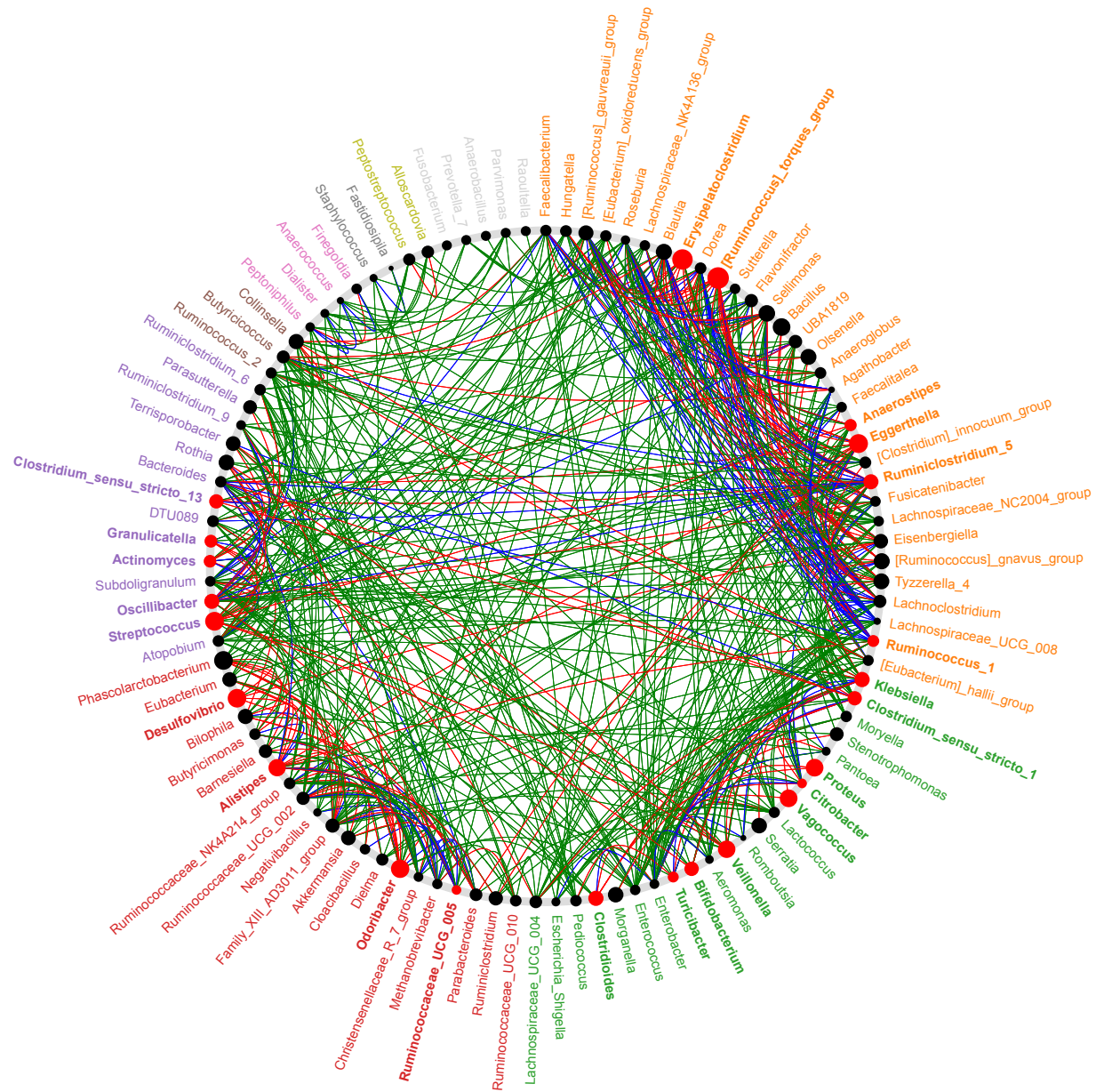
1043

1044

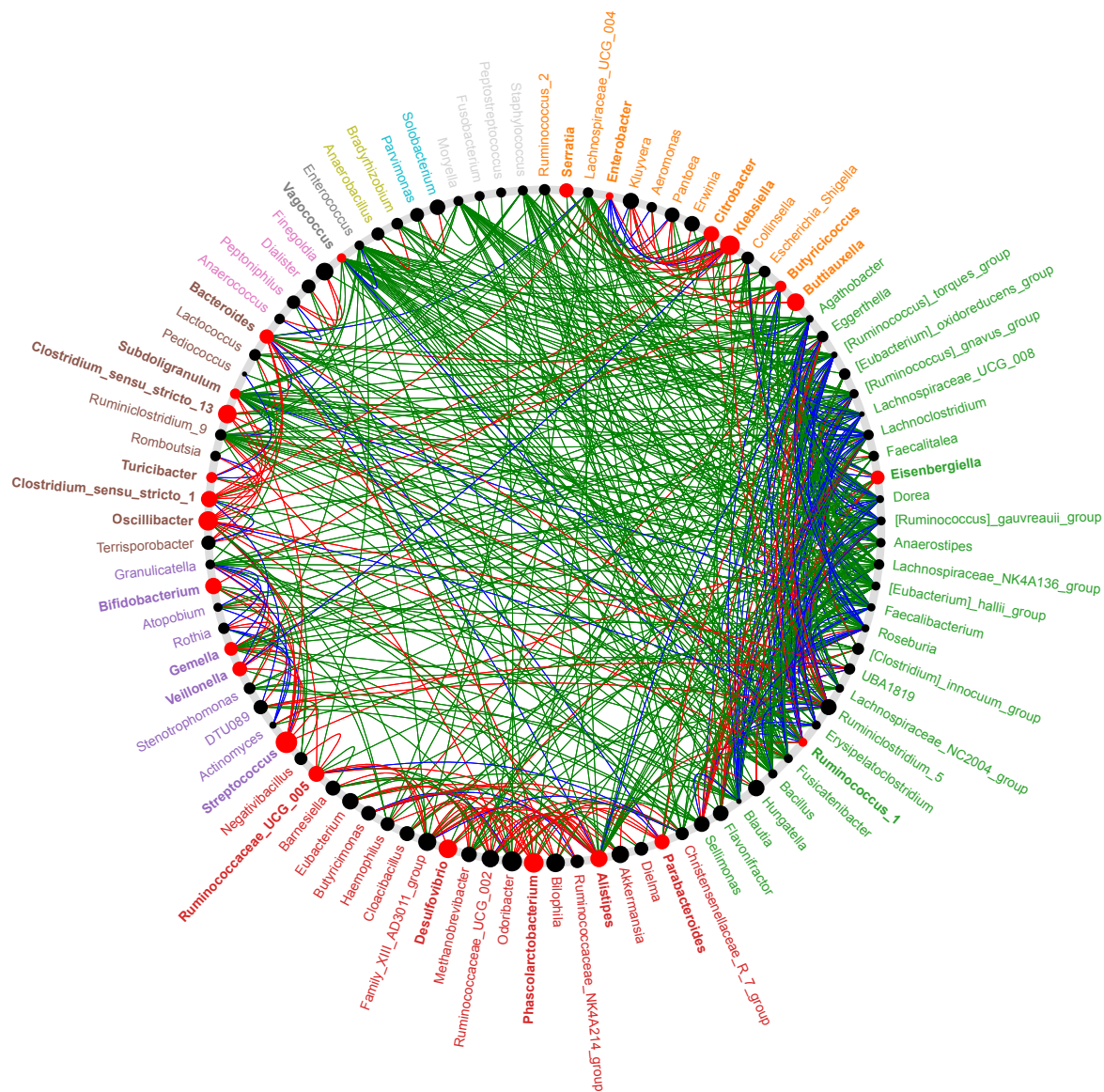
1045

1046

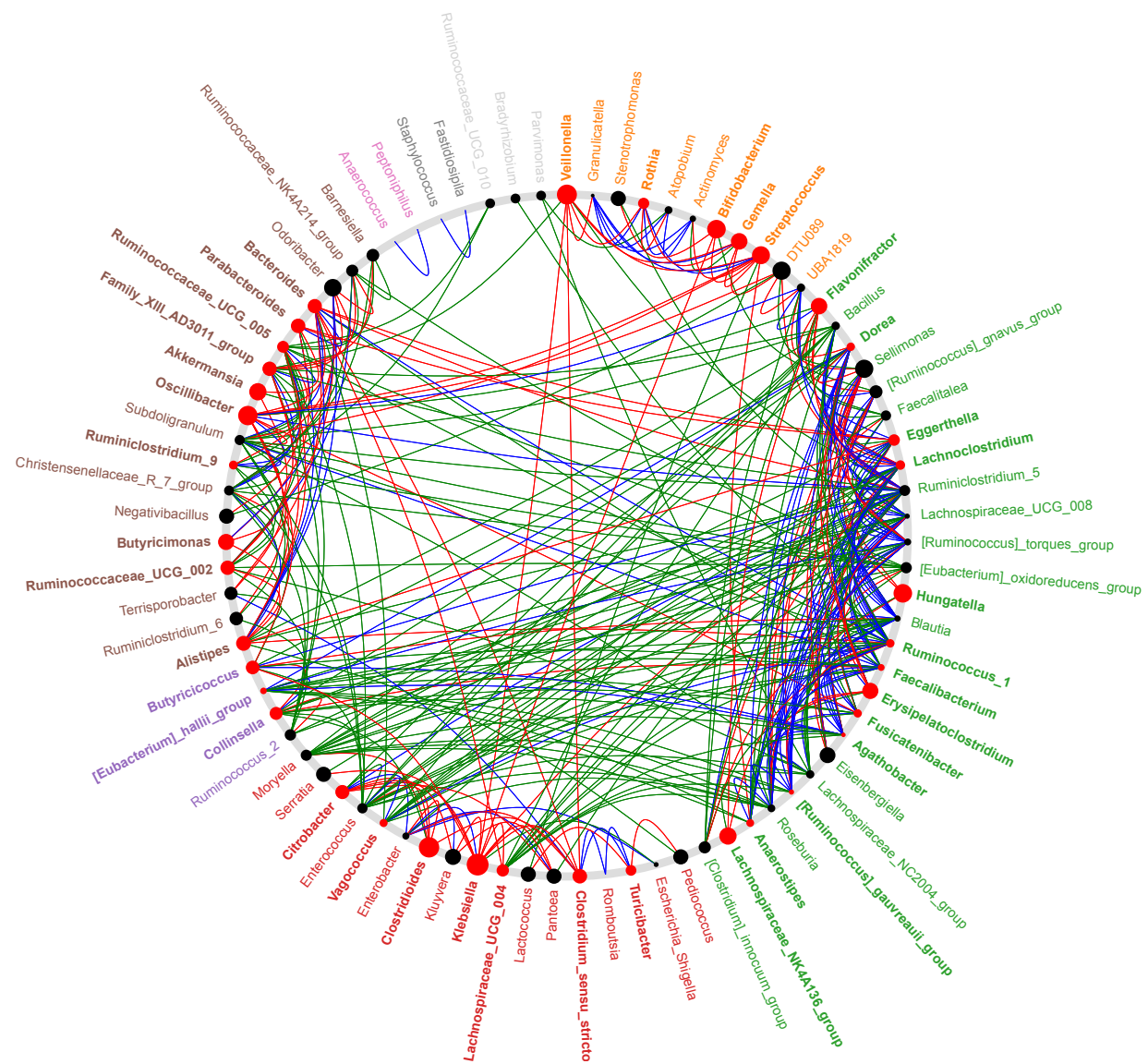
1047



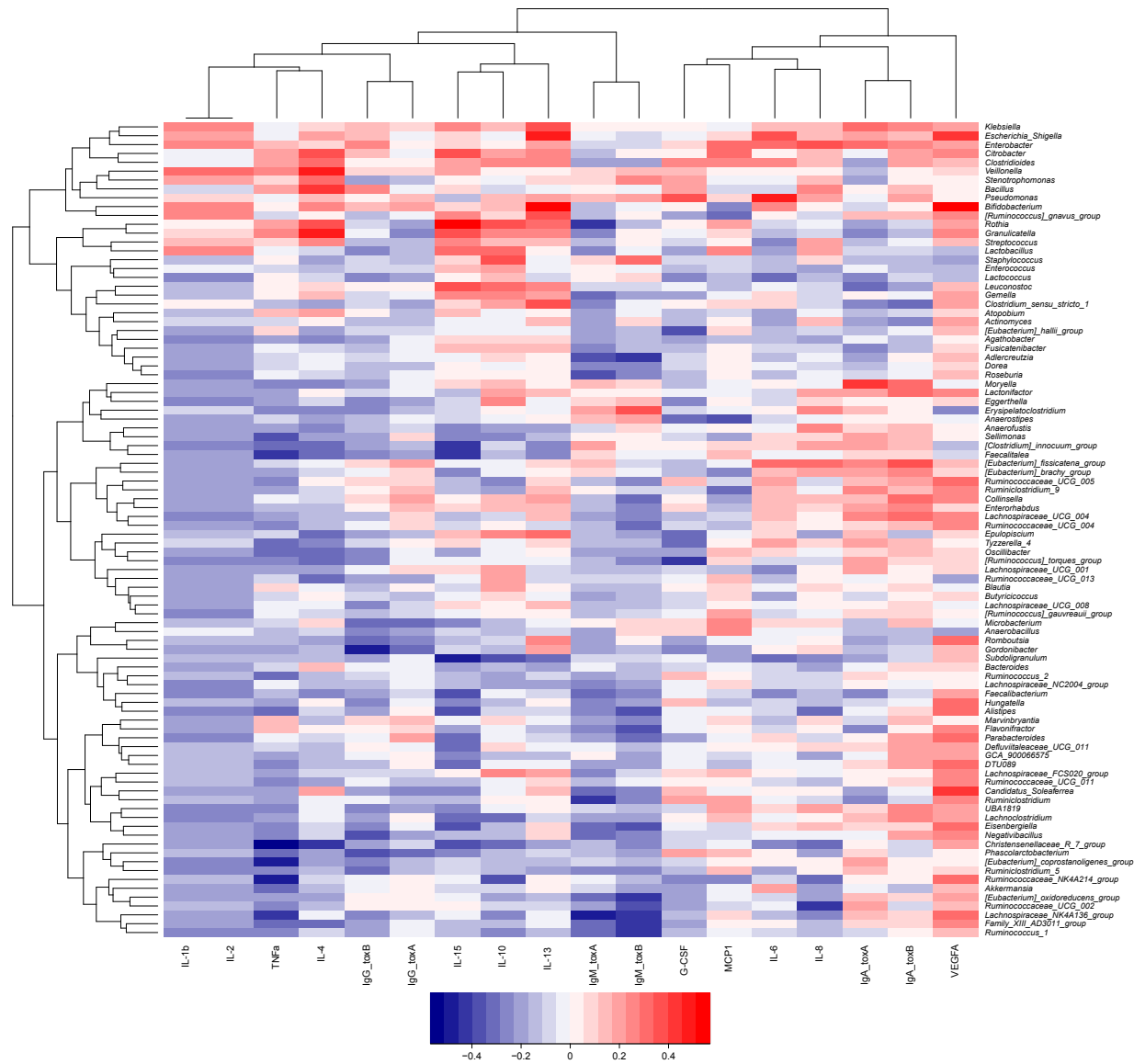
1048
1049 **Fig. S1. The “driver” taxa responsible for the change of microbial correlations between**
1050 **CDI and Asymptomatic Carriage.** Node sizes are proportional to their scaled neighbor shift
1051 (NESH) score (i.e., a score identifying important microbial taxa of microbial association
1052 networks) and a node is colored red if its betweenness increases when comparing microbial
1053 correlation networks of CDI with that of Asymptomatic Carriage. All taxa belonging to same
1054 community (common sub-network) are randomly assigned a color to their labels. Red (or green)
1055 edges represent microbial correlations that are only present in the CDI (or Asymptomatic
1056 Carriage) network, respectively. Blue edges present common microbial correlations that are
1057 present in both networks.



1058
1059 **Fig. S2. The “driver” taxa responsible for the change of microbial correlations between**
1060 **CDI and Non-CDI Diarrhea.** Node sizes are proportional to their scaled NESH score and a
1061 node is colored red if its betweenness increases when comparing microbial correlation networks
1062 of CDI with that of Non-CDI Diarrhea. All taxa belonging to same community (common sub-
1063 network) are randomly assigned a color to their labels. Red (or green) edges represent microbial
1064 correlations that are only present in the CDI (or Non-CDI Diarrhea) network, respectively. Blue
1065 edges present common microbial correlations that are present in both networks.
1066
1067

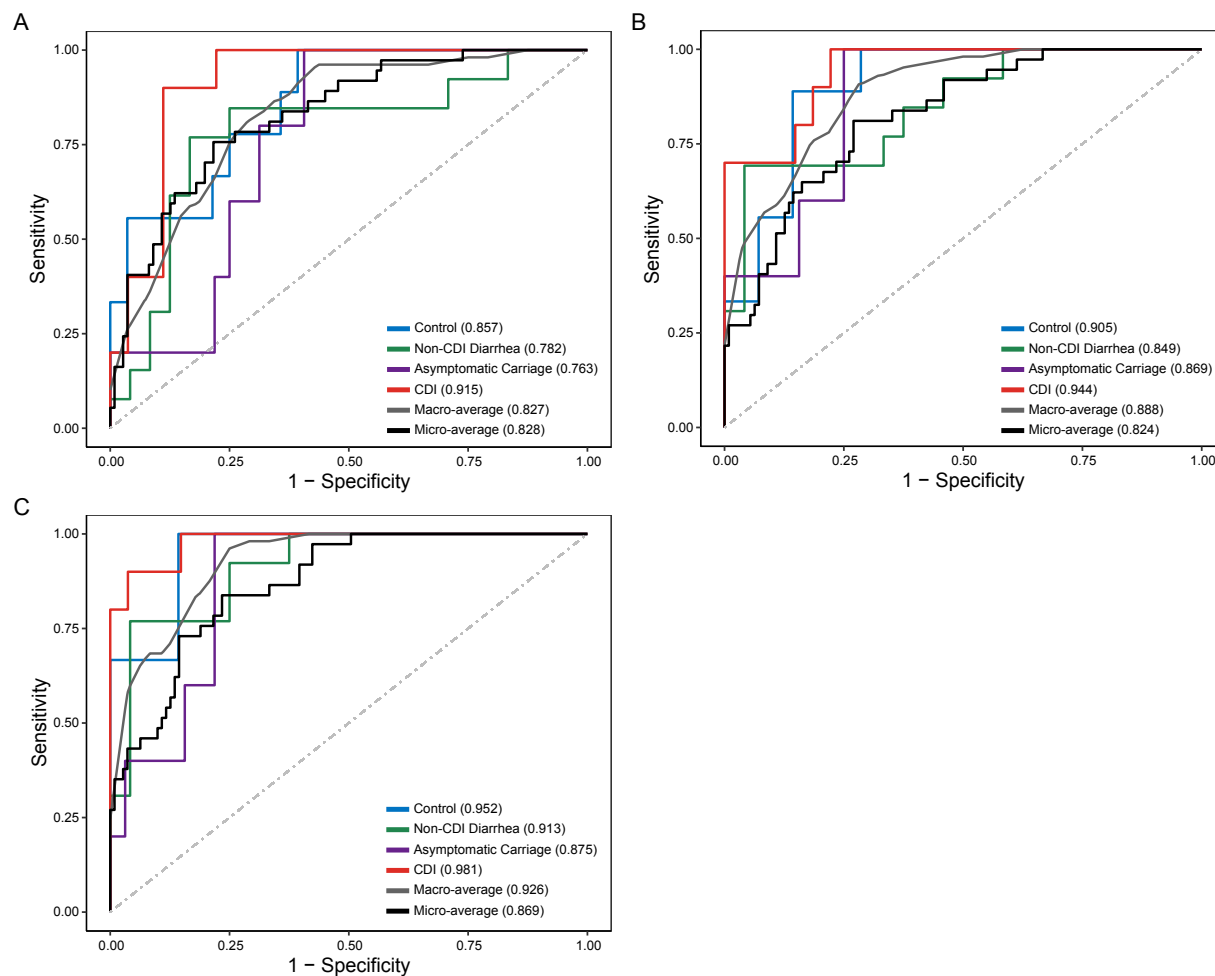


1068
1069 **Fig. S3. The potential “driver taxa” responsible for the change of microbial correlations**
1070 **between CDI and Non-CDI.** Node sizes are proportional to their scaled NESH score and a node
1071 is colored red if its betweenness increases when comparing microbial correlation networks of
1072 CDI with that of Non-CDI. All taxa belonging to same community (common sub-network) are
1073 randomly assigned a color to their labels. Red (or green) edges represent microbial correlations
1074 that are only present in the CDI (or Non-CDI) network, respectively. Blue edges present
1075 common microbial correlations that are present in both networks.
1076
1077
1078



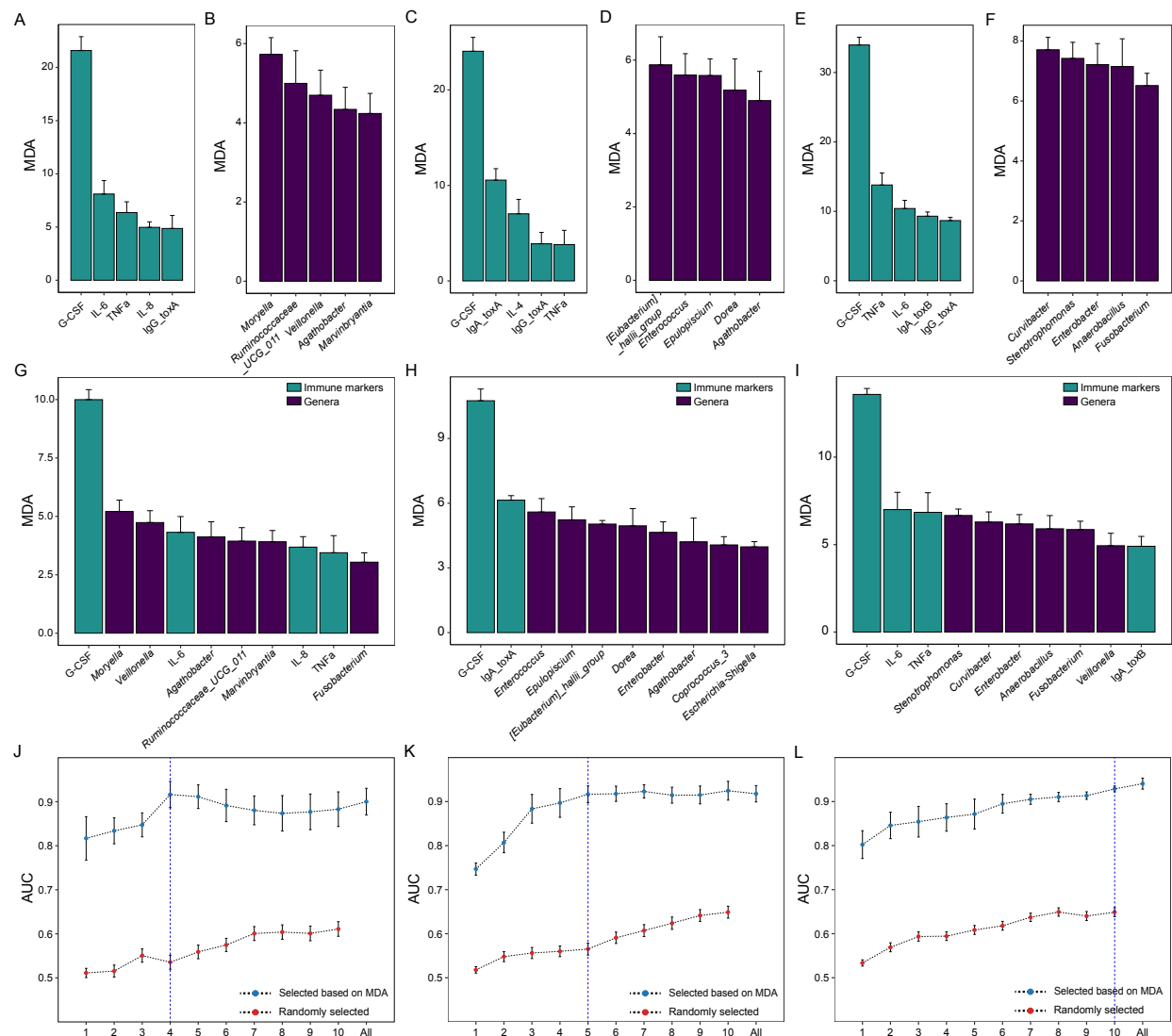
1079
 1080 **Fig. S4. Significant correlations between gut microbial abundances and host immune**
 1081 **markers in the Control group.** Gut microbial compositions and host immune markers were
 1082 clustered through hierarchical clustering. Rows correspond to bacterial taxa at genus level;
 1083 columns correspond to host immune markers. Red/blue represents positive/negative association,
 1084 respectively. The intensity of the colors denotes the strength of correlation between the genus
 1085 abundance and the immunological expression level.

1086
 1087
 1088
 1089
 1090
 1091
 1092

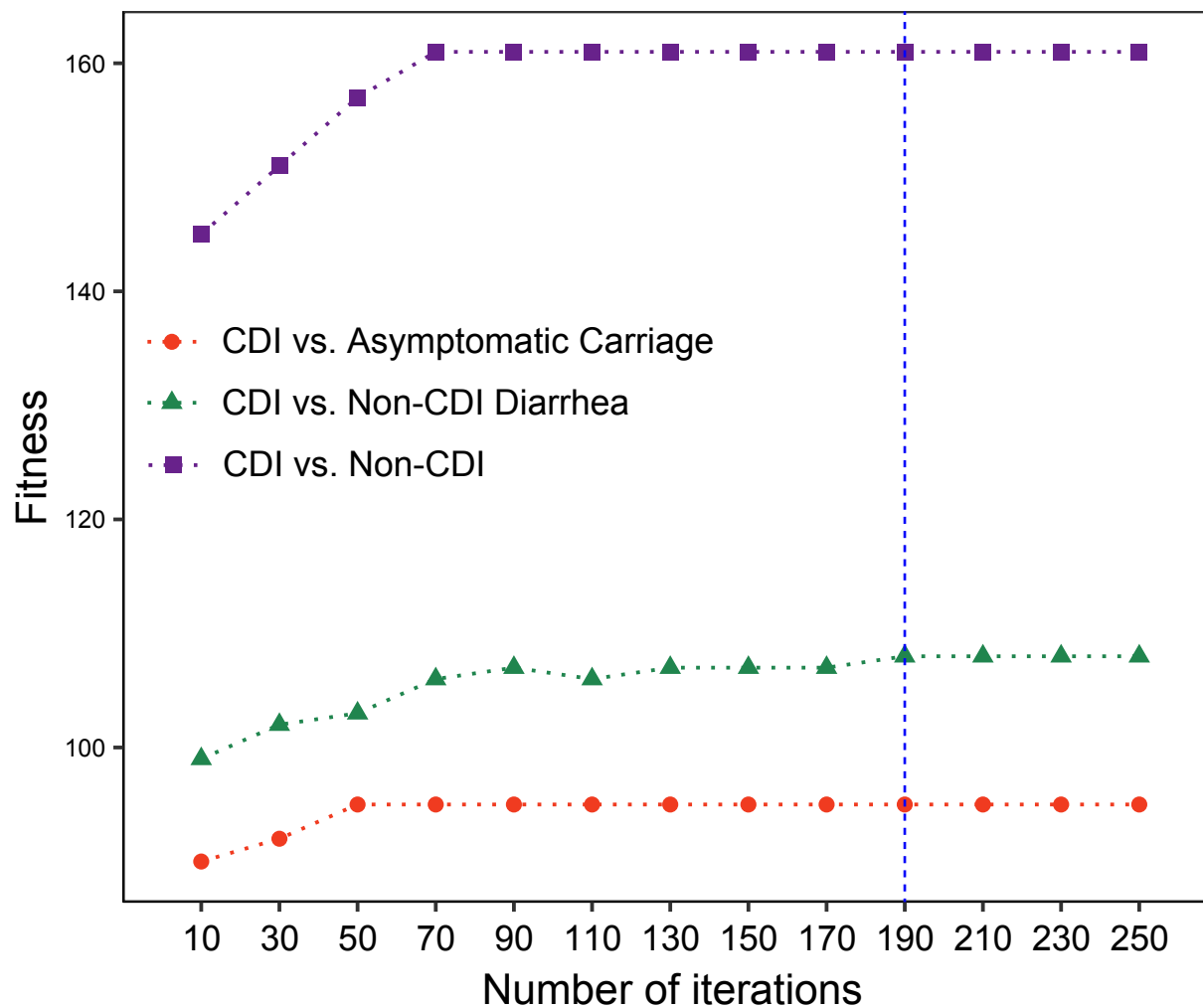


1093
1094 **Fig. S5. Gut microbiota and host immune markers can accurately differentiate different**
1095 **groups in multi-class classification models. (A) Use host immune markers alone. (B) Use gut**
1096 **microbiota data (at genus level) alone. (C) The integration of host immune markers and**
1097 **microbial data. The performance of each classifier is measured by the macro-average and micro-**
1098 **average AUCs.**

1099
1100
1101
1102
1103
1104
1105
1106



1107
 1108
 1109 **Fig. S6. Using the mean decrease accuracy (MDA) ranking and the 1-SE rule to select**
 1110 **features to distinguish CDI from other groups.** The most important features of cytokine data,
 1111 microbiome data, and the integration of cytokines and microbiome data in classifying CDI vs.
 1112 Asymptomatic Carriage (**A, B and G**), CDI vs. Non-CDI Diarrhea (**C, D and H**) and CDI vs.
 1113 Non-CDI (**E, F and I**). The performance of classifiers using different sets of integrated features:
 1114 selected based on MDA or randomly selected in CDI vs Asymptomatic Carriage (**J**), CDI vs
 1115 Non-CDI Diarrhea (**K**) and CDI vs Non-CDI (**L**). The minimum set of features selected based on
 1116 the MDA ranking and the 1-SE rule is highlighted by a vertical blue dashed line. Error bars
 1117 represent the standard errors of the means (SEM).
 1118
 1119
 1120
 1121



1122
1123 **Fig. S7. The fitness evolution during the symbolic classification based on genetic**
1124 **programming.** The fitness function is a maximization function, and the tree with highest fitness
1125 score in each iteration were plotted. The final selected number of generations is highlighted with
1126 a vertical blue dashed line.

1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138

1139 **Table S1. Sample sizes of different data types in different groups.**

1140

Characteristics	NAAT negative		NAAT positive		Total
	Control	Non-CDI Diarrhea	Asymptomatic Carriage	CDI	
Immunological data	45	44	35	99	223
Microbial data	41	42	33	91	207
Immunological & microbial data	39	42	28	78	187

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173 **Table S2. Permutational multivariate analysis of variance (PERMANOVA) in microbial**
1174 **compositions and immune markers.** CDI statuses: Control, Non-CDI Diarrhea, Asymptomatic
1175 Carriage, and CDI. Race: White, Native American, Asian, African American, Pacific Islander
1176 and mixed origin. Ethnicity: Hispanic and Not Hispanic. Here F represents the F-statistic: a
1177 larger F value indicate that the between-group variation is greater than within-group variation. R²
1178 represents the variation explained by the model. P represents the *P*-value calculated from
1179 permutation.

Test factors	Microbiome			Cytokines		
	F	R ²	P	F	R ²	P
CDI status	2.285	0.0388	0.001	3.351	0.052	0.016
Age	1.605	0.009	0.081	0.541	0.003	0.516
Sex	1.557	0.009	0.095	0.916	0.005	0.372
Race	0.881	0.031	0.832	1.595	0.050	0.153
Ethnicity	0.476	0.003	0.961	0.206	0.001	0.771

1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206

1207 **Table S3. Differentially abundant genera between CDI and Asymptomatic Carriage groups**
 1208 **detected by ANCOM, adjusted for age and sex.** For each genus, the first column represents its
 1209 W statistic, and subsequent four columns represent logical indicators of whether it is
 1210 differentially abundant under a series of cutoffs (0.9, 0.8, 0.7 and 0.6). The last two columns
 1211 represent its relative abundance (mean \pm standard deviation) in the two groups.
 1212

Genera	W_stat	Cutoff 0.9	Cutoff 0.8	Cutoff 0.7	Cutoff 0.6	Relative abundance (%) in CDI	Relative abundance (%) in Asymptomatic Carriage
<i>Veillonella</i>	204	TRUE	TRUE	TRUE	TRUE	1.86 \pm 5.55	0.06 \pm 0.20
<i>Enterobacter</i>	182	FALSE	TRUE	TRUE	TRUE	0.79 \pm 1.81	0.20 \pm 1.01
<i>Lactococcus</i>	179	FALSE	TRUE	TRUE	TRUE	0.10 \pm 0.30	0.36 \pm 0.69
<i>Dorea</i>	177	FALSE	TRUE	TRUE	TRUE	0.10 \pm 0.26	0.79 \pm 2.00
<i>Moryella</i>	174	FALSE	TRUE	TRUE	TRUE	0.20 \pm 1.19	0.13 \pm 0.24
<i>[Ruminococcus]_gauvreauii_group</i>	173	FALSE	TRUE	TRUE	TRUE	0.09 \pm 0.26	1.10 \pm 3.68
<i>Stenotrophomonas</i>	167	FALSE	TRUE	TRUE	TRUE	0.13 \pm 0.79	0.28 \pm 0.76
<i>Agathobacter</i>	158	FALSE	FALSE	TRUE	TRUE	0.07 \pm 0.19	0.25 \pm 0.42
<i>Granulicatella</i>	157	FALSE	FALSE	TRUE	TRUE	0.31 \pm 1.10	0.10 \pm 0.46
<i>Blautia</i>	154	FALSE	FALSE	TRUE	TRUE	5.30 \pm 7.99	10.18 \pm 14.05
<i>Sellimonas</i>	150	FALSE	FALSE	TRUE	TRUE	0.46 \pm 2.18	1.20 \pm 2.50
<i>Eggerthella</i>	147	FALSE	FALSE	TRUE	TRUE	1.39 \pm 2.16	2.98 \pm 4.06
<i>Faecalitalea</i>	145	FALSE	FALSE	TRUE	TRUE	0.91 \pm 2.07	1.41 \pm 3.05
<i>Dialister</i>	141	FALSE	FALSE	FALSE	TRUE	1.10 \pm 7.07	0.23 \pm 1.27
<i>Lachnospiraceae_UCG_008</i>	135	FALSE	FALSE	FALSE	TRUE	0.20 \pm 0.39	0.37 \pm 0.45

1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229

1230 **Table S4. Differentially abundant genera between CDI and Non-CDI Diarrhea groups**
 1231 **detected by ANCOM, adjusted for age and sex.** For each genus, the first column represents its
 1232 W statistic, and subsequent four columns represent logical indicators of whether it is
 1233 differentially abundant under a series of cutoffs (0.9, 0.8, 0.7 and 0.6). The last two columns
 1234 represent its relative abundance (mean \pm standard deviation) in the two groups.
 1235

Genera	W_stat	detected_0.9	detected_0.8	detected_0.7	detected_0.6	Relative abundance (%) in CDI	Relative abundance (%) in Non-CDI Diarrhea
<i>Clostridioides</i>	206	TRUE	TRUE	TRUE	TRUE	0.67 \pm 1.80	0.03 \pm 0.16
<i>[Eubacterium]_hallii_group</i>	199	TRUE	TRUE	TRUE	TRUE	0.40 \pm 2.12	1.14 \pm 2.05
<i>Collinsella</i>	195	TRUE	TRUE	TRUE	TRUE	0.57 \pm 1.59	2.57 \pm 5.33
<i>Enterobacter</i>	189	TRUE	TRUE	TRUE	TRUE	0.79 \pm 1.81	0.05 \pm 0.20
<i>Epulopiscium</i>	166	FALSE	TRUE	TRUE	TRUE	0.06 \pm 0.30	0.00 \pm 0.01
<i>Agathobacter</i>	165	FALSE	TRUE	TRUE	TRUE	0.07 \pm 0.19	0.42 \pm 0.93
<i>Dorea</i>	165	FALSE	TRUE	TRUE	TRUE	0.10 \pm 0.26	0.96 \pm 2.24
<i>Escherichia_Shigella</i>	163	FALSE	FALSE	TRUE	TRUE	3.54 \pm 6.46	1.84 \pm 5.01
<i>Eisenbergiella</i>	149	FALSE	FALSE	TRUE	TRUE	1.03 \pm 3.36	0.10 \pm 0.40
<i>Stenotrophomonas</i>	147	FALSE	FALSE	TRUE	TRUE	0.13 \pm 0.79	0.09 \pm 0.17
<i>Streptococcus</i>	147	FALSE	FALSE	TRUE	TRUE	6.16 \pm 13.27	7.00 \pm 7.39
<i>Dialister</i>	138	FALSE	FALSE	FALSE	TRUE	1.10 \pm 7.07	0.06 \pm 0.25
<i>Ruminiclostridium</i>	137	FALSE	FALSE	FALSE	TRUE	0.08 \pm 0.44	0.00 \pm 0.01
<i>Fusobacterium</i>	131	FALSE	FALSE	FALSE	TRUE	0.18 \pm 0.51	0.01 \pm 0.04
<i>Klebsiella</i>	131	FALSE	FALSE	FALSE	TRUE	1.75 \pm 6.94	0.58 \pm 2.74
<i>Veillonella</i>	125	FALSE	FALSE	FALSE	TRUE	1.86 \pm 5.55	0.27 \pm 0.78

1236
 1237
 1238
 1239
 1240
 1241
 1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251

1252 **Table S5. Differentially abundant genera between CDI and Non-CDI groups detected by**
 1253 **ANCOM, adjusted for age and sex.** For each genus, the first column represents its *W* statistic,
 1254 and subsequent four columns represent logical indicators of whether it is differentially abundant
 1255 under a series of cutoffs (0.9, 0.8, 0.7 and 0.6). The last two columns represent its relative
 1256 abundance (mean \pm standard deviation) in the two groups.

Genera	W_stat	detected _0.9	detected_ 0.8	detected_ 0.7	detected_ 0.6	Relative abundance (%) in CDI	Relative abundance (%) in Non-CDI
<i>Clostridioides</i>	201	TRUE	TRUE	TRUE	TRUE	0.67 \pm 1.81	0.08 \pm 0.26
<i>Veillonella</i>	201	TRUE	TRUE	TRUE	TRUE	1.86 \pm 5.58	0.14 \pm 0.50
<i>Enterobacter</i>	200	TRUE	TRUE	TRUE	TRUE	0.79 \pm 1.82	0.08 \pm 0.55
<i>Klebsiella</i>	196	TRUE	TRUE	TRUE	TRUE	1.75 \pm 6.98	0.49 \pm 2.25
<i>Collinsella</i>	194	TRUE	TRUE	TRUE	TRUE	0.57 \pm 1.60	2.29 \pm 4.64
<i>Fusobacterium</i>	194	TRUE	TRUE	TRUE	TRUE	0.18 \pm 0.51	0.04 \pm 0.25
<i>[Eubacterium]_hallii_group</i>	193	TRUE	TRUE	TRUE	TRUE	0.40 \pm 2.13	1.21 \pm 2.55
<i>Stenotrophomonas</i>	193	TRUE	TRUE	TRUE	TRUE	0.13 \pm 0.79	0.30 \pm 1.06
<i>Escherichia_Shigella</i>	191	TRUE	TRUE	TRUE	TRUE	3.54 \pm 6.50	1.96 \pm 5.31
<i>[Ruminococcus]_gnavus_group</i>	185	FALSE	TRUE	TRUE	TRUE	1.73 \pm 3.15	0.50 \pm 1.81
<i>Agathobacter</i>	179	FALSE	TRUE	TRUE	TRUE	0.07 \pm 0.19	0.50 \pm 1.37
<i>Dialister</i>	176	FALSE	TRUE	TRUE	TRUE	1.10 \pm 7.11	0.10 \pm 0.69
<i>Dorea</i>	170	FALSE	TRUE	TRUE	TRUE	0.10 \pm 0.26	0.86 \pm 2.17
<i>Lactococcus</i>	169	FALSE	TRUE	TRUE	TRUE	0.10 \pm 0.30	0.23 \pm 0.56
<i>Anaerobacillus</i>	164	FALSE	FALSE	TRUE	TRUE	0.02 \pm 0.06	0.00 \pm 0.01
<i>Moryella</i>	164	FALSE	FALSE	TRUE	TRUE	0.20 \pm 1.19	0.11 \pm 0.26
<i>Adlercreutzia</i>	162	FALSE	FALSE	TRUE	TRUE	0.00 \pm 0.02	0.08 \pm 0.33
<i>Family_XIII_AD3011_group</i>	161	FALSE	FALSE	TRUE	TRUE	0.03 \pm 0.07	0.09 \pm 0.16
<i>Erysipelatoclostridium</i>	160	FALSE	FALSE	TRUE	TRUE	3.56 \pm 7.56	0.80 \pm 1.79
<i>[Eubacterium]_brachy_group</i>	154	FALSE	FALSE	TRUE	TRUE	0.01 \pm 0.03	0.04 \pm 0.10
<i>Campylobacter</i>	154	FALSE	FALSE	TRUE	TRUE	0.03 \pm 0.11	0.00 \pm 0.00
<i>Citrobacter</i>	154	FALSE	FALSE	TRUE	TRUE	0.41 \pm 2.11	0.20 \pm 1.09
<i>Clostridium_sensu_stricto_1</i>	151	FALSE	FALSE	TRUE	TRUE	0.72 \pm 1.76	0.38 \pm 1.19
<i>Clostridium_sensu_stricto_13</i>	151	FALSE	FALSE	TRUE	TRUE	0.07 \pm 0.41	0.00 \pm 0.01
<i>Akkermansia</i>	149	FALSE	FALSE	TRUE	TRUE	3.3 \pm 9.68	6.41 \pm 12.03
<i>Alistipes</i>	142	FALSE	FALSE	FALSE	TRUE	1.28 \pm 3.05	1.30 \pm 2.34
<i>Bacillus</i>	139	FALSE	FALSE	FALSE	TRUE	0.01 \pm 0.03	0.00 \pm 0.01
<i>Enterorhabdus</i>	137	FALSE	FALSE	FALSE	TRUE	0.01 \pm 0.03	0.08 \pm 0.55
<i>Pantoea</i>	137	FALSE	FALSE	FALSE	TRUE	0.01 \pm 0.03	0.00 \pm 0.00
<i>Ruminococcaceae_UCG_004</i>	134	FALSE	FALSE	FALSE	TRUE	0.05 \pm 0.17	0.15 \pm 0.41
<i>Curvibacter</i>	132	FALSE	FALSE	FALSE	TRUE	0.01 \pm 0.01	0.00 \pm 0.00
<i>Granulicatella</i>	131	FALSE	FALSE	FALSE	TRUE	0.31 \pm 1.11	0.08 \pm 0.27
<i>Lachnospiraceae_NC2004_group</i>	131	FALSE	FALSE	FALSE	TRUE	0.02 \pm 0.04	0.05 \pm 0.09
<i>Ruminiclostridium_5</i>	131	FALSE	FALSE	FALSE	TRUE	0.50 \pm 1.11	1.33 \pm 2.79
<i>Epulopiscium</i>	130	FALSE	FALSE	FALSE	TRUE	0.06 \pm 0.30	0.01 \pm 0.04
<i>Robinsoniella</i>	130	FALSE	FALSE	FALSE	TRUE	0.06 \pm 0.33	0.01 \pm 0.07
<i>[Eubacterium]_coprostanoligenes_group</i>	128	FALSE	FALSE	FALSE	TRUE	0.13 \pm 0.39	0.39 \pm 2.19
<i>Eggerthella</i>	128	FALSE	FALSE	FALSE	TRUE	1.39 \pm 2.17	1.86 \pm 2.96
<i>Erwinia</i>	126	FALSE	FALSE	FALSE	TRUE	0.00 \pm 0.00	0.00 \pm 0.00
<i>[Ruminococcus]_gnavreuii_group</i>	124	FALSE	FALSE	FALSE	TRUE	0.09 \pm 0.26	0.45 \pm 2.04

1257 **Table S6. Characteristics of microbial correlation networks associated with different**
1258 **groups.**
1259

Groups	Average degree	Clustering coefficient	Edges	Graph density	Modularity	Nodes
Control	9.475	0.368	938	0.048	0.377	198
Non-CDI Diarrhea	11.314	0.474	1171	0.055	0.271	207
Asymptomatic Carriage	9.730	0.349	973	0.049	0.442	200
CDI	5.200	0.502	299	0.046	0.568	115

1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291

1292 **Table S7. Comparison of host immune markers in different groups. Mean (Q1, Q3); p-value**
 1293 **calculated with Mann-Whitney U test.**
 1294

Immune markers	Control (n=45)	Non-CDI Diarrhea (n=44)	Asymptomatic Carriage (n=35)	CDI (n=99)	P-value (CDI vs. Asymptomatic Carriage)	P-value (CDI vs. Non-CDI Diarrhea)	P-value (CDI vs. Non-CDI)
IgA_toxA	33.63 (7.24, 62.54)	16.94 (8.44, 16.96)	44.59 (10.21, 102.50)	48.41 (12.10, 104)	0.543	< 0.001	0.001
IgG_toxA	19.87 (9.75, 22.43)	24.20 (11.11, 27.77)	22.20 (11.48, 24.86)	40.09 (14.77, 59.18)	0.009	0.002	< 0.001
IgM_toxA	2.04 (0.00, 2.87)	2.45 (0.00, 3.36)	2.15 (0.00, 2.98)	2.09 (0.00, 2.84)	0.316	0.643	0.172
IgA_toxB	9.73 (2.68, 8.61)	18.78 (4.07, 19.18)	23.02 (5.80, 20.12)	35.07 (4.76, 67.79)	0.469	0.102	0.002
IgG_toxB	9.66 (4.30, 9.80)	11.97 (5.92, 13.46)	13.57 (3.98, 15.98)	14.61 (4.92, 17.56)	0.980	0.870	0.379
IgM_toxB	12.12 (2.28, 9.24)	13.72 (2.78, 12.77)	11.08 (2.74, 9.04)	14.99 (1.87, 10.47)	0.261	0.192	0.256
GCSF	11.27 (0.46, 14.17)	49.37 (2.18, 29.36)	20.01 (2.18, 20.49)	386.64 (22.56, 159.95)	< 0.001	< 0.001	< 0.001
IL-10	8.29 (0.00, 3.15)	33.22 (0.00, 14.06)	9.05 (0.00, 9.78)	35.17 (1.63, 27.99)	0.002	0.021	< 0.001
IL-13	1.38 (0.00, 0.00)	1.97 (0.00, 0.34)	5.22 (0.00, 0.00)	9.35 (0.00, 1.09)	0.529	0.593	0.167
IL-15	1.56 (0.00, 0.24)	2.82 (0.00, 3.28)	2.03 (0.00, 1.46)	5.22 (0.19, 5.33)	0.007	0.037	< 0.001
IL-1b	0.05 (0.00, 0.00)	0.11 (0.00, 0.00)	0.44 (0.00, 0.00)	0.7 (0.00, 0.00)	0.920	0.387	0.159
IL-2	0.04 (0.00, 0.00)	0.05 (0.00, 0.00)	0.45 (0.00, 0.00)	1.4 (0.00, 0.00)	0.630	0.151	0.051
IL-4	1.88 (0.00, 0.00)	9.58 (2.57, 12.44)	5.73 (0.00, 0.00)	11.54 (0.00, 9.12)	0.002	0.011	0.032
IL-6	15.78 (0.00, 3.77)	24.97 (0.00, 10.71)	9.52 (0.00, 5.46)	47.09 (2.52, 37.71)	< 0.001	< 0.001	< 0.001
IL-8	100.51 (12.37, 59.19)	80.85 (15.22, 73.45)	59.78 (10.33, 44.76)	128.05 (27.20, 122.24)	< 0.001	0.004	< 0.001
MCPI	477.00 (397.19, 545.00)	591.61 (399.55, 779.32)	613.28 (430.85, 791.37)	844.96 (522.41, 990.98)	0.053	0.020	< 0.001
TNFa	8.61 (6.20, 10.95)	21.76 (8.93, 21.22)	12.45 (4.91, 14.84)	26.68 (13.88, 28.94)	< 0.001	0.006	< 0.001
VEGFA	102.78 (35.76, 118.27)	109.85 (34.53, 127.69)	118.44 (27.54, 188.60)	125.93 (19.54, 140.49)	0.499	0.603	0.390

1295
 1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315

1316 **Table S8. Accuracy, Precision, Recall and F1-score of symbolic classification in CDI**
1317 **diagnosis.** CDI subjects were considered as either true positive or true negative. Results shown
1318 in the parenthesis represents the latter case. The performance of the symbolic classification
1319 model evaluated by cross-validation. We randomly split the dataset to form a training set (80%
1320 of the data) and a test set (20% of the data) in 10 different ways. Each time, for each
1321 classification task (diagnostic goal), we learned the SC model from the training dataset and
1322 evaluated it on the test dataset. Data represents as mean \pm standard deviation.
1323

Diagnostic goal	Accuracy	Precision	Recall	F1-score
CDI vs. Asymptomatic Carriage	0.873 \pm 0.085	0.878 \pm 0.101 (0.857 \pm 0.180)	0.963 \pm 0.043 (0.639 \pm 0.229)	0.915 \pm 0.061 (0.718 \pm 0.185)
CDI vs. Non- CDI Diarrhea	0.883 \pm 0.036	0.899 \pm 0.078 (0.865 \pm 0.098)	0.912 \pm 0.070 (0.847 \pm 0.110)	0.901 \pm 0.039 (0.846 \pm 0.055)
CDI vs. Non- CDI	0.818 \pm 0.040	0.770 \pm 0.102 (0.861 \pm 0.059)	0.784 \pm 0.120 (0.839 \pm 0.091)	0.769 \pm 0.074 (0.845 \pm 0.036)

1324
1325
1326