

1 **Microbial genes outperform species and SNVs as diagnostic markers**
2 **for Crohn's disease on multicohort fecal metagenomes empowered by**
3 **artificial intelligence**

4

5 Sheng Gao^{1,#}, Xiang Gao^{1,#}, Ruixin Zhu^{1,*}, Dingfeng Wu², Zhongsheng Feng¹, Na Jiao², Ruicong
6 Sun¹, Wenxing Gao¹, Qing He^{3,*}, Zhanju Liu^{4,*}, Lixin Zhu^{5,*}

7

8 *1 Department of Gastroenterology, The Shanghai Tenth People's Hospital, School of Medicine,*
9 *School of Life Sciences and Technology, Tongji University, Shanghai 200072, P. R. China.*

10 *2 National Clinical Research Center for Child Health, the Children's Hospital, Zhejiang University*
11 *School of Medicine, Hangzhou 310058, Zhejiang, P. R. China.*

12 *3 Departments of Gastroenterology and Nutrition, the Sixth Affiliated Hospital, Sun Yat-sen*
13 *University, Guangzhou 510655, P.R. China.*

14 *4 Center for IBD Research, Department of Gastroenterology, The Shanghai Tenth People's*
15 *Hospital, School of Medicine, Tongji University, Shanghai 200072, P. R. China.*

16 *5 Department of Colorectal Surgery, Guangdong Institute of Gastroenterology, Guangdong*
17 *Provincial Key Laboratory of Colorectal and Pelvic Floor Diseases, The Sixth Affiliated Hospital,*
18 *Sun Yat-sen University, Guangzhou 510655, P.R. China.*

19

20 #Co-first authors.

21

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

22 ***Corresponding authors:**

23 Lixin Zhu (zhulx6@mail.sysu.edu.cn)

24 Department of Colorectal Surgery, Guangdong Institute of Gastroenterology, Guangdong

25 Provincial Key Laboratory of Colorectal and Pelvic Floor Diseases, The Sixth Affiliated Hospital,

26 Sun Yat-sen University, Guangzhou 510655, P.R. China.

27 Tel: 86-199-4625-6235

28

29 Zhanju Liu (liuzhanju88@126.com)

30 Center for IBD Research, Department of Gastroenterology, The Shanghai Tenth People's Hospital,

31 School of Medicine, Tongji University, Shanghai 200072, P. R. China.

32 Tel: 86-21-6630-1164

33

34 Qing He (heqing5@mail.sysu.edu.cn)

35 Departments of Gastroenterology and Nutrition, the Sixth Affiliated Hospital, Sun Yat-sen

36 University, Guangzhou 510655, P.R. China.

37 Tel: 86-136-0906-1002

38

39 Ruixin Zhu (rxzhu@tongji.edu.cn)

40 Department of Gastroenterology, The Shanghai Tenth People's Hospital, School of Medicine,

41 School of Life Sciences and Technology, Tongji University, Shanghai 200072, P. R. China.

42 Tel: 86-21-6598-1041

43 **Abstract**

44 **Background:** Dysbiosis of gut microbial community is associated with the pathogenesis of CD and
45 may serve as a promising non-invasive diagnostic tool. We aimed to compare the performances of
46 the microbial markers of different biological levels by conducting a multidimensional analysis on
47 the microbial metagenomes of CD.

48

49 **Methods:** We collected fecal metagenomic datasets generated from eight cohorts that altogether
50 include 870 CD patients and 548 healthy controls. The microbial alterations in CD patients were
51 assessed at multidimensional levels including species-, gene- and SNV- level, and then diagnostic
52 models were constructed using artificial intelligence algorithm.

53

54 **Results:** A total of 227 species, 1047 microbial genes and 21877 microbial SNVs were identified
55 that differed between CD and controls. The species-, gene- and SNV- models achieved an average
56 AUC of 0.97, 0.95 and 0.77, respectively. Notably, the gene model exhibited superior diagnostic
57 capability, achieving average AUCs of 0.89 and 0.91 in internal and external validations,
58 respectively. Moreover, the gene model was specific for CD against other microbiome-related
59 diseases. Further, we found that phosphotransferase system (PTS) contributed substantially to the
60 diagnostic capability of the gene model. The outstanding performance of PTS was mainly explained
61 by genes *celB* and *manY*, which demonstrated high predictabilities for CD with the metagenomic
62 datasets and was validated in an independent cohort by qRT-PCR analysis.

63

64 **Conclusions:** Our global metagenomic analysis unravels the multidimensional alterations of the
65 microbial communities in CD, and identifies microbial genes as robust diagnostic biomarkers across
66 geographically and culturally distinct cohorts.

67

68 **Keywords:** Crohn's disease, microbiome biomarkers, non-invasive diagnosis, artificial intelligence,
69 phosphotransferase system

70

71 **Introduction**

72 Crohn's disease (CD), one of the two main forms of inflammatory bowel disease (IBD), is
73 characterized by skip lesions and transmural inflammation of the gastrointestinal tract. The
74 incidence of CD has risen globally in past two decades, causing substantial economic burdens for
75 patients and societies [1, 2]. Currently diagnosis of CD is mainly based on the combined evaluation
76 of endoscopic, radiographic and pathological findings [3, 4]. However, the diagnostic power of
77 endoscopy is often limited by patient compliance, bowel preparation quality and other
78 uncontrollable factors [5]. Therefore, a sensitive, specific and convenient non-invasive diagnostic
79 tool for CD is urgently needed.

80

81 Serologic and fecal biomarkers, such as C-reactive protein and fecal calprotectin, have been used
82 as indicators to evaluate inflammatory activity in IBD [6, 7]. However, the accuracy and specificity
83 of these biomarkers are not satisfactory. Recently, the diagnostic potential of the microbial
84 signatures has emerged as potential diagnostic markers for IBD [8-12]. For instance, Pascal et al.

85 constructed a diagnostic model using microbial species abundance and achieved a sensitivity of
86 81.8% for CD [11]. Similarly, Franzosa et al. reported a model that achieved an area under the ROC
87 curve (AUC) of 0.92 [12]. Along this line, future effort is needed to conduct similar analysis
88 incorporating multiple cohorts of distinct cultural and geographical background to identify markers
89 of universal value.

90

91 Notably, species abundance may not be an accurate representative of the microbial functions as
92 reflected by the fact that the nomenclatures of many gut microbial species are currently and
93 constantly being adjusted. In this regard, the diagnostic value of microbial genes and their
94 polymorphisms has become popular subjects of investigation [13-16] (**Fig. 1B**). For example,
95 microbial functional genes outperformed microbial species in distinguishing CRC from controls
96 [14]. Similarly, a recent study demonstrated high accuracy of microbial SNVs for diagnosing CD
97 [17]. Currently, an integrated investigation on multidimensional signatures of CD at species-, gene-
98 and SNV-levels is missing and seems to be warranted in the clinic.

99

100 In this study, with large numbers of whole metagenome sequencing (WMS) samples from
101 multiple cohorts, we constructed diagnostic models for CD and systematically assessed the
102 predictabilities of multidimensional signatures. Candidate biomarkers for CD diagnosis were
103 identified and further validated by qRT-PCR with an independent cohort. Collectively, these results
104 uncover the multidimensional alterations of the microbial communities in CD patients and provide
105 unbiased and robust biomarkers for CD diagnosis.

106

107 **Methods**

108 **Study inclusion and data acquisition**

109 For discovery dataset, we used PubMed to search for studies that published fecal shotgun
110 metagenomic data of CD patients and controls. Raw FASTQ files of 1241 fecal samples from four
111 studies were downloaded from the European Nucleotide Archive (ENA) including datasets ‘D1-D4’
112 **(Fig. 1A)**.

113

114 For validation dataset, the raw data of 177 samples from three studies were collected from the
115 ENA including datasets ‘V1-V3’ **(Fig. 1A)**. The clinical characteristics of patients were shown in
116 **Table S1**.

117

118 To evaluate whether the prediction model is specific for CD rather than non-CD diseases, we
119 further collected five cohorts of non-CD diseases including ulcerative colitis (UC), colorectal cancer
120 (CRC), type-2 diabetes (T2D), liver cirrhosis (LC) and Parkinson's disease (PD).

121

122 **Patient recruitment and sample collection of Chinese cohorts**

123 The Chinese cohort D5 consisted of 40 CD and 53 control samples **(Table S2)**. The CD patients
124 and controls were enrolled at the Sixth Affiliated Hospital of the Sun Yat-sen University,
125 Guangdong province, China.

126

127 For qRT-PCR validation, we enrolled CD patients and controls at the Shanghai Tenth People's
128 Hospital. Patients with diagnosis of CD were included in the study. Potential participants were
129 excluded if they were pregnant, were diagnosed with indeterminate colitis, had an acute
130 gastrointestinal infection, or had antibiotic therapy within 3 months. In total, we collected 73 fecal
131 samples (N = 37 for CD and N = 36 for control, **Table S3**) that were then stored at -80 °C before
132 DNA extraction. The study was approved by the Institutional Review Board at the Shanghai Tenth
133 People's Hospital, Tongji University, Shanghai (No. 20KT863), and each participant provided
134 informed consent.

135

136 **Quality control of WMS sequencing data**

137 For preprocessing of the WMS sequencing data, quality control was performed using KneadData
138 V0.6.0. Subsequently, reads with length lower than 50bp, or with low quality bases were filtered
139 out by Trimmomatic software (V0.32). Furthermore, reads that mapped to the mammalian genome,
140 bacterial plasmids, UNiVec sequences, and chimeric sequences were removed.

141

142 **Annotation and abundance estimation of microbial taxa, genes and SNVs**

143 For multi-kingdom species level analysis, The customized reference database was constructed with
144 18756 bacterial, 359 archaeal, 9346 viral reference genomes from the NCBI Refseq database
145 (accessed on January 2020), and 1094 fungal reference genomes from the NCBI Refseq database,
146 FungiDB (<http://fungidb.org>) and Ensemble (<http://fungi.ensembl.org>) (all accessed on January
147 2020). Quality-filtered reads were aligned and quantified by Kraken2 [18] and Bracken, respectively.

148 For microbial gene level analysis, we assembled the quality-filtered metagenomes into contigs
149 with Megahit (v1.2.9) [19] using ‘meta-sensitive’ parameters. Contigs shorter than 500-bp were
150 excluded for further analysis. Prodigal (v2.6.3) software [20] was used to predict genes at the
151 metagenome mode (-p meta). A non-redundant microbial gene reference was constructed with CD-
152 HIT [21] using a sequence identity cut-off of 0.95, and a minimum coverage cut-off of 0.9 for the
153 shorter sequences. The reference was annotated with EggNOG mapper (v2.0.1) based on EggNOG
154 orthology data. Subsequently, CoverM (V4.0) was used to estimate gene abundances by mapping
155 reads to the non-redundant reference and to calculate the coverage of genes in the original contigs.
156 The abundance of KEGG orthologous (KOs) groups were calculated by summing the expression of
157 genes annotated to the same KOs.

158 For SNV level analysis, MIDAS was used to perform microbial SNV annotation [22]. A
159 customized reference genome database was constructed to include 7 species with sufficient
160 coverage (>3X) in at least 20% of all samples. Then, the WMS reads were mapped to the database
161 for SNV calling. Subsequently, the SNV profiles of all samples were merged, with only bi-allelic
162 positions chosen. Other parameters were identical with those of the preset option ‘—core_snps’
163 (merge_midias.py snps —core_snps).

164

165 **Diagnostic model construction and evaluation**

166 ***Model construction***

167 Artificial intelligence (AI) algorithm called feedforward neural network (FNN) was employed to
168 construct the diagnostic model. In detail, the hidden layers were activated by rectified linear unit

169 (ReLU) activation function and the output layer was activated by sigmoid function. Subsequently,
170 we performed stratified ten-fold cross-validation to avoid overfitting issue and model estimation
171 using Scikit-learn 1.1.0. Finally, we trained the diagnostic model with well-optimized
172 hyperparameter combinations with TensorFlow 2.8.0. The feature importance was evaluated with
173 SHapley Additive exPlanations (SHAP) [23] to explain the output of machine learning model.

174

175 ***Model interpretation***

176 To better interpret the compositions and corresponding contributions of features in model, we
177 grouped KO genes by gene sets based on the priori knowledge of KEGG database. Subsequently,
178 we randomly shuffled the abundance values of KO genes of a gene set in validation dataset, and
179 performed predictions using the constructed diagnostic model. The decrease of AUC was
180 considered as the importance of gene set to the diagnostic model. The above procedure was repeated
181 for 50 times.

182

183 ***Evaluation of the model's robustness and generalization***

184 To test the robustness and generalization of selected optimal model among distinct cohorts, we
185 performed cohort-to-cohort transfer and leave-one-cohort-out (LOCO) validation as described in
186 our previous studies [24, 25]. For cohort-to-cohort transfer, diagnostic models were trained on one
187 single cohort and validated on each of the remaining cohorts. For LOCO validation, one single
188 cohort was set as the validation dataset while all other cohorts were pooled together as the discovery
189 dataset.

190

191 **Disease specificity assessment of prediction model**

192 Using non-CD diseases samples of UC, CRC, T2D, LC and PD, we evaluated the disease specificity
193 of the predictive model for CD, following the method described by Thomas et al [26]. In detail, we
194 randomly selected 10 control samples and 10 case samples from non-CD external data and added
195 them into the control group in the validation dataset. If the model is specific for CD, the model
196 would not perform worse with the addition of a case relative to the addition of the controls, because
197 the model does not cover the characteristics of non-CD diseases. We repeated the procedure for 50
198 times.

199

200 **Validation of microbial genes by qRT-PCR**

201 The gDNA was extracted with the TIANamp Stool DNA Kit (Cat# 4992205, TIANGEN) according
202 to the manufacturer's instructions. The primers used for validation are listed in **Table S4**. To
203 perform the qRT-PCR analysis, the reaction mixture contained the primer pair with concentrations
204 diluted to 0.2 μ M and 10 ng gDNA in a 10 μ l final volume with the SYBR Green qPCR Mix
205 (Thermo Fisher Scientific). The cycling program was set as indicated: pre-denaturation at 95 $^{\circ}$ C for
206 10 min; denaturation at 95 $^{\circ}$ C for 15 s and annealing at 60 $^{\circ}$ C for 60 s for 40 cycles, followed by
207 melting curve analysis. The qRT-PCR results were quantitated by calculating $-\Delta\Delta$ Ct values
208 between candidate genes and the 16S gene. The significance of the comparison between CD and
209 control samples was tested by a two- sided Wilcoxon rank-sum test ($P < 0.05$).

210

211 **Statistical analysis**

212 *Alpha and beta diversity analysis*

213 Alpha diversity of taxonomic profiles including Shannon, ACE, Simpson and Chao1 index were
214 calculated based on Bray-Curtis distance using R (V4.0.5) “vegan” (V2.5.7) package. Beta diversity
215 between groups were calculated by permutational multivariate analysis of variance (PERMANOVA)
216 called adonis test, and significance was evaluated with 999 permutations.

217

218 *Co-abundance analysis*

219 Firstly, we generated species abundance profiles of CD and controls, respectively. Then we
220 employed SparCC [27] to perform co-abundance analysis of differential multi-kingdom species.
221 Correlations between differential multi-kingdom species were determined with 50 iterations. Then
222 SparCC resampled the original dataset through bootstrap method to obtain random datasets. Later,
223 pseudo-p-values are calculated from these random data sets to assess the significance of the initial
224 observation scores. The statistical significance was calculated with 999 permutations. The network
225 was visualized with Gephi (V0.9.5).

226

227 *Multidimensional signatures association analysis*

228 To further explore the potential associations between multi-dimensional signatures, Hierarchical
229 All-against-All association testing (HALLA, V 0.8.20) [28] was performed. We generated species-,
230 gene-, and SNV-profiles of CD patients and controls, respectively. Subsequently, the associations
231 between the species-, gene-, SNV-signatures were calculated in pairs by HALLA. After that, we

232 merged the output correlation matrices. Correlations with $|\text{cor}| > 0.4$ and P -values < 0.05 were used
233 to constructed the network and visualized with Gephi (V0.9.5).

234

235 **Results**

236 **Characterization of multicohort WMS data and study design**

237 In this study, we collected eight fecal shotgun metagenomics datasets from published studies to
238 characterize the gut microbiome in CD patients compared to healthy controls (**Fig. 1A**). Patients
239 treated with antibiotics were excluded. In total, we included 785 samples from CD patients and 456
240 healthy control samples across geographically distinct regions from U.S. and China as the discovery
241 dataset. In addition, 85 CD samples and 92 controls from three independent cohorts from U.S.,
242 Spain and Netherlands were included as the validation dataset. The overall protocol for this study
243 (**Fig. 1C**) was based on the workflow of a previous study [24] with modifications.

244

245 **Multidimensional alterations in gut microbial profiles in CD patients**

246 At species level, we found that alpha and beta diversities were significantly differed between CD
247 patients and controls (**Fig. 2A-B**). A total of 80 bacterial species were identified with significantly
248 different abundances between CD and control, such as *Escherichia coli*, *Flavonifractor plautii*,
249 *Klebsiella pneumoniae* and *Bacteroides intestinalis*. (**Fig. 2C; Table S5**). Besides, 147 non-
250 bacterial species including 70 fungus, 42 viruses and 35 archaea exhibited differential abundances
251 between CD and controls, such as *Aspergillus rambellii*, *Capronia epimyces*, *Bacteroides phage*
252 *B124-14*, *Klebsiella virus KpV80* and *DPANN group archaeon LC1Nh* (**Fig. S1 and Table S5**).

253 Further, we investigated the differences in microbial interactions between CD and controls by
254 performing co-abundance analysis via SparCC. Interestingly, interactions among intra-kingdom
255 species were more frequently observed in the network of CD, compared to the network of controls
256 (**Fig. S2**), indicating large scale alterations in the structure and function of the gut microbiome in
257 CD.

258

259 Next, we assessed the microbial alterations at KO gene level, and identified 497 genes with
260 increased abundance and 1043 genes with decreased abundance in CD patients, such as the genes
261 encoding peptide/nickel transport system permease protein (*ABC.PE.P*), mannose PTS system EIIC
262 component (*manY*), flagellin (*fliC*) and cellobiose PTS system EIIC component (*celB*) (**Fig. 2E**;
263 **Table S6**). For better understanding of these differential KO genes, we performed gene set
264 enrichment analysis (GSEA). 59 enriched pathways, including 18 pathways with increased
265 abundances and 41 with decreased abundances in CD patients, were identified (**Fig. S3A and Table**
266 **S7**). Propanoate metabolism, quorum sensing, phosphotransferase system (PTS) and purine
267 metabolism exhibited increased abundances in CD, while biosynthesis of secondary metabolites,
268 pantothenate and CoA biosynthesis exhibited decreased abundances in CD.

269

270 For microbial SNV level analysis, a total of 7 commonly observed species that have sufficient
271 coverage (> 3X) in at least 20% of the samples were annotated, with the number of SNVs ranging
272 from 74 with *Bacteroides rodentium* to 99305 with *Bacteroides vulgatus* (**Fig. S3B and Fig. S4**).
273 A total of 21877 differential SNVs were identified in the seven annotated species (**Fig. S3C**). For

274 instance, *Bacteroides vulgatus*, belonging to the most commonly encountered *Bacteroides* species
275 in the human colon, had 11134 significantly differential SNVs that located on genes such as *panC*,
276 *rodA* and *ruvB*. (**Fig. 2D; Table S8**). These differential SNVs are potential candidates of risk factors
277 mediating abnormal gene functions. Collectively, we systemically assessed the multidimensional
278 microbial alterations in CD patients compared to controls, and identified differential signatures for
279 diagnostic model construction.

280

281 **Diagnostic models for CD based on microbial multidimensional signatures**

282 Based on all of the differential signatures at species-, gene- and microbial SNV-levels, we
283 constructed models using deep learning algorithm. At species level, we firstly evaluated the
284 capability of single-kingdom species for distinguishing CD from controls. The average AUCs of
285 cross-validation based on fungal, viral, archaeal signatures were 0.89, 0.81 and 0.76, respectively.
286 Compared to non-bacterial species, bacterial species demonstrated a better performance in disease
287 prediction (average AUC=0.94) (**Fig. S5A-D**). Furthermore, we merged single-kingdom signatures
288 together, and found that the species model based on multi-kingdom signatures had higher diagnostic
289 accuracy with an average AUC of 0.97 (**Fig. 3A; Fig. S5E**). Interestingly, we noticed that several
290 fungal species including *A. rambellii* and *A. ochraceoroseus* were top-ranking features of the model
291 with high SHAP values, suggesting their close association with CD pathology (**Fig. S6A and Table**
292 **S9**).

293

294 Subsequently, we constructed a diagnostic model based on 1047 differential KO genes. The gene
295 model achieved an average AUC of 0.95 in 10-fold cross-validation, slightly lower than that of the
296 multi-kingdom species model (**Fig. 3A**). From feature importance evaluation, we found that CDP-
297 abequeose synthase (*rfbJ*), type VI secretion system protein ImpB (*impB*), nitrite reductase (NO-
298 forming) (*nirK*), and *celB* were the most important KO genes with SHAP values ranged from 0.006
299 to 0.008 (**Table S10**). Notably, the KO gene *celB* was found to be significantly increased in CD
300 patients of each dataset (**Fig. S6B**), suggesting an outstanding contribution of *celB* gene to the
301 diagnostic power of the model.

302

303 Furthermore, we explored the diagnostic potentials of microbial SNVs. The SNV model achieved
304 an average AUCs of 0.77 in cross-validation (**Fig. 3A**). The most important SNVs were mainly from
305 *Bacteroides* species including *B. ovatus*, *vulgatus* and *uniformis* (**Fig. S7A and Table S11**). As the
306 most widely colonized microbes in the gut [29], *Bacteroides* species contributes to the major
307 diagnostic power of the SNV model in our results.

308

309 Finally, we constructed a model with the combination of species-, gene- and SNV-signatures (**Fig.**
310 **S8D**). The combined model achieved an average AUC of 0.95. Interestingly, the performance of
311 combined model was not significantly improved compared to species- and gene- models, and most
312 of the top-ranking features were from KO genes (**Fig. 3A; Fig. S7B**). These results suggest that the
313 gene signatures are the most powerful biomarkers for CD.

314

315 **Gene model achieves superior robustness and generalization**

316 To assessed the robustness and generalization of species-, gene-, SNV- and combined models, we
317 performed internal and external validations. With the internal validation cohorts, the gene model
318 achieved the highest average AUCs of 0.87 and 0.89 in cohort-to-cohort transfer and LOCO
319 validation, respectively (**Fig. 3C-D**), compared to other diagnostic models (**Fig. S9A-G**). In external
320 validation, the gene model also exhibited the best performance with an average AUC of 0.91 in
321 three independent cohorts (**Fig. 3B**). Taken together, the gene model demonstrated superior
322 robustness compared to the species-, SNV-model and even the combined model.

323

324 **Gene model is highly specific for CD**

325 To ascertain the discriminative power of the gene model, that is, the model is specific for CD but
326 not other microbiome-related diseases, we chose five microbiome-related diseases including UC,
327 CRC, PD, T2D and LC to evaluate the disease specificity of the gene model. Adding UC samples
328 into three independent validation cohorts decreased the AUC by 6.6, 10.1 and 12.9%, respectively
329 (**Fig. 3E**). These changes were not significant considering the baseline values of the altered AUCs
330 when adding CD samples in the validation dataset (decreased AUCs by 10.7, 17.5 and 20.5%,
331 respectively, **Fig. S9H**). With CRC cohorts, slight and insignificant changes of AUCs in validation
332 (decreased by 0.7, 1.1 and 1.3%, respectively) were observed. Similarly, slight and insignificant
333 changes of AUC were observed in validations with T2D (increased by 2.0, 3.4 and 4.0%,
334 respectively), liver cirrhosis (increased by 1.5, 2.5 and 3.0%, respectively) and PD (increased by
335 1.6, 2.6 and 3.1%, respectively). Altogether, the slight changes in AUCs suggest limited effects of

336 the samples with non-CD diseases on the CD model, indicating that our diagnostic model is specific
337 for CD.

338

339 **Outstanding contributions of phosphotransferase system to the diagnostic capability of the**
340 **gene model**

341 To evaluate the respective contributions of each gene set and of key gene feature in the gene model,
342 the KO gene features were grouped by gene set and the importance of each gene set was evaluated
343 as described in Methods section. Relative to the baseline AUC of 0.91, the abundance disturbance
344 of the gene sets quorum sensing, PTS and ABC transporters caused the greatest decrease of AUC
345 in the predictive model by 1.09 to 1.70 percent (**Fig. 4A**). Further, we performed recursive feature
346 elimination by gene sets and reconstructed diagnostic models. We found that the AUC of cross-
347 validation did not decrease significantly until the glycerolipid metabolism gene set was eliminated,
348 which confirmed the important contribution of quorum sensing, PTS, ABC transporters, fructose
349 and mannose metabolism and glycerolipid metabolism to the diagnostic model (**Fig. S10A**). To
350 further strengthen these results, we constructed a sub-model with genes of these five gene sets,
351 which achieved an AUC of 0.89 in cross-validation (**Fig. S10B**). The sub-model displayed decent
352 robustness in internal validations and achieved an average AUC of 0.81 in external validation (**Fig.**
353 **S10C**). Notably, we found that *celB* was the most important feature in the sub-model (**Table S12**).
354 These results suggest that the above identified gene sets are the key contributors to diagnostic
355 capabilities of the gene model.

356

357 Next, we assessed the prediction power of representative KO genes of each gene set (**Table S13**).
358 Notably, *celB* and *manY* displayed excellent diagnostic capabilities with AUCs of 0.74 and 0.71,
359 respectively (**Fig. 4B**). Since *celB* and *manY* (also a member of fructose and mannose metabolism)
360 are both members of PTS, the above results indicate that PTS gene set mediated the most significant
361 functional alterations of gut microbiome in CD patients. Finally, we validated the abundances of
362 *celB* and *manY* with an independent cohort of CD patients and controls by qRT-PCR. Consistent
363 with the metagenomic data (**Fig. 4C**), both *celB* and *manY* were significantly more abundant in CD
364 patients (**Fig. 4D**). Additionally, we validated the abundances of those genes that belong to
365 important pathways and with high feature importance by qRT-PCR (**Fig. S11**). These results
366 revealed the respective contributions of individual gene feature to the diagnostic capability of the
367 gene model, and identified *celB* and *manY* as the individual biomarkers with the highest predictive
368 power for diagnosing CD.

369

370 **Altered interactions within and between each level of microbial signatures in CD**

371 For a global understanding of the interactions among all the microbial signatures in CD, we
372 investigated the associations among all the microbial signatures via HALLA (**Fig. 5A-B**). In both
373 CD and control networks, considerable associations were observed between KO genes and species,
374 but few observed between SNVs and the other two levels ($|\text{correlation}| > 0.4$) (**Fig. 5B, E**). More
375 associations were observed in the network of CD (206 associations) (**Fig. 5D**), than in the network
376 of controls (163 associations) (**Fig. 5G; Table S14-15**). Interestingly, there were more negative
377 associations between the gene- and the species-signatures in the control network than that in the CD

378 network. For example, D-nopaline dehydrogenase (*nos*), type IV secretion system protein TrbJ (*trbJ*)
379 genes were negatively associated with *R. hominis*, *R. bassiana*, and *C. aerofaciens*. Notably, we
380 found that KO genes had stronger degree centrality than species in the CD network (**Fig. 5C**).
381 Moreover, compared to the control network, these KO genes in CD tended to form isolated clusters,
382 as exemplified by the independent module consisting *celB* and *manY* in the CD network (**Fig. 5A**).

383

384 **Discussion**

385 Here, for the first time, multidimensional microbial signatures of CD were systematically analyzed
386 with multiple cohorts of distinct cultural and geographical backgrounds. In comparison of the
387 diagnostic capabilities of the microbial signatures including differential species, genes and SNVs,
388 the gene model achieved superior accuracy and robustness in distinguishing CD from controls, and
389 the gene model was specific for CD against other microbiome-related diseases. Finally, the major
390 contributing genes in the gene model were identified and validated.

391

392 The multidimensional alterations of the gut microbiome in CD patients contain massive
393 information that could predict the disease state. Therefore, we employed deep learning method to
394 fit the underlying characteristics of gut microbiome in CD. With the microbial species models, while
395 bacterial species achieved the best performance among single-kingdom models, multi-kingdom
396 model with both bacteria and non-bacterial species achieved better accuracy than the single-
397 kingdom models, which is similar to our observations with the microbial models for colorectal
398 cancer [14].

399 Comparing models of three different types, the gene model demonstrated the best generalization
400 and robustness in model evaluations compared to the species-, SNV- and combined models. This is
401 reasonable, considering that the homologous genes of different microorganisms may contribute to
402 the same abnormalities in the gut microbiome in connection to specific pathological processes [30,
403 31].

404

405 In examining the contributions of individual gene set and gene to the diagnostic capabilities of
406 the gene model, we found that genes that belong to PTS gene set had great impacts on the model
407 accuracy in abundance disturbance analysis. The importance of the PTS gene set in the diagnosis
408 model was also demonstrated in recursive feature elimination analysis and in cross-validation of the
409 sub-model. In gut bacteria, the PTS is known as a system that catalyze sugar transport as well as
410 sugar phosphorylation [32, 33]. In addition, the PTS regulates a wide variety of transport, metabolic
411 processes, biofilm formation and virulence [34], thus it is considered as a comprehensive regulation
412 and coordination system. We observed that the CD patients exhibited increased abundance in PTS ,
413 and that the KO genes in PTS were associated with the differential species in CD (**Fig. S12A**), such
414 as *Streptococcus pneumoniae* and *A. ochraceoroseus* that are associated with gut diseases [35-37].
415 That is, PTS may participate in the pathogenesis of CD.

416 More importantly, the KO gene *celB* that encodes the enzyme IIC component (EIIC) of cellobiose
417 PTS, exhibited the top-ranking predictability among all the gene markers and an increased
418 abundance in CD patients. These observations support an outstanding potential for the microbial
419 gene *celB* of PTS to be used as a biomarker for non-invasive CD diagnosis. Moreover, *celB* was

420 associated with *K. pneumonia* (**Fig. S12B**), which is in line with the roles of *celB* component of
421 PTS in biofilm formation and the virulence of *K. pneumoniae* [38], and the roles of *K. pneumoniae*
422 in the initiation and perpetuation of the pathological damage of CD were also demonstrated [39].
423 We also observed a significant increase of *K. pneumoniae* in CD patients (**Fig. S12C and Table**
424 **S5**). Therefore, it is reasonable to hypothesize that the interaction between *celB* and *K. pneumoniae*
425 may contribute to the development of CD.

426 Another microbial gene *manY* was also identified as a biomarker for CD diagnosis. *manY* encodes
427 the EHC component of mannose PTS system (man-PTS) that is a part of the PTS regulatory network.
428 The hairpin tips of IIC in man-PTS is coordinated with mannose and mediates the mannose transport
429 [40]. Interestingly, previous studies found that man-PTS and cellobiose-PTS were upregulated in
430 gut microbes by changing from a low-fat diet to a high-fat, high-sugar diet [40, 41], suggesting that
431 the PTS of gut microbes is sensitive to the nutritional environment of mucosal surfaces. Thus, the
432 up-regulations of *celB* and *manY* in CD likely indicate the up-regulation in the biological activities
433 of cellobiose-PTS and man-PTS in association with CD pathology. That is, *manY* may also
434 participate in the pathogenesis of CD. However, the cause of these alterations in CD is not clear and
435 requires further investigation.

436

437 Our work takes advantages of excellent adaptability and learning ability of AI in large dataset
438 and provides an effective non-invasive diagnostic tool with improved discrimination power for CD.
439 However, the interpretability of the AI model is limited, which could be improved with a causal
440 analysis in the future.

441

442 **Conclusions**

443 Our global metagenomic analysis unravels the multidimensional alterations of the microbial
444 communities in CD and identifies microbial genes as robust diagnostic biomarkers across cohorts.

445 These genes are functionally related to the pathogenesis of CD. Future study on these genes may
446 lead to an effective non-invasive diagnostic tool for CD.

447

448 **Ethics approval and consent to participate**

449 All participants provided written informed consent prior to data collection. The study was approved
450 by the Institutional Review Board at the Shanghai Tenth People's Hospital, Tongji University,
451 Shanghai (No. 20KT863).

452

453 **Consent for publication**

454 Not applicable.

455

456 **Availability of data and materials**

457 The data that support the findings of this study are available from the corresponding author, upon
458 reasonable request. The code and scripts are available on GitHub
459 (<https://github.com/tjcadd2020/Diagnosis-for-CD>).

460

461 **Competing interests**

462 The authors declare that they have no competing interests.

463

464 **Authors' contributions**

465 LZ, ZL, QH and RZ conceived and designed the project. SG and XG drafted the manuscript. RZ,
466 DW, ZF, NJ, RS, WG, QH, ZL and LZ revised the manuscript. All authors read and approved the
467 final manuscript.

468

469 **Acknowledgements**

470 This work was supported by the National Natural Science Foundation of China (82170542 to RZ,
471 92251307 to RZ, 32200529 to DW, 82000536 to NJ), the National Key Research and Development
472 Program of China (2021YFF0703700/2021YFF0703702 to RZ), and Guangdong Province “Pearl
473 River Talent Plan” Innovation and Entrepreneurship Team Project (2019ZT08Y464 to LZ). The
474 funders had no role in study design, data collection and analysis, decision to publish, or preparation
475 of the manuscript.

476

477 **Abbreviations**

478 **ABC.PE.P**: peptide/nickel transport system permease protein; **agaF**: N-acetylgalactosamine PTS
479 system EIIA component; **AKR1A1**: alcohol dehydrogenase (NADP+); **ALDH**: aldehyde
480 dehydrogenase (NAD+); **allA**: ureidoglycolate lyase; **AUC**: area under the ROC curve; **CD**:
481 Crohn’s disease; **celB**: cellobiose PTS system EIIC component; **CRC**: colorectal cancer; **EIIC**:
482 enzyme IIC component; **ENA**: European Nucleotide Archive; **fliC**: flagellin; **FNN**: Feedforward
483 neural network; **GSEA**: gene set enrichment analysis; **IBD**: inflammatory bowel disease; **impB**:

484 type VI secretion system protein ImpB; **KO**: KEGG Orthology; **LC**: liver cirrhosis; **LOCO**: leave-
485 one-cohort-out; **maeB**: malate dehydrogenase (oxaloacetate-decarboxylating) (NADP+); **manY**:
486 mannose PTS system EIIC component; **nirK**: nitrite reductase (NO-forming); **pckA**:
487 phosphoenolpyruvate carboxykinase (GTP); **PD**: Parkinson's disease; **PTS**: phosphotransferase
488 system; **ReLU**: rectified linear unit; **rfbJ**: CDP-abequose synthase; **ROC**: receiver operating
489 characteristic; **SHAP**: SHapley Additive exPlanations; **SNVs**: single nucleotide variants; **sucD**:
490 succinyl-CoA synthetase alpha subunit; **T2D**: type-2 diabetes; **tcPp**: toxin coregulated pilus
491 biosynthesis protein P; **tmoC**: toluene monooxygenase system ferredoxin subunit; **trbJ**: type IV
492 secretion system protein TrbJ; **ttuC**: artrate dehydrogenase/decarboxylase / D-malate
493 dehydrogenase; **UC**: ulcerative colitis; **WMS**: whole metagenome sequencing

494

495 **References**

- 496 1. Ng SC, Shi HY, Hamidi N, Underwood FE, Tang W, Benchimol EI, et al. Worldwide incidence
497 and prevalence of inflammatory bowel disease in the 21st century: a systematic review of
498 population-based studies. *Lancet*. 2017;390(10114):2769-78.
- 499 2. The global, regional, and national burden of inflammatory bowel disease in 195 countries and
500 territories, 1990-2017: a systematic analysis for the Global Burden of Disease Study 2017. *Lancet*
501 *Gastroenterol Hepatol*. 2020;5(1):17-30.
- 502 3. Feuerstein JD, Cheifetz AS. Crohn Disease: Epidemiology, Diagnosis, and Management.
503 *Mayo Clin Proc*. 2017;92(7):1088-103.

- 504 4. Maaser C, Sturm A, Vavricka SR, Kucharzik T, Fiorino G, Annese V, et al. ECCO-ESGAR
505 Guideline for Diagnostic Assessment in IBD Part 1: Initial diagnosis, monitoring of known IBD,
506 detection of complications. *J Crohns Colitis*. 2019;13(2):144-64.
- 507 5. Feld LD, Kirk K, Feld AD. A High Quality Approach to Addressing Complications of
508 Endoscopy and Optimizing Risk Management Strategies. *Techniques and Innovations in*
509 *Gastrointestinal Endoscopy*. 2022.
- 510 6. Lewis JD. The utility of biomarkers in the diagnosis and therapy of inflammatory bowel
511 disease. *Gastroenterology*. 2011;140(6):1817-26 e2.
- 512 7. Mosli MH, Zou G, Garg SK, Feagan SG, MacDonald JK, Chande N, et al. C-Reactive Protein,
513 Fecal Calprotectin, and Stool Lactoferrin for Detection of Endoscopic Activity in Symptomatic
514 Inflammatory Bowel Disease Patients: A Systematic Review and Meta-Analysis. *Am J*
515 *Gastroenterol*. 2015;110(6):802-19; quiz 20.
- 516 8. Dubinsky M, Braun J. Diagnostic and Prognostic Microbial Biomarkers in Inflammatory
517 Bowel Diseases. *Gastroenterology*. 2015;149(5):1265-74.e3.
- 518 9. Lloyd-Price J, Arze C, Ananthakrishnan AN, Schirmer M, Avila-Pacheco J, Poon TW, et al.
519 Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature*.
520 2019;569(7758):655-62.
- 521 10. Huang Q, Zhang X, Hu Z. Application of Artificial Intelligence Modeling Technology Based
522 on Multi-Omics in Noninvasive Diagnosis of Inflammatory Bowel Disease. *J Inflamm Res*.
523 2021;14:1933-43.

- 524 11. Pascal V, Pozuelo M, Borruel N, Casellas F, Campos D, Santiago A, et al. A microbial
525 signature for Crohn's disease. *Gut*. 2017;66(5):813-22.
- 526 12. Franzosa EA, Sirota-Madi A, Avila-Pacheco J, Fornelos N, Haiser HJ, Reinker S, et al. Gut
527 microbiome structure and metabolic activity in inflammatory bowel disease. *Nat Microbiol*.
528 2019;4(2):293-305.
- 529 13. Tierney BT, Tan Y, Kostic AD, Patel CJ. Gene-level metagenomic architectures across
530 diseases yield high-resolution microbiome diagnostic indicators. *Nat Commun*. 2021;12(1):2907.
- 531 14. Liu NN, Jiao N, Tan JC, Wang Z, Wu D, Wang AJ, et al. Multi-kingdom microbiota analyses
532 identify bacterial-fungal interactions and biomarkers of colorectal cancer across cohorts. *Nat*
533 *Microbiol*. 2022;7(2):238-50.
- 534 15. Yan Y, Nguyen LH, Franzosa EA, Huttenhower C. Strain-level epidemiology of microbial
535 communities and the human microbiome. *Genome Med*. 2020;12(1):71.
- 536 16. Ma C, Chen K, Wang Y, Cen C, Zhai Q, Zhang J. Establishing a novel colorectal cancer
537 predictive model based on unique gut microbial single nucleotide variant markers. *Gut Microbes*.
538 2021;13(1):1-6.
- 539 17. Jiang S, Chen D, Ma C, Liu H, Huang S, Zhang J. Establishing a novel inflammatory bowel
540 disease prediction model based on gene markers identified from single nucleotide variants of the
541 intestinal microbiota. *iMeta*. 2022;1(3).
- 542 18. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome Biol*.
543 2019;20(1):257.

- 544 19. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node solution
545 for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*.
546 2015;31(10):1674-6.
- 547 20. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic
548 gene recognition and translation initiation site identification. *BMC Bioinformatics*. 2010;11(1):119.
- 549 21. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation
550 sequencing data. *Bioinformatics*. 2012;28(23):3150-2.
- 551 22. Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. An integrated metagenomics pipeline
552 for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome Res*.
553 2016;26(11):1612-25.
- 554 23. Scott M, Lundberg S-IL. A Unified Approach to Interpreting Model Predictions. *NIPS'17:*
555 *Proceedings of the 31st International Conference on Neural Information Processing Systems*.
556 2017:4768–77.
- 557 24. Gao W, Chen W, Yin W, Zhu X, Gao S, Liu L, et al. Identification and validation of microbial
558 biomarkers from cross-cohort datasets using xMarkerFinder. *Protocol Exchange*. 2022.
- 559 25. Wu Y, Jiao N, Zhu R, Zhang Y, Wu D, Wang AJ, et al. Identification of microbial markers
560 across populations in early detection of colorectal cancer. *Nat Commun*. 2021;12(1):3063.
- 561 26. Thomas AM, Manghi P, Asnicar F, Pasolli E, Armanini F, Zolfo M, et al. Metagenomic
562 analysis of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures and a
563 link with choline degradation. *Nat Med*. 2019;25(4):667-78.

- 564 27. Friedman J, Alm EJ. Inferring Correlation Networks from Genomic Survey Data. PLOS
565 Computational Biology. 2012;8(9):e1002687.
- 566 28. Ghazi AR, Sucipto K, Rahnavard G, Franzosa EA, McIver LJ, Lloyd-Price J, et al. High-
567 sensitivity pattern discovery in large, paired multi-omic datasets. bioRxiv. 2021:2021.11.11.468183.
- 568 29. Kim S, Covington A, Pamer EG. The intestinal microbiota: Antibiotics, colonization resistance,
569 and enteric pathogens. Immunol Rev. 2017;279(1):90-105.
- 570 30. Schirmer M, Garner A, Vlamakis H, Xavier RJ. Microbial genes and pathways
571 in inflammatory bowel disease. Nat Rev Microbiol. 2019;17(8):497-511.
- 572 31. Rowland I, Gibson G, Heinken A, Scott K, Swann J, Thiele I, et al. Gut microbiota functions:
573 metabolism of nutrients and other food components. Eur J Nutr. 2018;57(1):1-24.
- 574 32. Lengeler JW, Jahreis K. Bacterial PEP-dependent carbohydrate: phosphotransferase systems
575 couple sensing and global control mechanisms. Contrib Microbiol. 2009;16:65-87.
- 576 33. Västermark A, Saier MH, Jr. The involvement of transport proteins in transcriptional and
577 metabolic regulation. Curr Opin Microbiol. 2014;18:8-15.
- 578 34. Saier MH, Jr. The Bacterial Phosphotransferase System: New Frontiers 50 Years after Its
579 Discovery. J Mol Microbiol Biotechnol. 2015;25(2-3):73-8.
- 580 35. Aymeric L, Donnadieu F, Mulet C, du Merle L, Nigro G, Saffarian A, et al. Colorectal cancer
581 specific conditions promote *Streptococcus gallolyticus* gut colonization. Proc Natl Acad Sci U S A.
582 2018;115(2):E283-e91.

- 583 36. Lin Y, Lau HC, Liu Y, Kang X, Wang Y, Ting NL, et al. Altered Mycobiota Signatures and
584 Enriched Pathogenic *Aspergillus rambellii* Are Associated With Colorectal Cancer Based on
585 Multicohort Fecal Metagenomic Analyses. *Gastroenterology*. 2022;163(4):908-21.
- 586 37. Coker OO, Nakatsu G, Dai RZ, Wu WKK, Wong SH, Ng SC, et al. Enteric fungal microbiota
587 dysbiosis and ecological alterations in colorectal cancer. *Gut*. 2019;68(4):654-62.
- 588 38. Wu M-C, Chen Y-C, Lin T-L, Hsieh P-F, Wang J-T, Camilli A. Cellobiose-Specific
589 Phosphotransferase System of *Klebsiella pneumoniae* and Its Importance in Biofilm Formation and
590 Virulence. *Infection and Immunity*. 2012;80(7):2464-72.
- 591 39. Rashid T, Ebringer A, Wilson C. The role of *Klebsiella* in Crohn's disease with a potential for
592 the use of antimicrobial measures. *Int J Rheumatol*. 2013;2013:610393.
- 593 40. Jeckelmann JM, Erni B. The mannose phosphotransferase system (Man-PTS) - Mannose
594 transporter and receptor for bacteriocins and bacteriophages. *Biochim Biophys Acta Biomembr*.
595 2020;1862(11):183412.
- 596 41. Turnbaugh PJ, Ridaura VK, Faith JJ, Rey FE, Knight R, Gordon JI. The effect of diet on the
597 human gut microbiome: a metagenomic analysis in humanized gnotobiotic mice. *Sci Transl Med*.
598 2009;1(6):6ra14.

599

600 **Figure Legends**

601

602 **Fig. 1. Overview of the fecal samples included in this study and the analysis protocol.** A We
603 collected a total of 1418 samples from 8 cohorts with fecal shotgun metagenomic data. The
604 discovery dataset included D1, D2, D3, D4 and D5. The validation dataset included V1, V2 and V3.

605 **B** Three levels of analysis were conducted in this study: species-, gene- and microbial SNV- levels.

606 **C** The overall workflow of the study: Firstly, the microbial alterations were identified to retrieve

607 the differential multidimensional signatures of gut microbiome. Subsequently, diagnostic models

608 were constructed and the optimal model was selected according to the performances of the models

609 in internal and external validations. Finally, disease specificity of model was evaluated and model

610 interpretation was conducted for final determination of the microbial biomarker, and then

611 biomarkers were validated by qRT-PCR analysis.

612

613

614 **Fig. 2. Multidimensional alterations in the gut microbiome of CD patients at species-, gene-**

615 **and SNV-levels. A** Alpha diversity measured by Shannon, ACE, Simpson and Chao1 index of

616 patients with CD (orange, n = 785) and control individuals (blue, n = 456); * $P < 0.05$, ** $P < 0.01$,

617 *** $P < 0.001$ and **** $P < 0.0001$. **B** Principal coordinate analysis (PCoA) of samples from all five

618 cohorts based on Bray–Curtis distance, which shows that microbial compositions were different

619 between groups ($R^2 = 0.0265$, $P = 0.001$). P values of beta diversity based on Bray–Curtis distance

620 were calculated with PERMANOVA by 999 permutations (two-sided test). **C** Phylogenetic tree

621 showing the differential bacteria species, grouped by the phyla. The differential species in each

622 dataset are shown in each circle ‘D1-D5’ ($P < 0.05$, two-sided test); the meta-analysis results in

623 integrated dataset were marked by ‘All’. Increased and decreased abundances are indicated by red

624 and blue, respectively. **D** The chord diagram shows the distributions of annotated SNVs in

625 *Bacteroides vulgatus* genome. The outer circle represents the genome of *B. vulgatus*; the inner

626 circles represent the GC-content (cyan indigo lines), sequencing depth (purple lines) and sites of
627 differential SNVs (brown points) in the genome, respectively. **E** UpSet plot showing the number of
628 differential KO genes identified via MaAsLin2 in each dataset and those shared by the datasets.
629 The number above each column represents the intersection size of differential KO genes. The
630 connected dots represent the common differential genes across connected cohorts. The set size on
631 the right represents the number of differential genes in each cohort.

632

633 **Fig. 3. The performance of diagnostic models constructed with multidimensional signatures.**

634 **A** The ROC curves from ten-fold cross-validation of species-, gene-, SNV- and combined diagnostic
635 models. **B** The AUCs of species-, gene-, SNV- and combined diagnostic models in external
636 validation dataset. **C** The AUCs of each model in cohort-to-cohort validation. Each number
637 represents the average AUC of validation with the cohort specified by its column tag as the training
638 cohort, and all other cohorts as the validation cohorts. **D** The AUC of each model in LOCO
639 validation. Each number represents the resulting AUC of validation with the cohort specified by its
640 column tag as the validation cohort while the other cohorts combined as training cohort. **E**
641 Prediction performances as AUC values on the validation cohorts when adding an external set of
642 control and case samples from non-CD disease cohorts (ulcerative colitis (UC), colorectal cancer
643 (CRC), type-2 diabetes (T2D), liver cirrhosis (LC) and Parkinson's disease (PD)). Gray and colored
644 bars are the AUCs after adding control and case samples from the non-CD disease cohorts,
645 respectively.

646

647 **Fig. 4. The model interpretation of the gene model.** **A** The left column lists the average percent
648 change of AUC after shuffling the abundance values of the genes in each gene set in validation
649 dataset with the background color indicating the degrees of AUC change; the center left column
650 lists the number of KO genes in each gene set with the background color indicating the set size; the
651 center right column is the representative signature of each gene set; and the right column lists the
652 cross-validation AUC of the representative microbial gene with the background color indicating an
653 increased (red) or decreased (blue) AUC. The line plot shows the values of feature importance of
654 the representative signatures (upper horizontal axis); the box plot shows the AUCs of each gene set
655 in validation dataset with the dotted line representing the baseline AUC of 0.91 (lower horizontal
656 axis). **B** The ROC curve shows diagnostic performance of microbial gene *celB* and *manY*,
657 respectively. **C-D** The box plot shows the abundances of *celB* (upper) and *manY* (lower) in
658 metagenomic data (C) and qRT-PCR data (D) (N=37, CD; N=36, control), respectively. Data are
659 presented as mean \pm standard deviation. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ and **** $P < 0.0001$.
660

661 **Fig. 5. The cross-talk between multidimensional signatures.** **A-B** Correlations among species-,
662 gene- and SNV- signatures in the control (A) and the CD (B) networks. The colors of nodes indicate
663 signatures of different levels: species (green), KO genes (yellow), and SNVs (brown). Red line
664 indicates positive interaction; and blue line indicates negative interaction ($|\text{correlation}| > 0.4$, FDR <
665 0.05). **C** Density distribution of correlations between different levels of signatures (FDR < 0.05) in
666 the control network. **D** Density distribution of node degrees for different levels of signatures
667 ($|\text{correlation}| > 0.4$, FDR < 0.05) in the control network. **E** Numbers of edges ($|\text{correlation}| > 0.4$, FDR

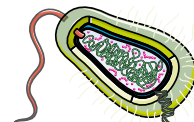
668 < 0.05) in the control network. **F** Density distribution of correlations between different levels of
669 signatures (FDR < 0.05) in the CD network. **G** Density distribution of node degrees for different
670 levels of signatures ($(|correlation|>0.4, FDR < 0.05)$ in the CD network. **H** Number of edges
671 ($(|correlation|>0.4, FDR < 0.05)$ in the CD network.
672

A

	Dataset	CD	Control	Country
Discovery dataset	D1	248	322	U.S.
	D2	89	21	U.S.
	D3	68	34	U.S.
	D4	340	26	U.S.
	D5	40	53	China (in-house)
Validation dataset	V1	44	49	U.S.
	V2	21	23	Spain
	V3	20	20	Netherlands

B

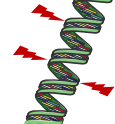
Species



Gene



SNV

**C**