

1 **Shared and disease-specific host gene-microbiome interactions across human**
2 **diseases**

3

4 Sambhawa Priya^{1,3}, Michael B. Burns⁴, Tonya Ward⁵, Ruben A. T. Mars⁶, Beth
5 Adamowicz¹, Eric F. Lock⁷, Purna C. Kashyap⁶, Dan Knights^{5,8}, Ran Blekhman^{1,2,*}

6

7 *Corresponding author: blekhman@umn.edu

8

9 ¹ Department of Genetics, Cell Biology and Development, University of Minnesota,
10 Minneapolis, MN 55455, USA

11 ² Department of Ecology, Evolution, and Behavior, University of Minnesota,
12 Minneapolis, MN 55455, USA

13 ³ Bioinformatics and Computational Biology, University of Minnesota, Minneapolis, MN
14 55455, USA

15 ⁴ Department of Biology, Loyola University Chicago, Chicago, IL 60660, USA

16 ⁵ BioTechnology Institute, College of Biological Sciences, University of Minnesota,
17 Minneapolis, MN 55455, USA

18 ⁶ Division of Gastroenterology and Hepatology, Department of Internal Medicine, Mayo
19 Clinic, Rochester, MN 55902, USA

20 ⁷ Division of Biostatistics, School of Public Health, University of Minnesota, Minneapolis,
21 MN 55455, USA

22 ⁸ Department of Computer Science and Engineering, University of Minnesota,
23 Minneapolis, MN 55455, USA

24

25 **Abstract**

26

27 While the gut microbiome and host gene regulation separately contribute to
28 gastrointestinal disorders, it is unclear how the two may interact to influence host
29 pathophysiology. Here, we developed a machine learning-based framework to jointly
30 analyze host transcriptomic and microbiome profiles from 416 colonic mucosal samples
31 of patients with colorectal cancer, inflammatory bowel disease, and irritable bowel
32 syndrome. We identified potential interactions between gut microbes and host genes
33 that are disease-specific, as well as interactions that are shared across the three
34 diseases, involving host genes and gut microbes previously implicated in
35 gastrointestinal inflammation, gut barrier protection, energy metabolism, and
36 tumorigenesis. In addition, we found that mucosal gut microbes that have been
37 associated with all three diseases, such as *Streptococcus*, interact with different host
38 pathways in each disease, suggesting that similar microbes can affect host
39 pathophysiology in a disease-specific manner through regulation of different host genes.

40

41

42

43

44 Introduction

45

46 The human gut microbiome plays a critical role in modulating human health and
47 disease. Variations in the composition of the human gut microbiome have been
48 associated with a wide variety of chronic diseases, including colorectal cancer (CRC),
49 inflammatory bowel disease (IBD), and irritable bowel syndrome (IBS). For example,
50 previous studies have reported an increase in abundance of *Fusobacterium nucleatum*
51 and *Parvimonas* in CRC ^{1,2}, reduced abundance of *Faecalibacterium prausnitzii* and
52 enrichment of enterotoxigenic *Bacteroides fragilis* in CRC and IBD ³⁻⁵, and
53 overrepresentation of Enterobacteriaceae and *Streptococcus* in IBD and IBS ⁶⁻⁸. In
54 addition to the gut microbiome, dysregulation of host gene expression and pathways
55 have also been implicated in these diseases. Researchers have reported disruption of
56 Notch and WNT signalling pathways in CRC ^{9,10}, activation of toll-like receptors (e.g.
57 TLR4) that induce NF- κ B and TNF- α signaling pathways in IBD ^{11,12}, and dysregulation
58 of immune response and intestinal antibacterial gene expression in IBS ^{8,13}. While host
59 transcription and gut microbiome have separately been identified as contributing factors
60 to these gastrointestinal (GI) diseases, it is unclear how the two may interact to
61 influence host pathophysiology ¹⁴.

62

63 Studies in model organisms have demonstrated that the modulation of host gene
64 expression by the gut microbiome is a potential mechanism by which microbes can
65 affect host physiology ¹⁵⁻²⁰. For example, in zebrafish, the gut microbiome negatively
66 regulates the transcription factor hepatocyte nuclear factor 4, leading to host gene
67 expression profiles associated with human IBD ¹⁸. In mice, the gut microbiota can alter
68 host epigenetic programming to modulate intestinal gene expression involved in
69 immune and metabolic processes ^{16,17}. Additionally, recent *in vitro* cell culture
70 experiments have shown that specific gut microbes can modify the gene expression in
71 interacting human colonic epithelial cells ^{21,22}. Given the evidence for crosstalk between
72 the gut microbiome and host gene regulation, characterizing the interplay between the
73 two factors is critical for unravelling their role in the pathogenesis of human intestinal
74 diseases.

75

76 A few recent studies have investigated interactions between the host transcriptome
77 and gut microbiome in specific human gut disorders, including IBD, CRC, and IBS. For
78 example, studies examining microbiome-host gene relationships in IBD have identified
79 mucosal microbiome associations with host transcripts enriched for
80 immunoinflammatory pathways ²³⁻²⁵. While investigating longitudinal host-microbiome
81 dynamics in IBD, Lloyd-Price et al. identified interactions between expression of
82 chemokine genes, including *CXCL6* and *DUOX2*, and abundance of gut microbes,
83 including *Streptococcus* and Ruminococcaceae ²⁵. Studies investigating the role of host

84 gene-microbiome interactions in CRC have found correlations between the abundance
85 of pathogenic mucosal bacteria and expression of host genes implicated in
86 gastrointestinal inflammation and tumorigenesis^{26,27}. In IBS, host genes implicated in
87 gut barrier function and peptidoglycan binding, such as *KIFC3* and *PGLYRP1*, are
88 associated with microbial abundance of Peptostreptococcaceae and *Intestinibacter*⁸.
89 While these studies have revealed important insights about host gene-microbiome
90 crosstalk in GI diseases, they are limited in several aspects. For example, to boost
91 statistical power, most studies have examined interactions between a limited subset of
92 host genes and gut microbes; for instance, by focusing only on differentially expressed
93 genes^{24,25,27}, genes associated with immune functions^{13,26}, or select microbes
94 representing bacterial clusters or co-abundance groups^{23,26}, thus characterizing only a
95 subset of potential interactions. In addition, the identification of host gene-microbe
96 interactions is based on testing for pairwise correlation between every host gene and
97 microbe using Spearman or Pearson correlation, thus ignoring the inherent multivariate
98 properties of these datasets^{24,25,27}. This approach may also decrease statistical power
99 to detect biologically meaningful associations due to the large number of statistical tests
100 performed. Additionally, most studies focus on examining interactions in a single
101 disease at a time; hence, common and unique patterns of host-microbiome interactions
102 across multiple disease states remain poorly characterized.

103

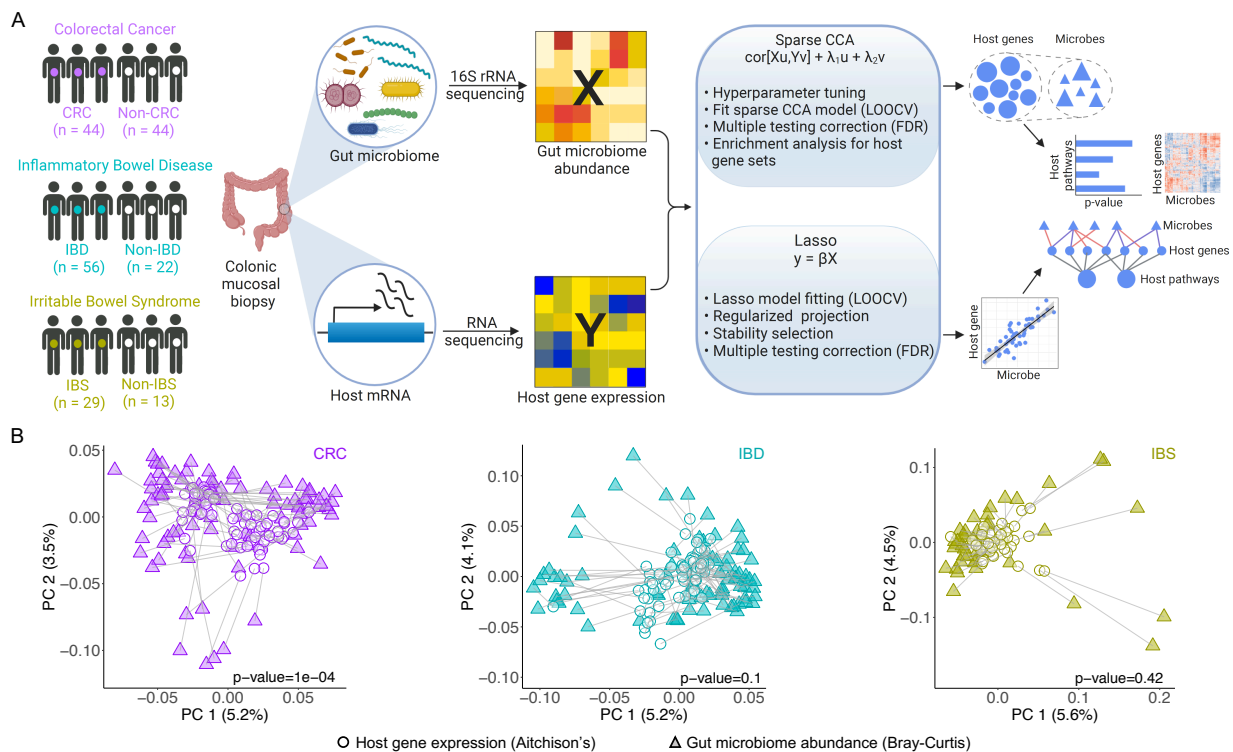
104 Here, we comprehensively characterized interactions between mucosal gene
105 expression and microbiome composition in patients with colorectal cancer, inflammatory
106 bowel disease, and irritable bowel syndrome, three GI disorders in which both host
107 gene regulation and gut microbiome have been implicated as contributing factors
108^{1,6,8,10,13}. We developed and applied a machine learning framework that overcomes
109 typical challenges in multi-omic integrations, including high-dimensionality, sparsity, and
110 multicollinearity, to identify biologically meaningful associations between gut microbes
111 and host genes and pathways in each disease. We leveraged our framework to
112 characterize disease-specific and shared host gene-microbiome interactions across the
113 three diseases that may facilitate new insights into the molecular mechanisms
114 underlying pathophysiology of these gastrointestinal diseases.

115

116 **Results**

117

118 **Integrating host gene expression and gut microbiome abundance in colorectal**
119 **cancer, inflammatory bowel disease and irritable bowel syndrome.**



120
121

122 **Figure 1.** Integrating host gene expression and gut microbiome abundance in colorectal cancer
123 (CRC), inflammatory bowel disease (IBD), and irritable bowel syndrome (IBS). **A.** Study design
124 representing disease cohorts, generation of host gene expression data and gut microbiome
125 abundance data from patient samples, overview of integration framework, and expected output
126 (left to right). For description of mathematical notations, please see Methods. **B.** Procrustes
127 analysis showing overall association between variation in host gene expression and gut
128 microbiome composition in CRC, IBD and IBS (left to right). We used Aitchison's distance for
129 host gene expression data (circle), and Bray-Curtis distance for gut microbiome data (triangle).

130

131 To study host-microbiome relationship across diseases, we used host gene
132 expression (RNA-seq) data and gut microbiome abundance (16S rRNA sequencing)
133 data generated from colonic mucosal biopsies obtained from patients with colorectal
134 cancer (CRC), inflammatory bowel disease (IBD), and irritable bowel syndrome (IBS)
135 (**Figure 1A**). Our study included 208 microbiome samples and 208 paired gene
136 expression samples (416 in total). These 208 paired samples include 88 pairs of
137 samples in the CRC cohort (44 tumor and 44 patient-matched normal), 78 pairs of
138 samples in the IBD cohort (56 patients and 22 controls)^{25,28}, and 42 pairs of samples in
139 the IBS cohort (29 patients and 13 controls; see Supplementary Table S1)⁸. Detailed
140 information on disease cohorts, samples, sequencing, quality control, and data
141 processing is available in Methods.

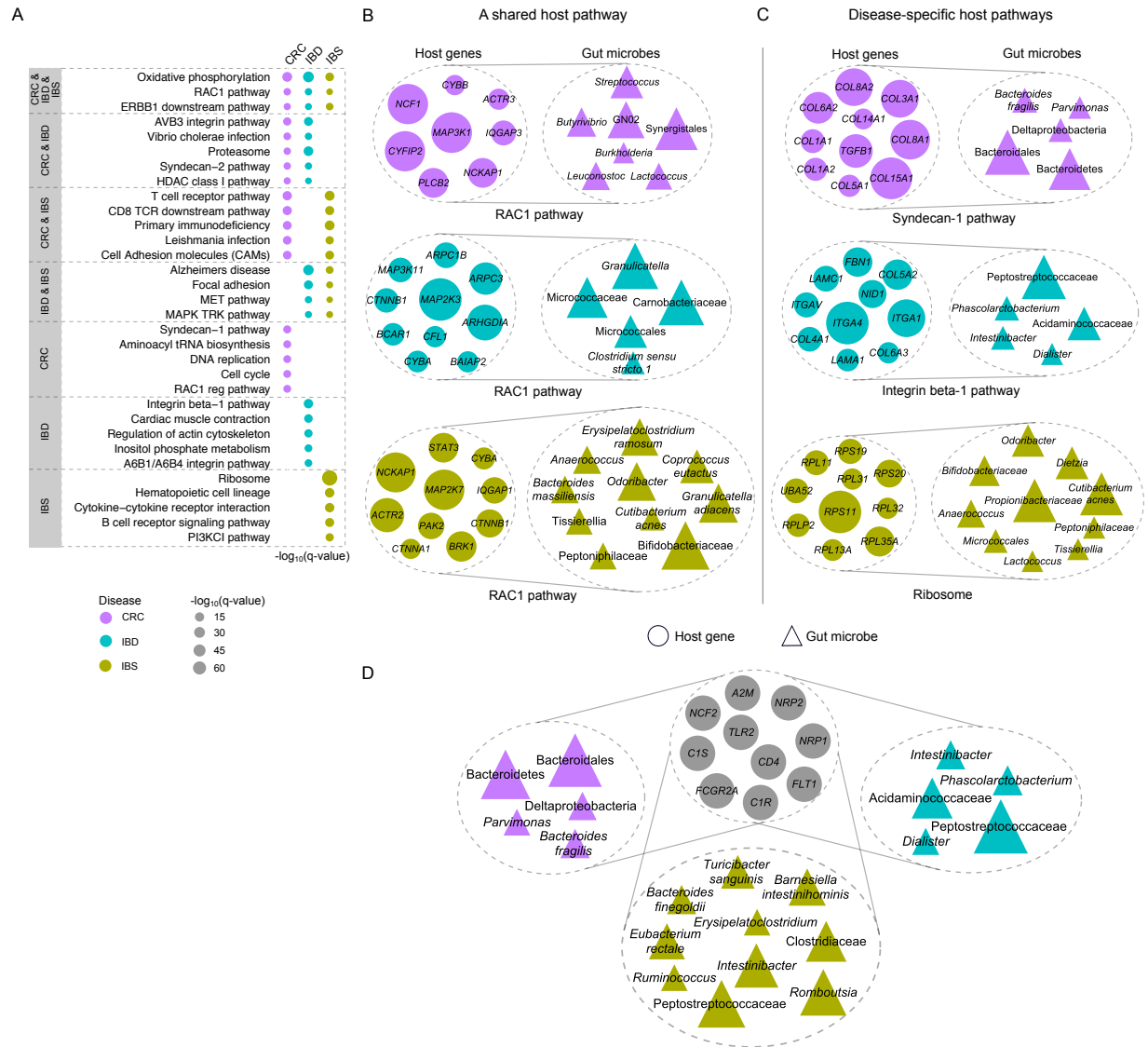
142

143 Previous studies have identified host gene-microbiome associations in human gut
144 disorders, including CRC, IBD, and IBS^{8,25,26}. Thus, one might expect intestinal gene
145 expression patterns and microbiome composition to be broadly correlated in these
146 diseases. To test for such an overall association between host gene expression and gut
147 microbiome composition, we performed Procrustes analysis using paired data for each
148 disease cohort. Our analysis showed significant correspondence between host gene
149 expression variation and gut microbiome composition across subjects in CRC (Monte
150 Carlo p-value = 0.0001). However, Procrustes agreement is not significant in IBD
151 (Monte Carlo p-value = 0.1) and IBS (Monte Carlo p-value = 0.42) (**Figure 1B**, see
152 Methods). This lack of significant overall correspondence between host transcriptome
153 and gut microbiome across diseases might suggest that, instead of an overall
154 association between the two, it is likely that only a subset of gut microbes interact with a
155 subset of host genes at the colonic epithelium^{15,17}. Hence, we need novel integration
156 approaches to characterize such host gene-microbiome interactions.

157
158 To this end, we developed a machine learning framework for integrating multi-omic
159 high-dimensional datasets, such as host gene expression and gut microbiome
160 abundance, to identify relevant host genes and pathways associated with gut microbes.
161 Our integration approach has two parts: (i) Sparse canonical correlation analysis
162 (sparse CCA)^{29,30} for identifying groups of host genes that associate with groups of gut
163 microbial taxa to characterize pathway-level interactions; and (ii) Lasso penalized
164 regression³¹, for identifying specific interactions between individual host genes and gut
165 microbial taxa (see **Figure 1**, Methods, Supplementary Figure S1). We applied our
166 integration analysis to matched host gene expression data and gut microbiome data for
167 each disease cohort separately to avoid any potential batch effects. For each disease
168 cohort dataset, we conducted the integration analysis separately for the patient data
169 (i.e. CRC, IBD, and IBS) and corresponding control data (non-CRC, non-IBD, and non-
170 IBS, respectively), and considered only associations that were found in patients and not
171 in controls. As opposed to the Procrustes analysis, our approach identified significant
172 and potentially biologically meaningful associations between gut microbiota and host
173 genes and pathways across the three diseases.

174

175 **Shared host pathways associate with disease-specific gut microbes across GI**
176 **diseases**



177
 178 **Figure 2.** Shared immunoregulatory and metabolic host pathways associate with disease-
 179 specific gut microbes across human diseases. **A.** Host pathways enriched for sparse CCA gene
 180 sets associated with gut microbiome composition across diseases (FDR < 0.1). Size of the dots
 181 represent the significance of enrichment for each pathway, and color of the dots denote the
 182 disease cohort in which this pathway is significantly associated with microbiome composition. **B.**
 183 Association between microbial taxa in CRC, IBD and IBS (top to bottom) and host genes in the
 184 RAC1 pathway, a shared host pathway (i.e., a pathway for which host gene expression
 185 correlates with gut microbes across disease cohorts). Size of circles and triangles represent the
 186 absolute value of sparse CCA coefficients of genes and microbes, respectively. **C.** Association
 187 between the set of host genes in disease-specific host pathways (i.e. host pathways for which
 188 gene expression correlates with gut microbes in only one of the disease cohorts) and group of
 189 gut bacteria in CRC, IBD and IBS (top to bottom). **D.** A common set of host genes (grey circles)

190 interact with disease-specific sets of microbes. These host genes are enriched for
191 immunoregulatory and acute inflammatory response pathways.

192

193 We hypothesized that host genes and gut microbial taxa involved in common
194 biological functions would act in a coordinated fashion, and, hence, would have
195 correlated expression and abundance patterns. To investigate this, we used sparse
196 CCA to characterize group-level association between host transcriptome and gut
197 microbiome in each of the three diseases^{29,30}. We fit the sparse CCA model for each
198 dataset to identify subsets of significantly correlated host genes and gut microbes,
199 known as components (see Methods, and Supplementary Tables S2–S4). We then
200 performed pathway enrichment analysis on the set of host genes in each significant
201 component to determine host pathways that interact with gut microbes in a disease. We
202 identified *shared* pathways, namely host pathways for which gene expression correlates
203 with gut microbes across disease cohorts, and *disease-specific* pathways, namely host
204 pathways for which gene expression correlates with gut microbes in only one of the
205 three disease cohorts (**Figure 2A**; Fisher's exact test, Benjamini-Hochberg FDR < 0.1,
206 Supplementary Table S5). For simplicity, we focused on the top five most significant
207 shared and disease-specific pathways (**Figure 2A**). We found three pathways shared
208 across CRC, IBD, and IBS that are known to regulate gastrointestinal tract inflammation
209 and gut barrier protection and repair. For example, oxidative phosphorylation, which is
210 the process of energy metabolism in the mitochondria, is known to be upregulated in
211 IBD and CRC, and contributes to tumorigenesis and drug resistance in CRC^{32–35}.
212 Interestingly, the gut microbiome can signal mitochondria of gut mucosal immune cells
213 to alter mitochondrial metabolism, including oxidative phosphorylation processes; this
214 can lead to impaired epithelial barrier function and chronic intestinal inflammation in IBD
215 and CRC³⁶. We also found overlapping host pathways between disease pairs (see
216 CRC & IBD, CRC & IBS, and IBD & IBS in **Figure 2A**), including immunoregulatory
217 pathways and cell-surface receptors like integrin pathway, cell and focal adhesion, and
218 proteasome.

219

220 In addition, we identified 102 disease-specific host pathways that are associated
221 with gut microbes, including 52 CRC-specific, 25 IBD-specific, and 25 IBS-specific
222 pathways (Supplementary Table S5, **Figure 2A**). While IBD-specific host pathways
223 include A6B1/A6B4 integrin pathway and Integrin beta-1 pathway that regulate
224 leukocyte recruitment in GI inflammation^{37,38}, IBS-specific pathways include immune
225 response pathways, including B cell receptor signaling pathway, and ribosome pathway.

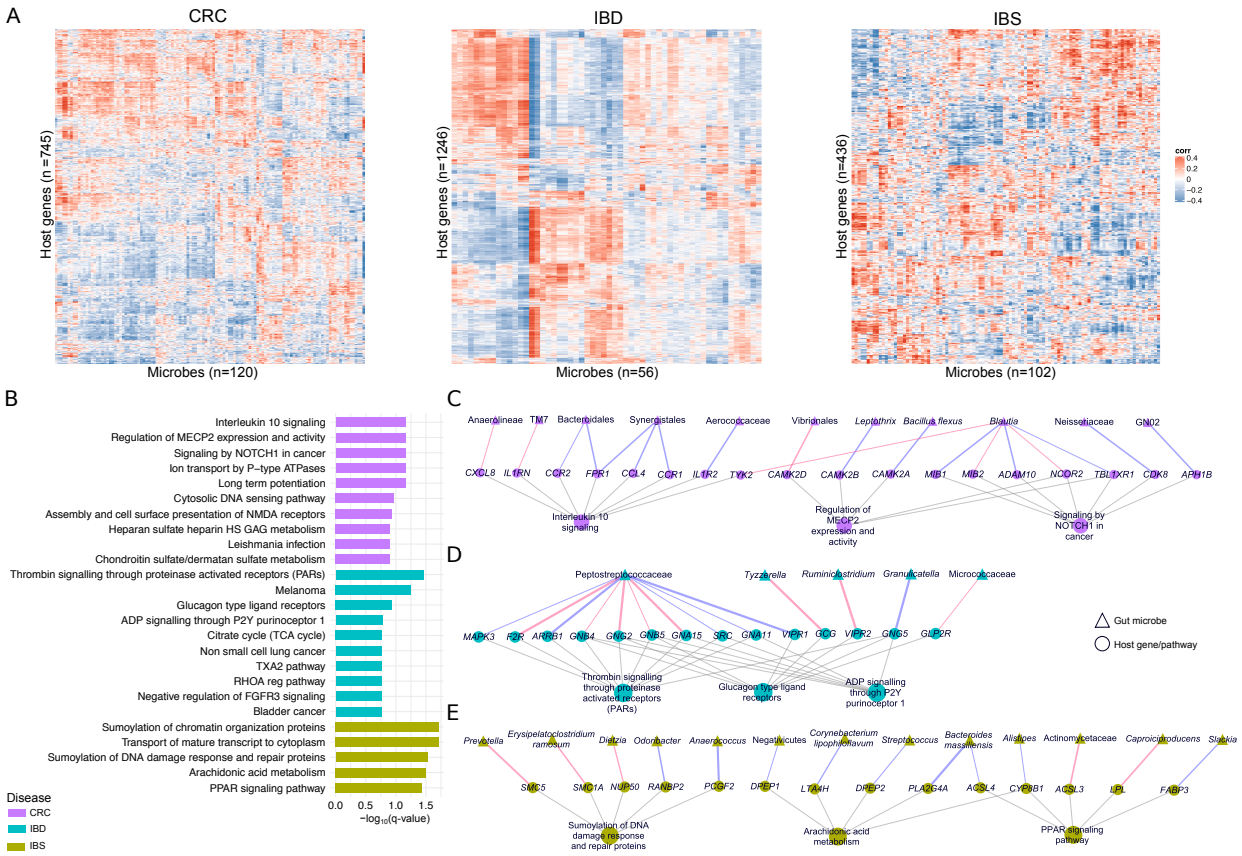
226

227 To better understand the host gene-microbe interactions that underlie common
228 associations, we focused on the RAC1 pathway, where host gene expression is
229 associated with microbiome composition in CRC, IBD, and IBS. The RAC1 pathway is
230 known to regulate immune response and intestinal mucosal repair, and has previously

231 been implicated in IBD and CRC^{39–41} (**Figure 2B**). As expected, we observed
232 overlapping host genes for this shared pathway across the three diseases. However,
233 the microbial taxa they are correlated with are disease-specific. In CRC, the RAC1
234 pathway is associated with oral bacterial taxa such as *Streptococcus*, *Synergistales*,
235 and *GN02*, where *Streptococcus* species are known to be associated with colorectal
236 carcinogenesis^{42,43}. In IBD, the RAC1 host pathway is associated with microbial taxa
237 previously implicated in IBD, including *Granulicatella*^{44–46}, and *Clostridium sensu stricto*
238 *1*, a microbe associated with chronic enteropathy similar to IBD⁴⁷. In IBS, this pathway
239 is associated with bacteria such as *Bacteroides massiliensis*, that has been shown to be
240 prevalent in colitis⁴⁸, and *Bifidobacterium* and *Odoribacter*, that are known to be
241 depleted in IBS^{49–52}.

242
243 To investigate disease-specific associations, we considered unique host pathways
244 for which host gene expression correlates with gut microbes only in one of the three
245 diseases (**Figure 2C**). For example, the Syndecan-1 pathway, which we found to be
246 associated with gut microbial taxa only in CRC, has been previously shown to regulate
247 the tumorigenic activity of cancer cells by altering extracellular matrix adhesion and cell
248 morphology^{53–55}. Host gene expression in this pathway is associated with microbial
249 taxa such as *Parvimonas* and *Bacteroides fragilis* that are known to promote intestinal
250 carcinogenesis and are considered biomarkers of CRC^{1,56–58}. The integrin-1 pathway, a
251 disease-specific host pathway in IBD, is found to be associated with
252 *Peptostreptococcaceae*, *Intestinibacter*, and *Phascolarctobacterium*, microbial taxa that
253 are have been implicated in IBD by previous studies^{59–62}. To assess similarities in host
254 gene components across diseases, we identified a set of host genes that are common
255 between components across the three diseases, and we found that these genes are
256 enriched for immune response pathways in gut epithelium, including vascular
257 endothelial growth factor (VEGF), complementation and coagulation cascades, and
258 cytokine-cytokine receptor interaction (**Figure 2D**; Fisher's exact test, Benjamini-
259 Hochberg FDR < 0.1). While this set of host genes is associated with disease-specific
260 groups of microbes, we also found overlapping microbes between IBD and IBS, such as
261 *Peptostreptococcaceae* and *Intestinibacter*, taxa that are found in high abundance in
262 gastrointestinal inflammation^{59,61,62}.

263
264 **Specific gut microbes interact with individual host genes and pathways in each**
265 **disease**



266

267

268 **Figure 3.** Specific gut microbes interact with individual host genes and pathways in each
 269 disease. **A.** Heatmap showing the overall pattern of interactions between significant and
 270 stability-selected host genes and gut microbial taxa identified by the Lasso model in CRC, IBD,
 271 and IBS (FDR < 0.1). **B.** Host pathways enriched among genes that are correlated with specific
 272 gut microbes in CRC (purple), IBD (green), and IBS (yellow). **C–D.** Networks showing specific
 273 gut microbes correlated with specific host genes enriched for disease-specific host pathways in
 274 CRC (**C**), IBD (**D**), and IBS (**E**). Triangular nodes represent gut microbes, circular nodes
 275 represent host genes and pathways. Edge color represents positive (blue) or negative (red)
 276 association, and edge width represents strength of association (spearman rho). Grey edges
 277 represent host gene - pathway associations.

278

279 Previous studies have shown that specific microbial taxa can regulate expression of
 280 individual host genes^{15,22}. Therefore, we explored interactions between individual host
 281 genes and gut microbes in each disease. To do so, we used Lasso penalized
 282 regression models to identify specific gut microbial taxa whose abundance is associated
 283 with the expression of a host gene³¹. We fit these models in a gene-wise manner, using
 284 expression for each host gene as response and abundance of gut microbial taxa as
 285 predictors. We then applied stability selection to identify robust associations (see
 286 Methods). Using this approach, we found 755, 1295, and 441 significant and stability-
 287 selected host gene-taxa associations in CRC, IBD, and IBS, respectively (**Figure 3**;

288 Tables S6–S8; FDR < 0.1). These represent interactions between 745 host genes and
289 120 gut microbes in CRC (Supplementary Table S6), between 1246 host genes and 56
290 gut microbes in IBD (Supplementary Table S7), and between 436 host genes and 102
291 gut microbes in IBS (Supplementary Table S8) (**Figure 3A**). Examples of specific host
292 gene-microbe interactions can be found in Supplementary Figure S2. Overall, we
293 observed disease-specific patterns in host gene-taxa interactions.

294

295 To characterize the biological functions represented by the host genes that interact
296 with specific gut microbes, we applied enrichment analysis on the set of gut microbiota-
297 associated host genes in each disease (see Methods). This is complementary to our
298 group-level approach (Figure 2) in that these host pathways are enriched among
299 individual host gene-microbe pairs. We identified 87 host pathways that are unique to
300 each disease, including 22 CRC-specific, 60 IBD-specific, and 5 IBS-specific pathways
301 that interact with unique gut bacteria, of which we visualized top 10 most significant host
302 pathways per disease (**Figure 3B**, Fisher's exact test, Benjamini-Hochberg FDR < 0.2,
303 Supplementary Table S9, see Methods). The host pathways enriched for CRC-specific
304 interactions are known to modulate tumor growth, progression and metastasis in CRC,
305 such as Interleukin-10 signaling, signaling by *NOTCH1* in cancer, and regulation of
306 MECP2 expression and activity^{63–66}. The host pathways we identified as enriched for
307 IBD-specific interactions are known to be responsible for maintenance of gastric
308 mucosa integrity, inflammatory response, and host defence against invading pathogens,
309 such as thrombin signalling through proteinase activated receptors (PARs), and
310 glucagon type ligand receptors^{67,68}. For IBS-specific interactions, the enriched host
311 pathways identified here have been shown to regulate homeostasis of intestinal tissue
312 and proinflammatory mechanisms in IBS, such as sumoylation of DNA damage
313 response and repair proteins, and arachidonic acid metabolism^{69–71}.

314

315 To characterize the potential mechanism of host gene-microbe interactions, we
316 further investigated the gut microbial taxa associated with host genes in these pathways
317 (**Figure 3C-E**). In CRC, we found that Anaerolineae and *TM7*, oral microbes that also
318 inhabit the human gastrointestinal tract, and are known to promote oral and colorectal
319 tumorigenesis^{72–76}, are negatively correlated with host genes enriched for tumor-
320 promoting Interleukin-10 signaling pathway, such *CXCL8* and *IL1RN* (**Figure 3C** and
321 Supplementary Figure S2). *CXCL8* is known to be overexpressed in CRC, and *IL1RN* is
322 centrally involved in immune and inflammatory response, and its polymorphisms are
323 implicated in colorectal carcinogenesis^{77–79}. Other host genes in Interleukin-10
324 signaling, such as *CCR2* and *FPR1*, are positively correlated with Bacteroidales (**Figure**
325 **3C** and Supplementary Figure S2). *CCR2* and *FPR1* are overexpressed in colorectal
326 tumors, while Bacteroidales are enriched in CRC and associated with tumorigenesis<sup>80–
327 82</sup>.

328

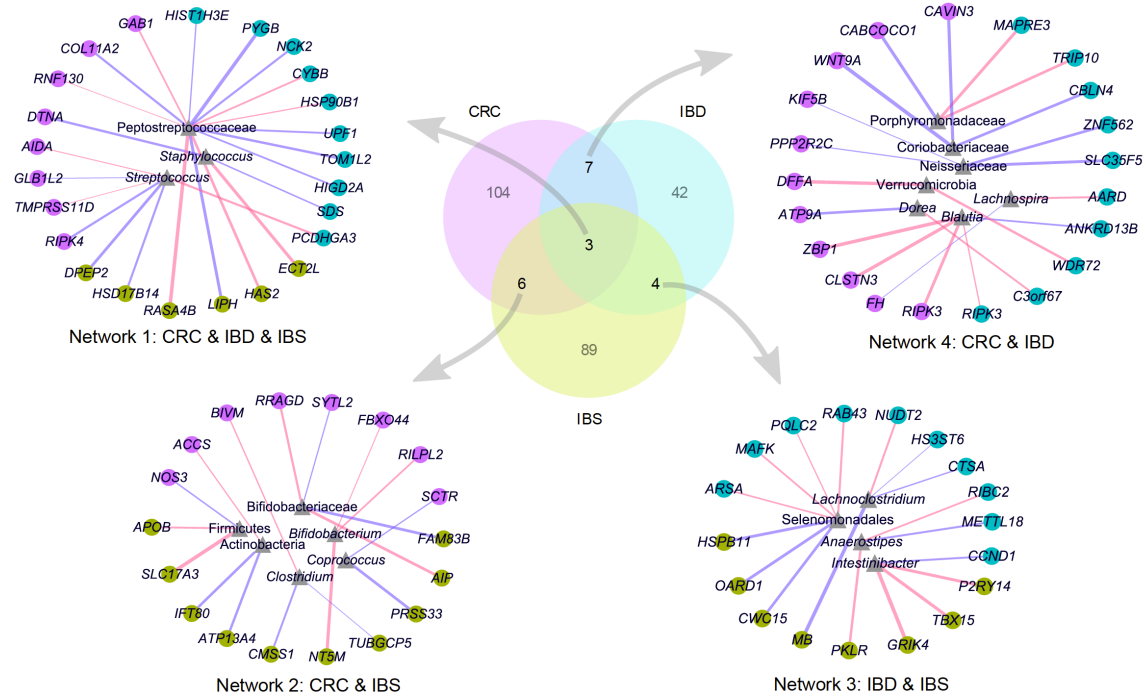
329 We observed that Peptostreptococcaceae, which is prevalent in patients with IBD
330 ^{62,83,84}, is associated with multiple host genes and pathways in IBD (**Figure 3D**). For
331 example, its abundance is positively correlated with the expression of host genes
332 *MAPK3* and *VIPR1*, involved in thrombin signalling through proteinase activated
333 receptors (PARs) and glucagon type ligand receptors pathways, respectively. *MAPK3* is
334 known to play a role in progression and development of IBD, and *VIPR1* is over-
335 expressed in inflamed mucosa ^{85,86}. The abundance of Micrococcaceae, which is known
336 to be increased in IBD, is negatively associated with the expression of *GLP2R*, a
337 glucagon receptor involved in maintenance of gut barrier integrity ^{68,87}. In IBS-specific
338 interactions, we found that the levels of *Prevotella*, which is known to be
339 overrepresented in individuals with loose stool, to be negatively associated with
340 expression of *SMC5*, which is involved in the sumoylation pathway ^{88,89,90} (**Figure 3E**).
341 Previous studies have shown that gut pathogens can target the host sumoylation
342 machinery that regulates inflammatory cascade in epithelial cells in inflammatory bowel
343 disease ⁶⁹. We also found the expression of *PLA2G4A*, a host gene that plays an
344 important role in arachidonic acid metabolism and is an integral member of
345 prostaglandin biosynthesis pathway that modulates gut epithelial homeostasis ^{91,92}, is
346 positively correlated with the abundance of *Bacteroides massiliensis* in IBS, a gut
347 microbe known to be prevalent in patients with gut malignancies, including ulcerative
348 colitis and colorectal carcinoma ^{48,93} (**Figure 3E**). Taken together, these findings
349 demonstrate that interactions between specific gut microbial taxa and specific host
350 genes and pathways vary by disease state.

351

352 **Disease-specific gut microbe-host gene crosstalk**

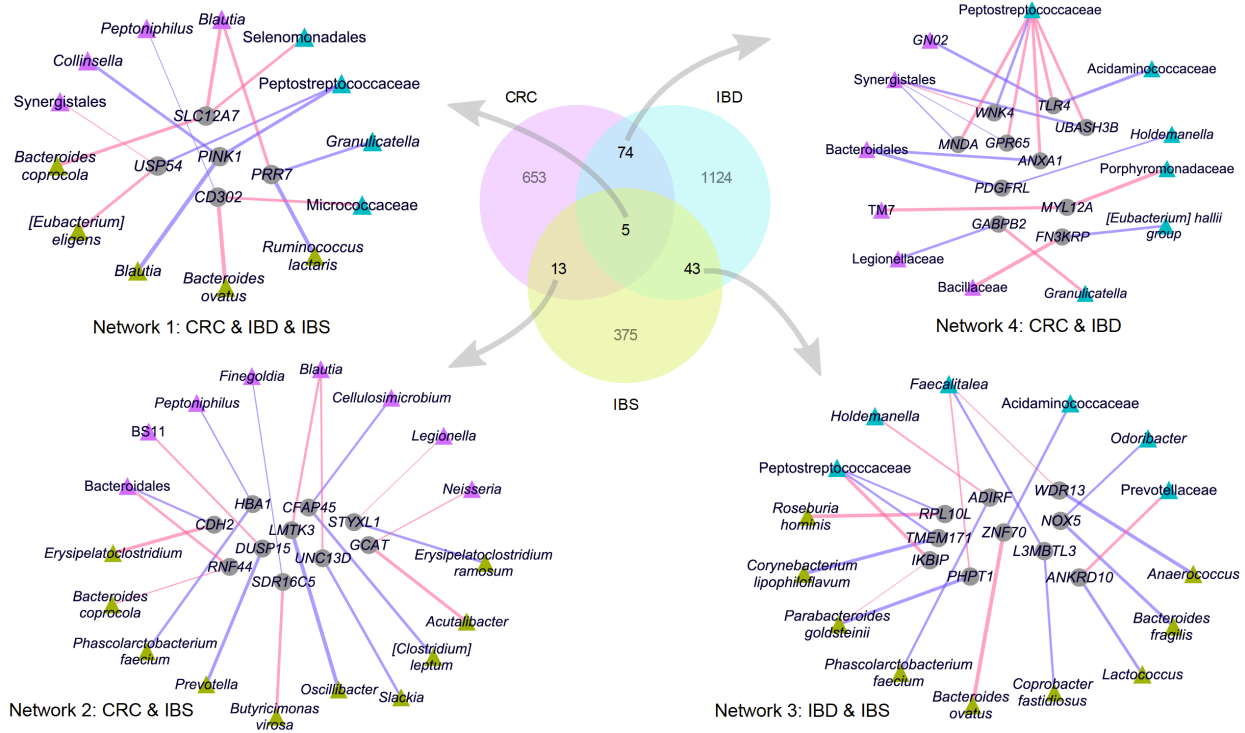
A

Shared gut microbes



B

Shared host genes



353
354

355 **Figure 4.** Associations for shared gut microbes and shared host genes reveal disease-specific
356 host-microbiome crosstalk. **A.** (*center*) Venn diagram showing overlap between gut microbes
357 associated with host genes in CRC, IBD and IBS, (*counter-clockwise*) networks showing host
358 gene-microbe interactions for gut microbes shared across CRC, IBD and IBS (**Network 1**),
359 between CRC and IBS (**Network 2**), between IBD and IBS (**Network 3**), and between CRC and
360 IBD (**Network 4**). **B.** (*center*) Venn diagram showing overlap between host genes associated
361 with gut microbes in CRC, IBD, and IBS, (*counter-clockwise*) networks showing host gene-
362 microbe interactions for host genes shared across CRC, IBD and IBS (**Network 1**), between
363 CRC and IBS (**Network 2**), between IBD and IBS (**Network 3**), and between CRC and IBD
364 (**Network 4**). Circular nodes represent host genes, triangular nodes represent gut microbes.
365 Colored nodes represent specific disease (purple: CRC, green: IBD, yellow: IBS), grey nodes
366 represent gut microbes (A) and host genes (B) shared across diseases. Edge color represents
367 positive (blue) or negative (red) association, and edge width represents strength of association
368 (spearman rho). All interactions were determined at FDR < 0.1.

369
370 To understand how shared gut microbes may interact with specific host genes
371 across diseases, we explored the overlaps between host gene-microbe associations in
372 CRC, IBD, and IBS (**Figure 4A**, Lasso regression, Benjamini-Hochberg FDR < 0.1,
373 Supplementary Table S10). We found that the abundance of 3 gut microbes,
374 Peptostreptococcaceae, *Streptococcus*, and *Staphylococcus*, is correlated with host
375 gene expression in all three diseases (**Figure 4A; Network 1**). Previous studies have
376 revealed that Peptostreptococcaceae and *Streptococcus* spp. are found at elevated
377 levels in CRC, IBD, and IBS^{8,43,62,94–97}. While traditionally considered nasal- or skin-
378 associated bacteria, *Staphylococcus* spp. also colonize the human gastrointestinal tract
379 and include opportunistic pathogens that can cause acute intestinal infections in
380 patients with CRC and IBD^{98–103}, and are associated with increased risk of IBS and
381 CRC^{97,103,104}. We found that the abundance of Peptostreptococcaceae is positively
382 correlated with the expression of host genes *PYGB* and *NCK2* in IBD, whereas it is
383 negatively correlated with the expression of host gene *HAS2* in IBS. *PYGB* and *NCK2*
384 are both upregulated in IBD, where *PYGB* is known to regulate Wnt/ β -catenin pathway,
385 and *NCK2* is involved in integrin and epidermal growth factor receptor signaling^{105–109}.
386 In contrast, *HAS2* is known to have a protective effect on the colonic epithelium through
387 regulation of intestinal homeostasis and inflammation^{110–112}. In CRC, we found that the
388 abundance of Peptostreptococcaceae is negatively associated with the expression of
389 *GAB1*, a host gene for which overexpression stimulates tumor growth in colon cancer
390 cells¹¹³. *Streptococcus* also shows a disease-specific pattern of association with host
391 gene expression. In CRC, its abundance is correlated with the expression of *RIPK4*,
392 which regulates WNT signaling and NF- κ B pathway, and is upregulated in several
393 cancer types, including colon cancer^{114–116}. Similarly, in IBS, *Streptococcus* abundance
394 is correlated with the expression of *DPEP2*, which is known to modulate macrophage
395 inflammatory response¹¹⁷, and is involved in arachidonic acid metabolism that is known
396 to be dysregulated in IBS^{70,71} (**Figure 4A; Network 1**).

397

398 To elucidate potential host gene-microbe interactions for gut microbes shared
399 between diseases, we visualized networks of most significant associations (**Figure 4A;**
400 **Networks 2–4**, Lasso regression, Benjamini-Hochberg FDR < 0.1, Supplementary
401 Table S10). We found 20 microbes for which abundance is associated with the
402 expression of host gene in at least two diseases. Notably, the abundance of *Blautia*, a
403 butyrate-producing beneficial microbe, is found to be negatively correlated with the
404 expression of *RIPK3* in both CRC and IBD (**Figure 4A; Network 4** and Supplementary
405 Figure S3). *RIPK3* promotes intestinal inflammation in IBD, and colon tumorigenesis^{118–}
406¹²². Interestingly, in CRC, *Blautia* is also associated with *ZBP1* (**Figure 4A; Network 4**),
407 a host gene that recruits *RIPK3* to induce NF-κB activation, and regulates innate
408 immune response to mediate host defense against tumors and pathogens^{123–125}.

409

410 Conversely, to explore how shared host genes may interact with gut microbes
411 across all diseases, we identified host genes for which expression is correlated with the
412 abundance of specific gut microbes in CRC, IBD, and IBS (**Figure 4B**, Lasso
413 regression, FDR < 0.1, Supplementary Table S11). We identified 5 such host genes that
414 interact with 4 gut microbes in CRC, 5 gut microbes in IBS, and 4 gut microbes in IBD
415 (**Figure 4B; Network 1**, Supplementary Table S11). Of note, the expression of *PINK1*,
416 a host gene that regulates mitochondrial homeostasis and activates PI3-kinase/AKT
417 signaling, contributing to intestinal inflammation in IBD, and tumorigenesis^{126–128}, is
418 associated with the abundance of *Collinsella* in CRC, Peptostreptococcaceae in IBD,
419 and *Blautia* in IBS. Previous studies have found that *Collinsella* is increased in
420 abundance in CRC, and has been shown to induce inflammation via altering gut
421 permeability^{129–131}. However, *Blautia* has been found to be both positively and
422 negatively correlated with IBS symptoms^{8,51,132}.

423

424 In addition, we identified 135 host genes for which expression is associated with
425 abundance of microbial taxa in at least two of the three diseases, and visualized the
426 network of most significant associations (**Figure 4B; Networks 2–4**, Lasso regression,
427 FDR < 0.1, Supplementary Table S11). We found that the host genes whose expression
428 is correlated with gut microbes in both CRC and IBD are enriched for pathways involved
429 in immune response, including natural killer cell mediated toxicity, *Leishmania* infection,
430 and leukocyte transendothelial migration (**Figure 4B; Network 4**, Fisher's exact test,
431 Benjamini-Hochberg FDR < 0.1). Some notable associations for these shared host
432 genes include host genes and taxa previously implicated in CRC and IBD. For example,
433 expression of Annexin A1 or *ANXA1*, a host gene known to regulate intestinal mucosal
434 injury and repair and found dysregulated in CRC and IBD^{133–135}, is positively correlated
435 with Bacteroidales in CRC, while negatively correlated with Peptostreptococcaceae in
436 IBD (**Figure 4B; Network 4**). Bacteroidales species are known to modulate maturation

437 of the host immune system and gut barrier integrity^{136,137}. *TLR4*, a host gene known to
438 modulate inflammatory response in intestinal epithelium through recognition of bacterial
439 lipopolysaccharide^{138,139}, and previously implicated in IBD and CRC^{140,141}, is found
440 associated with an oral microbe *GN02* in CRC¹⁴², whereas in IBD, it interacts with
441 Acidaminococcaceae, a gut microbe found increased in abundance in patients with
442 Crohn's disease¹⁴³ (**Figure 4B; Network 4**). Overall, our analysis shows that shared
443 gut microbial taxa and shared host genes depict disease-specific host-microbe
444 crosstalk, thus suggesting that the mechanism of host gene-microbiome interaction
445 might be specific to the disease.

446

447 Discussion

448

449 While gut microbial communities and host gene expression have separately been
450 implicated with human health and disease, the role of the interaction between gut
451 microbes and host gene regulation in the pathogenesis of human gastrointestinal
452 diseases remains largely unknown. Here, we comprehensively characterized
453 interactions between gut microbiome composition and host gene expression from 416
454 colonic mucosal samples taken from patients with colorectal cancer, inflammatory bowel
455 disease, and irritable bowel syndrome, in addition to non-disease controls. To overcome
456 the challenges associated with integrating high-dimensional multi-omic datasets, we
457 developed and applied a machine learning framework to characterize interactions
458 between the gut microbiome and host transcriptome in each disease. We identified a
459 common set of host genes and pathways, including pathways that regulate
460 gastrointestinal inflammation, gut barrier protection, and energy metabolism, that are
461 associated with gut microbiome composition in all three diseases. We also found that
462 gut microbes that have been previously associated with all three diseases, including
463 *Streptococcus*, interact with different host pathways in each disease. This suggests that
464 both common and disease-specific interplay between gut microbes and host gene
465 regulation may contribute to the underlying pathophysiology of GI disorders.

466

467 Previous studies have found common microbial signatures across CRC, IBD, and
468 IBS. For example, all three diseases exhibit an overrepresentation of
469 Peptostreptococcaceae and *Streptococcus* spp^{8,43,62,95}. In addition, both CRC and IBD
470 microbiomes are denoted by a loss of butyrate producing gut bacteria, including *Blautia*,
471 and an enrichment of enterotoxigenic *Bacteroides fragilis*^{5,58,95,144}. In contrast to these
472 microbiome similarities, host gene regulation shows distinct alterations across the three
473 GI disorders; for example, unique antibacterial gene expression profile and disruption of
474 purine salvage pathway are specific to IBS, deregulation of proinflammatory IL-23–IL-17
475 signaling is unique to IBD, and prominent activation of oncogenic pathways like Notch
476 and WNT signaling is a hallmark of CRC^{8,13,145,146}. Here, we found that common

477 disease-related gut microbes can interact with host genes and pathways in a disease-
478 specific manner. Thus, it is compelling to hypothesize that although diseases can be
479 characterized by similar microbial perturbations, these microbes can impact different
480 pathophysiological processes through interaction with different host genes in each
481 disease. For example, we found that, in CRC, *Streptococcus* is correlated with the
482 expression of host genes that regulate WNT signaling and NF- κ B pathway, whereas in
483 IBS, *Streptococcus* is correlated with host genes that modulate macrophage
484 inflammatory response, thus suggesting that this gut microbe may perturb distinct host
485 pathways in CRC and IBS. Of course, since our results are based on correlational
486 analysis, it is challenging to assess directionality. While it is possible that these disease-
487 specific interactions have a role in disease pathogenesis, it is also possible that the
488 disease-transformed colonic mucosa renders it more conducive to the same microbial
489 taxa.

490
491 We also identified a common set of host genes and pathways that are associated
492 with gut microbiome composition in all three diseases. These included pathways that
493 regulate gastrointestinal inflammation, immune response, and energy metabolism, and
494 have been previously implicated in these diseases^{33,147–149}. Our analysis shows that
495 these common host genes and pathways correlate with disease-specific gut microbes in
496 CRC, IBD, and IBS. For example, the expression of host gene *PINK1* that regulates the
497 PI3-kinase/AKT signaling pathway¹²⁶ is associated with the abundance of *Collinsella* in
498 CRC, Peptostreptococcaceae in IBD, and *Blautia* in IBS. This suggests that in some
499 cases, distinct gut microbes may modulate host genes and pathways that are commonly
500 dysregulated across different gut pathologies. At the same time, we also found disease-
501 specific host gene-microbe interactions. For example, in CRC, Syndecan-1 pathway, a
502 host pathway that modulates tumor growth and progression, is correlated with microbial
503 taxa such as *Parvimonas* and *Bacteroides fragilis* that are known to promote intestinal
504 carcinogenesis^{53,56,57,63}. These associations are not found in IBD or IBS, and are
505 unique to CRC. Taken together, our results indicate that GI disorders are characterized
506 by a complex network of interactions between microbes and host genes. Although these
507 interactions can be disease-specific, we find cases where the same microbial taxon is
508 associated with different host genes in each disease, and vice-versa: cases where the
509 same host pathway is associated with different microbes in each disease. Although
510 much effort in microbiome research has been directed towards identifying specific
511 microbial taxa that are responsible for the pathogenesis of disease, our findings indicate
512 that without incorporating data on host gene-microbe interactions, studies may be
513 missing the full picture. Host omic data provides invaluable information on the potential
514 mechanisms through which microbes can affect health.

515

516 An important contribution of our work is a machine learning-based integrative
517 framework for characterization of complex host gene-microbe interactions across
518 human diseases. Although few recent studies have investigated associations between
519 host transcriptome and gut microbiome in human gut disorders, our analysis uses an
520 innovative analytical technique that has several advantages^{24–27}. First, as opposed to
521 analyses that rely on calculating pairwise correlations between features (e.g. Dayama et
522 al.), our approach does not require restricting the data to a predetermined subset of
523 taxa or genes of interest to increase statistical power. In addition, compared to
524 Procrustes analysis, which is commonly used for finding overall correspondence
525 between paired datasets, our approach does not only detect overall association, but can
526 also find specific associations between gut microbial taxa and host genes (using Lasso)
527 and pathways (using sparse CCA), allowing identification of specific interactions and
528 shedding light on potential biological mechanisms of interaction. Furthermore, our
529 approach can be applied to other types of multi-omic dataset, including microbial
530 metabolomic and metagenomic data⁸. Lastly, our project incorporates data from across
531 several diseases, identifying commonalities across conditions as well as disease-
532 specific patterns.

533
534 Despite these advantages, our study has several limitations. While we report the
535 potential role of host gene-microbiome interactions in the pathophysiology of GI
536 disorders, our study identifies correlations, and we cannot infer causality here. Given
537 the challenges associated with studying causal mechanisms in humans, future studies
538 using cell culture or animal models would be useful in elucidating the causal role and
539 directionality of interactions between the gut microbiome and host gene regulation in
540 these diseases¹⁵⁰. Another caveat of our study is that it includes three different disease
541 cohorts with disparate sample collection and sequencing protocols, which can lead to
542 potential batch effects. To mitigate this issue, we performed our integration analysis in
543 each disease cohort separately, including cases and internal controls within each
544 cohort. Additionally, our analysis focused only on the taxonomic composition of the
545 microbiome, and hence we could not characterize interactions involving microbial genes
546 and pathways. Lastly, there are several environmental variables that could potentially
547 influence the microbiome, including diet and medication history, which are not available
548 across our disease cohorts.

549
550 Overall, our work demonstrates the power of integrating gut microbiome and host
551 gene expression data to provide insights into their combined role in GI diseases,
552 including CRC, IBD, and IBS. We find disease-specific and shared gut microbe-host
553 gene interactions across these gut disorders, involving gut microbes and host genes
554 implicated in gastrointestinal inflammation, gut barrier protection, and metabolic
555 functions. We also found that the same gut microbes interact with different host genes

556 in different diseases, suggesting potential mechanisms by which similar gut microbes
557 can affect different disease pathologies. These results represent an important step
558 towards characterizing the crosstalk between gut microbiome and host gene regulation
559 and understanding the contribution to disease etiology.

560

561 **Acknowledgements**

562

563 We would like to thank the IBD HMP2 consortium for making the dataset publicly
564 available. We are thankful to the Blekhman Lab members for their comments and
565 suggestions on the manuscript. We thank Dr. Wen Wang and Dr. Gabriel Al-Ghalith for
566 their feedback. This work is supported by NIH grant R35-GM128716 (to R.B.), a
567 University of Minnesota Doctoral Dissertation Fellowship (to S.P), and by NIH grant
568 R01-GM130622 (to E.F.L). This work was carried out, in part, by resources provided by
569 the Minnesota Supercomputing Institute.

570 **Methods**

571

572 **Overall study design, samples and data**

573 We obtained mucosal microbiome (16S rRNA) and host gene expression (RNA-seq)
574 data from colonic mucosal biopsy samples collected from the patients from three
575 disease cohorts: colorectal cancer (CRC), inflammatory bowel disease (IBD), and
576 irritable bowel syndrome (IBS). Except the host gene expression (RNA-seq) data for
577 CRC, all the other datasets were generated and described in detail in previous studies
578 ^{3,8,25}. Below, we describe the sample collection, sequencing, and quality control for host
579 RNA-seq data for CRC cohort, and summarize data acquisition process for other
580 datasets:

581

582 **CRC samples and data**

583 We used 88 pairs of colonic mucosal samples from 44 patients, with primary tumor
584 and normal tissue samples from each individual. These samples were characterized
585 and described in a previous study ³. Detailed cohort characteristics are included in
586 Supplementary Table S1.

587

588 Host RNA-seq sequencing, alignment and quality control. Total RNA was extracted
589 using a previously established protocol ^{3,151}. Approximately 100mg of flash-frozen tissue
590 per sample were lysed by placing the tissue in 1mL of Qiazol lysis reagent (Qiagen Inc.,
591 Valencia, CA, USA) and sonicating in a 65° C water bath for 1-2 hours. Nucleic acids
592 were purified from the lysates using the Qiagen AllPrep DNA/RNA mini kit (Qiagen Inc.,
593 Valencia, CA, USA), quantified using a Nanodrop 2000 spectrophotometer (Thermo
594 Fisher Scientific, Waltham, MA USA), and submitted for RNA sequencing to the
595 University of Minnesota Genomics Center. Total eukaryotic RNA isolates were
596 quantified using a fluorimetric RiboGreen assay, and once the samples passed the
597 initial QC step (≥ 1 microgram and RIN ≥ 8), they were converted to Illumina sequencing
598 libraries using Illumina's TruSeq Stranded Total RNA Library Prep (for details, see
599 www.illumina.com). Truseq libraries were hybridized to a paired-end flow cell and
600 individual fragments were clonally amplified by bridge amplification on the Illumina cBot.
601 Once clustering was complete, the flow cell was loaded on the HiSeq 2500 and
602 sequenced using Illumina's SBS chemistry. Base call (.bcl) files for each cycle of
603 sequencing were generated by Illumina Real Time Analysis (RTA) software. Primary
604 analysis and index de-multiplexing are performed using Illumina's bcl2fastq v2.20.0.422,
605 which output the demultiplexed FASTQ files.

606

607 A quality check of raw sequence FASTQ files was performed using FastQC software
608 (version 0.11.5) ¹⁵². Quality trimming was performed to remove sequence adaptors and
609 low quality bases using Trimmomatic with 3bp sliding window trimming from 3' end

610 requiring minimum Q16 (phred33)¹⁵³. FastQC was run on the resulting trimmed files to
611 ensure good quality of sequences. The paired-end reads were mapped to NCBI v38 *H.*
612 *sapiens* reference genome using HISAT2¹⁵⁴, resulting in an average alignment rate of
613 87.11% overall for 88 samples. We obtained a range of read counts between
614 14,365,657 and 31,530,487 aligned reads per sample, with an average of 22,475,688.2
615 and 22,697,605.5 aligned reads per sample. SAMtools was used for sorting and
616 indexing the aligned bam files. After alignment, the *Subread* package (version 1.4.6)
617 within the *featureCounts* program was used to generate transcript abundance file¹⁵⁵
618 (Supplementary Figure S4).

619

620 16S rRNA data acquisition. The microbiome dataset used in this study was generated
621 and characterized previously³. We used the unnormalized and unfiltered OTU table in
622 tab-delimited format, representing mucosal microbiome data from 44 tumor and 44
623 patient-matched colon tissue samples.

624

625 **IBD samples and data**

626 We used previously generated and described host gene expression (RNA-seq) and
627 mucosal gut microbiome (16S rRNA) data for the IBD cohort generated as part of the
628 HMP2 project^{25,28} (for detailed protocols, see <http://ibdmdb.org/protocols>). These
629 include data from colonic biopsy samples collected from 78 individuals, including 56
630 individuals with IBD, and 22 individuals without IBD (“non-IBD” in HMP2). Out of 56 IBD
631 patients, 34 patients had Crohn’s disease (CD) and 22 patients had ulcerative colitis
632 (UC). Detailed cohort characteristics are included in Supplementary Table S1. We
633 downloaded metadata, host RNA-seq data, and microbiome data for these samples
634 from <http://ibdmdb.org> in July 2018. We downloaded the unnormalized and unfiltered
635 OTU table and host transcript read counts files in tab-delimited format. We describe the
636 filtering and preprocessing steps for host gene expression and microbiome data below.

637

638 **IBS samples and data**

639 We used previously generated and characterized host gene expression (RNA-seq)
640 and mucosal gut microbiome (16S rRNA) data for the IBS cohort⁸. These include data
641 from colonic biopsy samples collected from 42 individuals, including 29 individuals with
642 IBS, and 13 healthy individuals (non-IBS). Detailed cohort characteristics are included in
643 Supplementary Table S1. We obtained the unnormalized and unfiltered OTU table and
644 host transcript read count files in tab-delimited format via personal communication with
645 authors of the paper⁸. For some individuals, samples were collected at two time points.
646 For these cases, we averaged the gene expression levels and microbiome abundance
647 measurements across the two time points. This is supported by a recent study showing
648 that “omics” methods are more accurate when using averages over multiple sampling

649 time points¹⁵⁶. We describe the filtering and preprocessing steps for host gene
650 expression and microbiome data below.

651

652 **Preprocessing host gene expression data**

653 For host gene expression data for each disease cohort, we used *biomaRt* R
654 package (version 2.37.4) to only keep data for protein-coding genes¹⁵⁷. We filtered out
655 low expressed genes to retain genes that are expressed in at least half of the samples
656 in each disease cohort. We performed variance stabilizing transformation using the R
657 package *DESeq2* (version 1.14.1) on the filtered gene expression read count data¹⁵⁸.
658 We filtered out genes with low variance, using 25% quantile of variance across samples
659 in each disease cohort as cutoff. Performing these steps for RNA-seq data for each
660 disease cohort separately resulted in a unique host gene expression matrix per disease
661 for downstream analysis, including 12513 genes in the CRC dataset, 11985 genes in
662 IBD dataset, and 12429 genes in IBS dataset.

663

664 **Preprocessing microbiome data**

665 We performed the following steps for microbiome data from each disease cohort
666 separately. First, sequences that were classified as either having originated from
667 archaea, chloroplasts, known contaminants originating from laboratory reagents or kits,
668 and soil or water-associated environmental contaminants were removed from the OTU
669 table as described previously¹⁵⁹. Next, we summarized the OTU table at the species (if
670 present), genus, family, order, class, and phylum taxonomic levels, and performed
671 prevalence and abundance-based filtering to retain taxa found at 0.001 relative
672 abundance in at least 10% of the samples. We then concatenated these summarized
673 taxa matrices (count data) into one combined taxa matrix for each disease dataset. We
674 applied centered log ratio (CLR) transform on the filtered taxa count matrix to account
675 for compositionality effects. These steps resulted in a taxonomic abundance matrix for
676 each disease cohort, which included 235 taxa in the CRC dataset, 121 taxa in the IBD
677 dataset, and 238 taxa in the IBS dataset.

678

679 **Procrustes analysis**

680 To assess overall correspondence between host gene regulation and gut
681 microbiome composition in CRC, IBD, and IBS, we performed Procrustes analysis in R
682 using the *vegan* package (version 2.4-5)¹⁶⁰. For each disease cohort, we used
683 Aitchison's distance on host gene expression data and Bray Curtis distance on gut
684 microbiome data as input to the Procrustes analysis¹⁶¹. The significance of rotation
685 agreement was obtained using the *protest()* function with 9,999 permutations.

686

687 **Sparse Canonical Correlation Analysis**

688 We used sparse canonical correlation analysis (sparse CCA) to identify group-level
689 correlations between paired host gene expression and gut microbiome data in each
690 disease cohort. Canonical correlation analysis (CCA) identifies linear projection of two
691 sets of observations into shared latent space that maximizes correlation between the
692 two datasets¹⁶². Sparse CCA is adapted from CCA for high dimensional settings to
693 incorporate feature selection by utilizing $L1$ or lasso penalty in CCA²⁹. The objective
694 function of sparse CCA can be expressed as follows:

$$695 \quad \text{maximize}_{u,v} u^T X^T Y v \text{ subject to } u^T X^T X u \leq 1, v^T Y^T Y v \leq 1, \|u\|_1 \leq \lambda_1, \|v\|_1 \leq \lambda_2$$

696
697 where, X and Y denote two data matrices with same number of samples, but different
698 number of features (representing gut microbiome taxonomic composition data and host
699 gene expression data, respectively); u and v are canonical loading vectors of X and Y
700 respectively; λ_1 and λ_2 control lasso penalties of u and v , respectively.
701
702

703 For each disease cohort separately, we applied sparse CCA using R (version 3.3.3)
704 package *PMA* (version 1.1) with gut microbiome taxonomic composition and host gene
705 expression as two sets of variables to be correlated¹⁶³. Below, we describe details on
706 hyperparameter tuning, fitting sparse CCA models, computing significance of correlation
707 for sparse CCA components, enrichment analysis, and visualization of sparse CCA
708 output.
709

710 **Hyperparameter tuning and fitting for Sparse CCA model**

711 We performed hyperparameter tuning to identify the sparsity penalty parameters for
712 gut microbiome abundance (λ_1) and host gene expression (λ_2) data. Since the
713 permutation search provided in the *PMA* package only performs a one-dimensional
714 search in the tuning parameter space, we implemented a grid-search approach using
715 leave-one-out cross-validation in R (version 3.3.3) for hyperparameter tuning. We
716 selected penalty parameters which had the highest correlation under cross-validation.
717 Using this approach, we identified λ_1 as 0.15 and λ_2 as 0.2 for CRC data, λ_1 as 0.177
718 and λ_2 as 0.333 for IBD data, λ_1 as 0.4 and λ_2 as 0.1 for IBS data.
719

720 After identifying sparsity parameters, we fit the sparse CCA model to obtain subsets
721 of correlated host genes and gut microbes, known as components. Each sparse CCA
722 component includes non-zero weights (or canonical loadings) on gut microbes, and
723 non-zero weights on a subset of host genes correlated with those gut microbes to
724 capture joint variation in the two sets of observations. We computed the first 10 sparse
725 CCA components for each disease cohort, performing a separate computation for case
726 and control samples. Sparse CCA components are computed iteratively, informed by

727 previously computed components, thus, resulting in uncorrelated components¹⁶⁴. Next,
728 we assessed the significance of sparse CCA components as described below.

729

730 **Significance of correlation for sparse CCA components**

731 We computed the significance of each pair of canonical variables (or a component)
732 using leave-one-out cross-validation approach in R (version 3.3.3). For a given
733 component, we first used the penalty parameters determined above to compute the
734 sparse CCA output with one sample held out. We then computed the scores for the
735 held-out sample, i.e., we computed $scoreX_i = X_i u_{-i}$ and $scoreY_i = Y_i v_{-i}$, where i is
736 the held-out sample, X_i and Y_i denote the values for i^{th} sample of the input data
737 matrices X and Y , and u_{-i} and v_{-i} are the canonical loadings estimated from the
738 sparse CCA computation without the i^{th} sample. We repeated this n times, where n is
739 the total number of samples in the data, to obtain the vector of held-out scores. To
740 assess the true strength of association and its significance, we used $cor.test()$ on the
741 scores computed for the held-out samples. We corrected the p-values for multiple
742 hypothesis testing using Benjamini-Hochberg (FDR) method within each disease cohort,
743 and determined significant components at $FDR < 0.1$.

744

745 Using this approach, we identified 7 significant components in CRC, with an average
746 of 828 host genes and 8 gut microbes; 4 significant components in IBD, with an average
747 of 2095 host genes and 6 gut microbes; and 6 significant components in IBS, with an
748 average of 577 host genes and 61 gut microbes ($FDR < 0.1$, Supplementary Tables S2-
749 S4).

750

751 **Enrichment analysis for sparse CCA**

752 To characterize host pathways enriched for the set of host genes associated with
753 microbes in each component, we implemented an enrichment analysis in R (version
754 3.3.3). We implemented Fisher's exact test to assess pathway enrichment, where we
755 used the set of host genes input to the sparse CCA analysis as background genes, and
756 set of host genes in a component as the genes of interest. We used KEGG and PID
757 gene sets from MsigDB canonical pathways collection^{165,166}. To avoid pathways that
758 are too large to provide any specific biological insights or too small to provide adequate
759 statistical power, we excluded any pathway with either (1) fewer than 25 genes, (2)
760 more than 300 genes, or (3) fewer than 5 genes that overlapped between the genes of
761 interest and the pathway. We combined the set of enriched host pathways for all
762 significant components for a given disease dataset, corrected for multiple hypothesis
763 testing within each disease cohort using Benjamini-Hochberg (FDR) approach, and
764 determined significant host pathways at $FDR < 0.1$. This analysis was performed
765 separately for case and control data for each disease.

766

767 To identify case-specific host pathways, we used a two-part approach: (1) first, we
768 identified pathways that are only significantly enriched in cases (FDR < 0.1) and not in
769 controls. (2) In addition, we identified pathways that are significant in both the cases and
770 controls at FDR < 0.1. For these pathways, we performed differential enrichment in
771 cases versus controls by implementing a comparative log odds-ratio approach in R
772 ^{167,168}. To do so, we first computed the z-score for the odds ratio for i -th pathway in
773 cases:

$$774 z_{i,case} = \log(\delta_i)/SE(\delta_i)$$

776 where, δ_i is the odds-ratio for i -th pathway in cases, and $SE(\delta_i)$ is the standard error for
777 i -th pathway in cases, which is computed using the four elements, n_1 to n_4 , of the 2x2
778 contingency table used in the enrichment analysis for the i -th pathway as follows:
779

$$780 SE(\delta_i) = \sqrt{1/n_1 + 1/n_2 + 1/n_3 + 1/n_4}$$

782 Similarly, we computed $z_{i,ctrl}$ for the same pathway in the controls. For a given
783 pathway, we compare enrichment for the component that gives highest significance for
784 the cases with that which gives highest significance for the controls. Next, we compute
785 a comparative log odds-ratio for i -th pathway overlapping between cases and controls
786 as follows:
787

$$788 z_{i,case-ctrl} = \frac{\log(\delta_{i,case}) - \log(\delta_{i,ctrl})}{SE(\delta_{i,case,ctrl})}$$

790 The greater the value of $z_{i,case-ctrl}$, the greater the odds a pathway is differentially
791 enriched in case versus control than by chance. P-values were inferred assuming
792 normal approximations, and corrected for multiple hypothesis testing using Benjamini-
793 Hochberg (FDR) approach. As the last step of part (2), we retained pathways that were
794 differentially enriched in cases versus controls at FDR < 0.2. Finally, we combined the
795 pathways from part (1) and (2) to obtain case-specific pathways.
796

797 **Visualizing disease-specific and shared host pathways and components from** 798 **sparse CCA**

801 To determine *shared* host pathways, i.e. host pathways for which gene expression
802 correlates with gut microbes across all three disease cohorts, and *disease-specific* host
803 pathways, i.e. host pathways for which gene expression correlates with gut microbes in
804 only one of the three disease cohorts, we computed overlaps between significant case-
805 specific host pathways determined above across the three disease cohorts. Given the

806 overlap across the curated gene sets from MsigDB, we controlled for redundancy
807 across pathways for visualization purposes. To do this, we identified similar pathways
808 based on their relative overlap in terms of the set of genes using an overlap coefficient.
809 The overlap coefficient between two pathways is defined as the number of common
810 genes between the pathways divided by the number of genes in the pathway with fewer
811 genes. Specifically, the overlap coefficient is represented as follows:
812

$$813 \quad \text{overlap}(X, Y) = \frac{|X \cap Y|}{\min(|X|, |Y|)}$$

814
815 For the top 15 most significant host pathways (FDR < 0.1) discovered for each
816 shared and disease-specific set (Supplementary Table S12), we computed pairwise
817 similarity between pathways as overlap coefficients and used a maximum allowed
818 similarity score of 0.5 as a cutoff. Using the pairs of pathways that satisfied the cutoff,
819 we computed the connected components to identify clusters of overlapping pathways.
820 For visualization purposes, we selected a representative pathway from each connected
821 component, prioritising the pathway with the highest number of genes (Figure 2A,
822 Supplementary Table S4). We visualized host pathway enrichment using the R package
823 *ggplot* (version 3.2.1).
824

825 For visualizing components corresponding to selected host pathways or common
826 host genes across diseases (Figures 2B – D), we ordered host genes and taxa by their
827 absolute coefficients in the component, and selected the top 10 host genes and taxa for
828 representation. If multiple taxa originating from the same lineage occurred in a
829 component, we selected the one with the highest coefficient to reduce redundancy, thus
830 representing the taxa with most contribution from a given lineage. The size of host
831 genes and gut microbial taxa are scaled by the absolute value of their corresponding
832 coefficients in a given component. All sparse CCA components were visualized using
833 Cytoscape (version 3.5.1) ¹⁶⁹.
834

835 **Lasso regression analysis**

836 We used Lasso penalized regression to identify specific interactions between
837 individual host genes and gut microbial taxa within each disease cohort ³¹. We
838 implemented a gene-wise model using expression for each host gene as response and
839 abundances of microbiome taxa as predictors, to identify microbial taxa that are
840 correlated with a host gene. An ordinary least squares (OLS) regression is not suitable
841 for this task, since OLS results in unstable solutions under high-dimensional settings or
842 when $p \gg n$, i.e. number of predictors p is much larger than number of samples n .
843 Additionally, we expect the abundance of only a few microbial taxa to correlate with the
844 expression of each host gene. To address this, we used lasso regression, which is

845 similar to multivariate OLS, except that it uses shrinkage or regularization to perform
846 variable selection, thus picking only a few taxa that associate with a host gene's
847 expression.

848

849 To account for other factors that can impact host gene expression or microbiome
850 composition, each model also included covariates in the predictor matrix (i.e.
851 microbiome abundance table) for gender (male or female), disease-subtype for IBD
852 (Crohn's Disease or ulcerative colitis), disease-subtype for IBS (constipation (IBS-C) or
853 diarrhea (IBS-D)).

854

855 The lasso model estimates the lasso regression coefficient $\hat{\beta}$ by minimizing the
856 following:

$$857 \quad \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j|$$

858 where, n = number of samples; p = number of predictors (taxa and other covariates);
859 $1 \leq i \leq n, 1 \leq j \leq p$; y = response (host gene expression); x = predictor (taxa
860 abundance and other covariates); λ = tuning parameter, $\lambda \geq 0$.

861

862 In addition to minimizing the residual sum of squares (first term in the equation),
863 lasso minimizes the l_1 norm of the coefficients (second term in the equation), which has
864 an effect of forcing some of the coefficients to zero as the value of tuning parameter, λ ,
865 increases. Thus, lasso performs feature or variable selection that leads to sparse
866 models.

867

868 We implemented a lasso regression framework using R (version 3.3.3) package
869 `glmnet` (version 2.0-13), which uses cyclical coordinate descent to compute
870 regularization path¹⁷⁰. Our framework executes a lasso regression for each host gene's
871 expression as response and abundances of microbial taxa and values of other
872 covariates as predictors. We used leave-one-out cross-validation to estimate the tuning
873 parameter, λ , which was used to fit the final model on a given disease dataset.

874

875 We then performed inference for the lasso model using a regularized projection
876 approach known as desparsified lasso. The desparsified lasso uses the asymptotic
877 normality of a bias-corrected version of the lasso estimator to obtain 95% confidence
878 intervals and p-values for the coefficient of each predictor (microbe) associated with a
879 given host gene¹⁷¹. We used the R package `hdi` (version 0.1-7) that implements the
880 desparsified lasso approach for estimation of confidence intervals and hypothesis
881 testing in high dimensional and sparse settings^{171,172}. We then corrected for multiple
882 hypothesis testing using Benjamini-Hochberg (FDR) method.

883

884 **Stability selection for Lasso model**

885 Since the lasso model is sensitive to small variations of the predictor variable, we
886 used stability selection to pick out robust microbes associated with a host gene ¹⁷³.

887 Stability selection is a resampling-based method that can be combined with different
888 variable selection procedures in high dimensional settings, including lasso. Briefly,
889 stability selection with lasso proceeds as follows:

890

891 Step 1. Select a random subset of the data.

892 Step 2. Fit the lasso model with a randomly perturbed penalty term in the
893 neighborhood of the “best” penalty λ . Record the set of selected variables (microbes).

894 Step 3. Repeat steps 1) and 2) K times.

895 Step 4. Compute the frequency of selection, f_i per variable (microbe) across all trials.

896 Step 5. Select the variables (microbes) that are selected with a frequency of at least
897 f_{thr} , a pre-specified threshold value. Thus, we select a set of stable variables
898 (microbes) such that $f_i \geq f_{thr}$.

899

900 The overall idea is that, if the same variables (microbes) are repeatedly selected when
901 the parameters are perturbed, then they are robust variables. Stability selection also
902 controls for family-wise error rate, thus controlling for false positives in addition to the
903 FDR approach mentioned above ¹⁷³. In our analysis, we used the R package *stabs*
904 (version 0.6-3) to perform stability selection ¹⁷⁴. Specifically, we used the following
905 parameters in the process described above: in Step 1, a random subset of size $n/2$ of
906 data is selected, where n is total number of samples, in Step 3, $K = 100$, and in Step 5,
907 $f_{thr} = 0.6$, i.e. a predictor (microbe) selected in at least 60% of the fitted models is
908 considered stable. The choice of these parameters are in accordance with the proposal
909 of stability selection by Meinshausen and Bühlmann ¹⁷³.

910

911 Finally, we performed an intersection between associations identified by stability
912 selection here and associations identified at $FDR < 0.1$ by the lasso model described
913 above. We removed any host gene-gender and host gene-disease-subtype associations
914 to obtain the significant and stability selected host gene-microbe associations at $FDR <$
915 0.1 .

916

917 **Parallel execution of lasso analysis on supercomputing nodes**

918 We implemented a parallel framework for executing the gene-wise lasso analysis,
919 where we parallelized execution of lasso models on host genes across multiple nodes
920 and cores on a compute cluster from Minnesota Supercomputing Institute. We used job
921 arrays to parallelize our analysis on multiple nodes on the cluster. Additionally, we used

922 R packages *doParallel* (version 1.0.15) and *foreach* (version 1.4.7) to run parallel
923 processes on multiple cores of each compute node.

924

925 **Enrichment analysis for lasso output**

926 To characterize biological functions for the host genes that were found associated
927 with specific gut microbes in a disease cohort by the lasso framework, we implemented
928 an enrichment analysis in R (version 3.3.3) using Fisher's exact test. We used the set of
929 expressed genes input to the lasso analysis as the background genes, and the set of
930 host genes associated with gut microbes in a patient samples as genes of interest. We
931 used the KEGG, PID, and REACTOME gene sets from MsigDB canonical pathways
932 collection ^{165,166}. To avoid too large or too small pathways, we excluded from our
933 analysis any pathway with fewer than 25 genes, greater than 85 genes, or fewer than 5
934 genes that overlap between the genes of interest and the pathway. The p-values
935 obtained from Fisher's exact test were adjusted for multiple testing using Benjamini-
936 Hochberg (FDR) approach. We identified 87 host pathways that are unique to each
937 disease, including 22 CRC-specific, 60 IBD-specific, and 5 IBS-specific pathways that
938 interact with unique gut bacteria (FDR < 0.2, Supplementary Table S8). Here, we used
939 a more relaxed FDR threshold of 0.2 to present a larger number of biologically relevant
940 host pathways.

941

942 **Visualization of shared genes and taxa interactions for lasso output**

943 In Figure 4, we visualized host gene-microbe interactions for gut microbes and host
944 genes shared across diseases. For visualizing interactions for shared gut microbes
945 (Figure 4A), we identified shared microbes between all possible overlaps between
946 diseases (Figure 4A; Networks 1–4), host genes that interact with these common
947 microbes in each disease, and all interactions involving these microbes and genes in
948 each disease (FDR < 0.1). Next, we grouped gene-taxa interactions identified per
949 disease by shared taxa, sorted them by FDR value, and picked top gene-taxa
950 association per shared taxa until we obtained at most 10 interactions per disease (FDR
951 < 0.1, Supplementary Table S9).

952

953 Similarly, for visualizing interactions for shared host genes (Figure 4B), we identified
954 shared host genes between all possible overlaps between diseases (Figure 4B;
955 Networks 1–4), and host gene-taxa interactions per disease for these host genes
956 shared across diseases. We sorted the interactions by FDR adjusted q-values (ordered
957 first by q-value in CRC associations, followed by q-value in IBD associations, and finally
958 by q-value in IBS associations, depending on the overlapping set under consideration).
959 We picked top 10 genes from this merged output, and identified at most top 10
960 associations involving these genes in each disease for the overlapping set under
961 consideration (FDR < 0.1, Supplementary Table S10). Since lasso gives biased

962 estimates of the coefficients, we used Spearman correlation coefficient (ρ) to depict
963 strength of association for visualizing host gene-taxa associations. All the associations
964 in Figure 4 were visualized using Cytoscape v3.5.1, where shared features are in grey
965 and disease-specific features in disease-specific colors ¹⁶⁹.

966

967 **Data and Software Availability**

968

969 Raw data for host RNA-seq for CRC cohort is available on the NCBI Sequence
970 Read Archive (SRA) under submission ID: SUB9143781. For raw data from 16S rRNA
971 sequencing for CRC cohort, RNA-seq and 16S rRNA sequencing for IBD cohort, and
972 RNA-seq and 16S rRNA sequencing for IBS cohort, please refer to data accession
973 details published previously ^{8,25,28}. Processed data tables for host transcriptomics and
974 microbiome data for each disease cohort have been included as supplemental tables
975 (Supplementary Tables S13–S18). Code used for integration analyses performed in the
976 paper is available at:

977 https://github.com/blekhmanlab/host_gene_microbiome_interactions

978

979 **Ethics statement**

980

981 For the colorectal cancer cohort, all research conformed to the Helsinki Declaration
982 and was approved by the University of Minnesota Institutional Review Board, protocol
983 1310E44403. For the inflammatory bowel disease and irritable bowel syndrome cohorts,
984 ethical approval is described in their respective publications ^{8,25,28}.

985

986 **Declaration of Interests**

987

988 D.K. serves as Senior Scientific Advisor to Diversigen, a company involved in the
989 commercialization of microbiome analysis. P.C.K. is an ad hoc consultant for Otsuka
990 Pharmaceuticals, Pendulum Therapeutics, IP group inc. and Novome Biotechnologies.

991

992 **References**

- 993 1. Sobhani, I. *et al.* Colorectal cancer-associated microbiota contributes to oncogenic
994 epigenetic signatures. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 24285–24295 (2019).
- 995 2. Castellarin, M. *et al.* Fusobacterium nucleatum infection is prevalent in human
996 colorectal carcinoma. *Genome Res.* **22**, 299–306 (2012).
- 997 3. Burns, M. B., Lynch, J., Starr, T. K., Knights, D. & Blekhman, R. Virulence genes

- 998 are a signature of the microbiome in the colorectal tumor microenvironment.
999 *Genome Med.* **7**, 55 (2015).
- 1000 4. Irrazábal, T., Belcheva, A., Girardin, S. E., Martin, A. & Philpott, D. J. The
1001 multifaceted role of the intestinal microbiota in colon cancer. *Mol. Cell* **54**, 309–320
1002 (2014).
- 1003 5. Swidsinski, A., Weber, J., Loening-Baucke, V., Hale, L. P. & Lochs, H. Spatial
1004 Organization and Composition of the Mucosal Flora in Patients with Inflammatory
1005 Bowel Disease. *Journal of Clinical Microbiology* vol. 43 3380–3389 (2005).
- 1006 6. McIlroy, J., Ianiro, G., Mukhopadhyay, I., Hansen, R. & Hold, G. L. Review article:
1007 the gut microbiome in inflammatory bowel disease-avenues for microbial
1008 management. *Aliment. Pharmacol. Ther.* **47**, 26–42 (2018).
- 1009 7. Distrutti, E., Monaldi, L., Ricci, P. & Fiorucci, S. Gut microbiota role in irritable
1010 bowel syndrome: New therapeutic strategies. *World J. Gastroenterol.* **22**, 2219–
1011 2241 (2016).
- 1012 8. Mars, R. A. T. *et al.* Longitudinal Multi-omics Reveals Subset-Specific Mechanisms
1013 Underlying Irritable Bowel Syndrome. *Cell* **0**, (2020).
- 1014 9. Schatoff, E. M., Leach, B. I. & Dow, L. E. Wnt Signaling and Colorectal Cancer.
1015 *Curr. Colorectal Cancer Rep.* **13**, 101–110 (2017).
- 1016 10. Koveitypour, Z. *et al.* Signaling pathways involved in colorectal cancer progression.
1017 *Cell Biosci.* **9**, 97 (2019).
- 1018 11. Pedersen, J., Coskun, M., Soendergaard, C., Salem, M. & Nielsen, O. H.
1019 Inflammatory pathways of importance for management of inflammatory bowel
1020 disease. *World J. Gastroenterol.* **20**, 64–77 (2014).

- 1021 12. Khor, B., Gardet, A. & Xavier, R. J. Genetics and pathogenesis of inflammatory
1022 bowel disease. *Nature* **474**, 307–317 (2011).
- 1023 13. Bennet, S. M. P. *et al.* Altered intestinal antibacterial gene expression response
1024 profile in irritable bowel syndrome is linked to bacterial composition and immune
1025 activation. *Neurogastroenterol. Motil.* **30**, e13468 (2018).
- 1026 14. Nichols, R. G. & Davenport, E. R. The relationship between the gut microbiome and
1027 host gene expression: a review. *Hum. Genet.* (2020) doi:10.1007/s00439-020-
1028 02237-0.
- 1029 15. Camp, J. G. *et al.* Microbiota modulate transcription in the intestinal epithelium
1030 without remodeling the accessible chromatin landscape. *Genome Res.* **24**, 1504–
1031 1516 (2014).
- 1032 16. Sommer, F., Nookaew, I., Sommer, N., Fogelstrand, P. & Bäckhed, F. Site-specific
1033 programming of the host epithelial transcriptome by the gut microbiota. *Genome*
1034 *Biol.* **16**, 62 (2015).
- 1035 17. Pan, W.-H. *et al.* Exposure to the gut microbiota drives distinct methylome and
1036 transcriptome changes in intestinal epithelial cells during postnatal development.
1037 *Genome Med.* **10**, 27 (2018).
- 1038 18. Davison, J. M. *et al.* Microbiota regulate intestinal epithelial gene expression by
1039 suppressing the transcription factor Hepatocyte nuclear factor 4 alpha. *Genome*
1040 *Res.* **27**, 1195–1206 (2017).
- 1041 19. Murdoch, C. C. & Rawls, J. F. Commensal Microbiota Regulate Vertebrate Innate
1042 Immunity-Insights From the Zebrafish. *Front. Immunol.* **10**, 2100 (2019).
- 1043 20. Broderick, N. A., Buchon, N. & Lemaitre, B. Microbiota-induced changes in

- 1044 drosophila melanogaster host gene expression and gut morphology. *MBio* **5**,
1045 e01117–14 (2014).
- 1046 21. Richards, A. L. *et al.* Genetic and transcriptional analysis of human host response
1047 to healthy gut microbiota. *mSystems* **1**, (2016).
- 1048 22. Richards, A. L. *et al.* Gut Microbiota Has a Widespread and Modifiable Effect on
1049 Host Gene Regulation. *mSystems* **4**, (2019).
- 1050 23. Morgan, X. C. *et al.* Associations between host gene expression, the mucosal
1051 microbiome, and clinical outcome in the pelvic pouch of patients with inflammatory
1052 bowel disease. *Genome Biol.* **16**, 67 (2015).
- 1053 24. Häslér, R. *et al.* Uncoupling of mucosal gene regulation, mRNA splicing and
1054 adherent microbiota signatures in inflammatory bowel disease. *Gut* **66**, 2087–2097
1055 (2017).
- 1056 25. Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory
1057 bowel diseases. *Nature* **569**, 655–662 (2019).
- 1058 26. Flemer, B. *et al.* Tumour-associated and non-tumour-associated microbiota in
1059 colorectal cancer. *Gut* **66**, 633–643 (2017).
- 1060 27. Dayama, G., Priya, S., Niccum, D. E., Khoruts, A. & Blekhman, R. Interactions
1061 between the gut microbiome and host gene regulation in cystic fibrosis. *Genome*
1062 *Med.* **12**, 12 (2020).
- 1063 28. Integrative HMP (iHMP) Research Network Consortium. The Integrative Human
1064 Microbiome Project. *Nature* **569**, 641–648 (2019).
- 1065 29. Witten, D. M., Tibshirani, R. & Hastie, T. A penalized matrix decomposition, with
1066 applications to sparse principal components and canonical correlation analysis.

- 1067 *Biostatistics* **10**, 515–534 (2009).
- 1068 30. Witten, D. M. & Tibshirani, R. J. Extensions of sparse canonical correlation analysis
1069 with applications to genomic data. *Stat. Appl. Genet. Mol. Biol.* **8**, Article28 (2009).
- 1070 31. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. *J. R. Stat. Soc.*
1071 *Series B Stat. Methodol.* **58**, 267–288 (1996).
- 1072 32. Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: the next generation. *Cell* **144**,
1073 646–674 (2011).
- 1074 33. Huang, Q. *et al.* LYRM2 directly regulates complex I activity to support tumor
1075 growth in colorectal cancer by oxidative phosphorylation. *Cancer Lett.* **455**, 36–47
1076 (2019).
- 1077 34. Lin, C.-S. *et al.* Role of mitochondrial function in the invasiveness of human colon
1078 cancer cells. *Oncol. Rep.* **39**, 316–330 (2018).
- 1079 35. Vellinga, T. T. *et al.* SIRT1/PGC1 α -Dependent Increase in Oxidative
1080 Phosphorylation Supports Chemotherapy Resistance of Colon Cancer. *Clin.*
1081 *Cancer Res.* **21**, 2870–2879 (2015).
- 1082 36. Jackson, D. N. & Theiss, A. L. Gut bacteria signaling to mitochondria in intestinal
1083 inflammation and cancer. *Gut Microbes* 1–20 (2019).
- 1084 37. Dotan, I. *et al.* The role of integrins in the pathogenesis of inflammatory bowel
1085 disease: Approved and investigational anti-integrin therapies. *Med. Res. Rev.* **40**,
1086 245–262 (2020).
- 1087 38. Integrins: Signalling and Disease. in *eLS* (ed. John Wiley & Sons Ltd) vol. 11 48
1088 (John Wiley & Sons, Ltd, 2001).
- 1089 39. Seinen, M. L., van Nieuw Amerongen, G. P., de Boer, N. K. H. & van Bodegraven,

- 1090 A. A. Rac Attack: Modulation of the Small GTPase Rac in Inflammatory Bowel
1091 Disease and Thiopurine Therapy. *Mol. Diagn. Ther.* **20**, 551–557 (2016).
- 1092 40. Kotelevets, L. & Chastre, E. Rac1 Signaling: From Intestinal Homeostasis to
1093 Colorectal Cancer Metastasis. *Cancers* **12**, (2020).
- 1094 41. Muise, A. M. *et al.* Single nucleotide polymorphisms that increase expression of the
1095 guanosine triphosphatase RAC1 are associated with ulcerative colitis.
1096 *Gastroenterology* **141**, 633–641 (2011).
- 1097 42. Abdulamir, A. S., Hafidh, R. R. & Abu Bakar, F. The association of *Streptococcus*
1098 *bovis/gallolyticus* with colorectal tumors: the nature and the underlying mechanisms
1099 of its etiological role. *J. Exp. Clin. Cancer Res.* **30**, 11 (2011).
- 1100 43. Aymeric, L. *et al.* Colorectal cancer specific conditions promote *Streptococcus*
1101 *gallolyticus* gut colonization. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E283–E291
1102 (2018).
- 1103 44. Qiu, Z. *et al.* Targeted Metagenome Based Analyses Show Gut Microbial Diversity
1104 of Inflammatory Bowel Disease patients. *Indian J. Microbiol.* **57**, 307–315 (2017).
- 1105 45. Dinakaran, V. *et al.* Identification of Specific Oral and Gut Pathogens in Full
1106 Thickness Colon of Colitis Patients: Implications for Colon Motility. *Front. Microbiol.*
1107 **9**, 3220 (2018).
- 1108 46. Ma, H.-Q., Yu, T.-T., Zhao, X.-J., Zhang, Y. & Zhang, H.-J. Fecal microbial
1109 dysbiosis in Chinese patients with inflammatory bowel disease. *World J.*
1110 *Gastroenterol.* **24**, 1464–1477 (2018).
- 1111 47. Wang, S. *et al.* Diet-induced remission in chronic enteropathy is associated with
1112 altered microbial community structure and synthesis of secondary bile acids.

- 1113 *Microbiome* **7**, 126 (2019).
- 1114 48. Lucke, K., Miehlike, S., Jacobs, E. & Schuppler, M. Prevalence of Bacteroides and
1115 Prevotella spp. in ulcerative colitis. *J. Med. Microbiol.* **55**, 617–624 (2006).
- 1116 49. Pittayanon, R. *et al.* Gut Microbiota in Patients With Irritable Bowel Syndrome-A
1117 Systematic Review. *Gastroenterology* **157**, 97–108 (2019).
- 1118 50. Duboc, H. *et al.* Increase in fecal primary bile acids and dysbiosis in patients with
1119 diarrhea-predominant irritable bowel syndrome: Bile acids and dysbiosis in IBS-D
1120 patients. *Neurogastroenterology & Motility* **24**, 513–e247 (2012).
- 1121 51. Rajilić-Stojanović, M. *et al.* Global and deep molecular analysis of microbiota
1122 signatures in fecal samples from patients with irritable bowel syndrome.
1123 *Gastroenterology* **141**, 1792–1801 (2011).
- 1124 52. Jeffery, I. B. *et al.* An irritable bowel syndrome subtype defined by species-specific
1125 alterations in faecal microbiota. *Gut* **61**, 997–1006 (2012).
- 1126 53. Szatmári, T., Ötvös, R., Hjerpe, A. & Dobra, K. Syndecan-1 in Cancer: Implications
1127 for Cell Signaling, Differentiation, and Prognostication. *Dis. Markers* **2015**, 796052
1128 (2015).
- 1129 54. Vicente, C. M. *et al.* Heparan Sulfate Proteoglycans in Human Colorectal Cancer.
1130 *Anal. Cell. Pathol.* **2018**, 8389595 (2018).
- 1131 55. Wei, H.-T., Guo, E.-N., Dong, B.-G. & Chen, L.-S. Prognostic and clinical
1132 significance of syndecan-1 in colorectal cancer: a meta-analysis. *BMC*
1133 *Gastroenterol.* **15**, 152 (2015).
- 1134 56. Purcell, R. V., Visnovska, M., Biggs, P. J., Schmeier, S. & Frizelle, F. A. Distinct gut
1135 microbiome patterns associate with consensus molecular subtypes of colorectal

- 1136 cancer. *Sci. Rep.* **7**, 11590 (2017).
- 1137 57. Haghi, F., Goli, E., Mirzaei, B. & Zeighami, H. The association between fecal
1138 enterotoxigenic *B. fragilis* with colorectal cancer. *BMC Cancer* **19**, 879 (2019).
- 1139 58. Ulger Toprak, N. *et al.* A possible role of *Bacteroides fragilis* enterotoxin in the
1140 aetiology of colorectal cancer. *Clin. Microbiol. Infect.* **12**, 782–786 (2006).
- 1141 59. Forbes, J. D. *et al.* A comparative study of the gut microbiota in immune-mediated
1142 inflammatory diseases—does a common dysbiosis exist? *Microbiome* vol. 6 (2018).
- 1143 60. Knights, D., Lassen, K. G. & Xavier, R. J. Advances in inflammatory bowel disease
1144 pathogenesis: linking host genetics and the microbiome. *Gut* **62**, 1505–1510
1145 (2013).
- 1146 61. Morgan, X. C. *et al.* Dysfunction of the intestinal microbiome in inflammatory bowel
1147 disease and treatment. *Genome Biol.* **13**, R79 (2012).
- 1148 62. Lavelle, A. *et al.* Spatial variation of the colonic microbiota in patients with
1149 ulcerative colitis and control volunteers. *Gut* **64**, 1553–1561 (2015).
- 1150 63. Mager, L. F., Wasmer, M.-H., Rau, T. T. & Krebs, P. Cytokine-Induced Modulation
1151 of Colorectal Cancer. *Front. Oncol.* **6**, 96 (2016).
- 1152 64. Fender, A. W., Nutter, J. M., Fitzgerald, T. L., Bertrand, F. E. & Sigounas, G. Notch-
1153 1 promotes stemness and epithelial to mesenchymal transition in colorectal cancer.
1154 *J. Cell. Biochem.* **116**, 2517–2527 (2015).
- 1155 65. Reedijk, M. *et al.* Activation of Notch signaling in human colon adenocarcinoma. *Int.*
1156 *J. Oncol.* **33**, 1223–1229 (2008).
- 1157 66. Luo, D. & Ge, W. MeCP2 Promotes Colorectal Cancer Metastasis by Modulating
1158 ZEB1 Transcription. *Cancers* **12**, (2020).

- 1159 67. Sébert, M., Sola-Tapias, N., Mas, E., Barreau, F. & Ferrand, A. Protease-Activated
1160 Receptors in the Intestine: Focus on Inflammation and Cancer. *Front. Endocrinol.*
1161 **10**, 717 (2019).
- 1162 68. Duan, L., Rao, X., Braunstein, Z., Toomey, A. C. & Zhong, J. Role of Incretin Axis
1163 in Inflammatory Bowel Disease. *Front. Immunol.* **8**, 1734 (2017).
- 1164 69. Mustafa, S. A. *et al.* SUMOylation pathway alteration coupled with downregulation of
1165 SUMO E2 enzyme at mucosal epithelium modulates inflammation in inflammatory
1166 bowel disease. *Open Biol.* **7**, (2017).
- 1167 70. Clarke, G. *et al.* Marked elevations in pro-inflammatory polyunsaturated fatty acid
1168 metabolites in females with irritable bowel syndrome. *J. Lipid Res.* **51**, 1186–1192
1169 (2010).
- 1170 71. Nielsen, O. H., Ahnfelt-Rønne, I. & Elmgreen, J. Abnormal metabolism of
1171 arachidonic acid in chronic inflammatory bowel disease: enhanced release of
1172 leucotriene B4 from activated neutrophils. *Gut* **28**, 181–185 (1987).
- 1173 72. Campbell, A. G. *et al.* Diversity and genomic insights into the uncultured Chloroflexi
1174 from the human microbiota. *Environ. Microbiol.* **16**, 2635–2643 (2014).
- 1175 73. Shao, T. *et al.* Combined Signature of the Fecal Microbiome and Metabolome in
1176 Patients with Gout. *Front. Microbiol.* **8**, 268 (2017).
- 1177 74. Idris, A., Hasnain, S. Z., Huat, L. Z. & Koh, D. Human diseases, immunity and the
1178 oral microbiota—Insights gained from metagenomic studies. *Oral Science*
1179 *International* **14**, 27–32 (2017).
- 1180 75. Narayanan, V., Peppelenbosch, M. P. & Konstantinov, S. R. Human fecal
1181 microbiome-based biomarkers for colorectal cancer. *Cancer Prev. Res.* **7**, 1108–

- 1182 1111 (2014).
- 1183 76. Lu, Y. *et al.* Mucosal adherent bacterial dysbiosis in patients with colorectal
1184 adenomas. *Sci. Rep.* **6**, 26337 (2016).
- 1185 77. Li, J. *et al.* Transcriptional Profiling Reveals the Regulatory Role of CXCL8 in
1186 Promoting Colorectal Cancer. *Front. Genet.* **10**, 1360 (2019).
- 1187 78. Viet, H. T., Wågsäter, D., Hugander, A. & Dimberg, J. Interleukin-1 receptor
1188 antagonist gene polymorphism in human colorectal cancer. *Oncol. Rep.* **14**, 915–
1189 918 (2005).
- 1190 79. Burada, F. *et al.* IL-1RN +2018T>C polymorphism is correlated with colorectal
1191 cancer. *Mol. Biol. Rep.* **40**, 2851–2857 (2013).
- 1192 80. Wolf, M. J. *et al.* Endothelial CCR2 signaling induced by colon carcinoma cells
1193 enables extravasation via the JAK2-Stat5 and p38MAPK pathway. *Cancer Cell* **22**,
1194 91–105 (2012).
- 1195 81. Li, S.-Q. *et al.* The Expression of Formyl Peptide Receptor 1 is Correlated with
1196 Tumor Invasion of Human Colorectal Cancer. *Sci. Rep.* **7**, 5918 (2017).
- 1197 82. Baxter, N. T., Zackular, J. P., Chen, G. Y. & Schloss, P. D. Structure of the gut
1198 microbiome following colonization with human feces determines colonic tumor
1199 burden. *Microbiome* **2**, 20 (2014).
- 1200 83. Chen, L. *et al.* Characteristics of fecal and mucosa-associated microbiota in
1201 Chinese patients with inflammatory bowel disease. *Medicine* **93**, e51 (2014).
- 1202 84. Gao, X. *et al.* Chronic stress promotes colitis by disturbing the gut microbiota and
1203 triggering immune system response. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E2960–
1204 E2969 (2018).

- 1205 85. Li, X. L. *et al.* Bioinformatic analysis of potential candidates for therapy of
1206 inflammatory bowel disease. *Eur. Rev. Med. Pharmacol. Sci.* **19**, 4275–4284
1207 (2015).
- 1208 86. Yukawa, T. *et al.* Differential expression of vasoactive intestinal peptide receptor 1
1209 expression in inflammatory bowel disease. *Int. J. Mol. Med.* **20**, 161–167 (2007).
- 1210 87. Imhann, F. *et al.* Interplay of host genetics and gut microbiota underlying the onset
1211 and clinical presentation of inflammatory bowel disease. *Gut* **67**, 108–119 (2018).
- 1212 88. Potts, P. R. & Yu, H. The SMC5/6 complex maintains telomere length in ALT
1213 cancer cells through SUMOylation of telomere-binding proteins. *Nat. Struct. Mol.*
1214 *Biol.* **14**, 581–590 (2007).
- 1215 89. Vandeputte, D. *et al.* Stool consistency is strongly associated with gut microbiota
1216 richness and composition, enterotypes and bacterial growth rates. *Gut* **65**, 57–62
1217 (2016).
- 1218 90. Su, T. *et al.* Altered Intestinal Microbiota with Increased Abundance of *Prevotella* Is
1219 Associated with High Risk of Diarrhea-Predominant Irritable Bowel Syndrome.
1220 *Gastroenterol. Res. Pract.* **2018**, 6961783 (2018).
- 1221 91. Pazmandi, J., Kalinichenko, A., Ardy, R. C. & Boztug, K. Early-onset inflammatory
1222 bowel disease as a model disease to identify key regulators of immune
1223 homeostasis mechanisms. *Immunol. Rev.* **287**, 162–185 (2019).
- 1224 92. Schewe, M. *et al.* Secreted Phospholipases A2 Are Intestinal Stem Cell Niche
1225 Factors with Distinct Roles in Homeostasis, Inflammation, and Cancer. *Cell Stem*
1226 *Cell* **19**, 38–51 (2016).
- 1227 93. Feng, Q. *et al.* Gut microbiome development along the colorectal adenoma–

- 1228 carcinoma sequence. *Nat. Commun.* **6**, 6528 (2015).
- 1229 94. Gagnière, J. *et al.* Gut microbiota imbalance and colorectal cancer. *World J.*
1230 *Gastroenterol.* **22**, 501–518 (2016).
- 1231 95. Chen, W., Liu, F., Ling, Z., Tong, X. & Xiang, C. Human intestinal lumen and
1232 mucosa-associated microbiota in patients with colorectal cancer. *PLoS One* **7**,
1233 e39743 (2012).
- 1234 96. Kojima, A. *et al.* Infection of specific strains of *Streptococcus mutans*, oral bacteria,
1235 confers a risk of ulcerative colitis. *Sci. Rep.* **2**, 332 (2012).
- 1236 97. Bennet, S. M. P., Ohman, L. & Simren, M. Gut microbiota as potential orchestrators
1237 of irritable bowel syndrome. *Gut Liver* **9**, 318–331 (2015).
- 1238 98. Piewngam, P. *et al.* Pathogen elimination by probiotic *Bacillus* via signalling
1239 interference. *Nature* **562**, 532–537 (2018).
- 1240 99. Acton, D. S., Plat-Sinnige, M. J. T., van Wamel, W., de Groot, N. & van Belkum, A.
1241 Intestinal carriage of *Staphylococcus aureus*: how does its frequency compare with
1242 that of nasal carriage and what is its clinical impact? *Eur. J. Clin. Microbiol. Infect.*
1243 *Dis.* **28**, 115–127 (2009).
- 1244 100. Bettenworth, D. *et al.* Crohn's disease complicated by intestinal infection with
1245 methicillin-resistant *Staphylococcus aureus*. *World J. Gastroenterol.* **19**, 4418–4421
1246 (2013).
- 1247 101. Chiba, M. *et al.* *Staphylococcus aureus* in inflammatory bowel disease. *Scand. J.*
1248 *Gastroenterol.* **36**, 615–620 (2001).
- 1249 102. Attiê, R., Chinen, L. T. D., Yoshioka, E. M., Silva, M. C. F. & de Lima, V. C. C.
1250 Acute bacterial infection negatively impacts cancer specific survival of colorectal

- 1251 cancer patients. *World J. Gastroenterol.* **20**, 13930–13935 (2014).
- 1252 103.Noguchi, N. *et al.* Specific clones of *Staphylococcus lugdunensis* may be
1253 associated with colon carcinoma. *J. Infect. Public Health* **11**, 39–42 (2018).
- 1254 104.Rinttilä, T., Lyra, A., Krogius-Kurikka, L. & Palva, A. Real-time PCR analysis of
1255 enteric pathogens from fecal samples of irritable bowel syndrome subjects. *Gut*
1256 *Pathog.* **3**, 6 (2011).
- 1257 105.Comelli, E. M. *et al.* Biomarkers of human gastrointestinal tract regions. *Mamm.*
1258 *Genome* **20**, 516–527 (2009).
- 1259 106.Xia, B., Zhang, K. & Liu, C. PYGB Promoted Tumor Progression by Regulating
1260 Wnt/ β -Catenin Pathway in Gastric Cancer. *Technol. Cancer Res. Treat.* **19**,
1261 1533033820926592 (2020).
- 1262 107.Xu, H., Cao, H. & Xiao, G. Signaling via PINCH: Functions, binding partners and
1263 implications in human diseases. *Gene* **594**, 10–15 (2016).
- 1264 108.Hehlhans, S., Haase, M. & Cordes, N. Signalling via integrins: implications for cell
1265 survival and anticancer strategies. *Biochim. Biophys. Acta* **1775**, 163–180 (2007).
- 1266 109.Burczynski, M. E. *et al.* Molecular classification of Crohn's disease and ulcerative
1267 colitis patients using transcriptional profiles in peripheral blood mononuclear cells.
1268 *J. Mol. Diagn.* **8**, 51–61 (2006).
- 1269 110.de la Motte, C. A. Hyaluronan in intestinal homeostasis and inflammation:
1270 implications for fibrosis. *Am. J. Physiol. Gastrointest. Liver Physiol.* **301**, G945–9
1271 (2011).
- 1272 111.Kim, Y. & de la Motte, C. A. The Role of Hyaluronan Treatment in Intestinal Innate
1273 Host Defense. *Front. Immunol.* **11**, 569 (2020).

- 1274 112.Zheng, L., Riehl, T. E. & Stenson, W. F. Regulation of colonic epithelial repair in
1275 mice by Toll-like receptors and hyaluronic acid. *Gastroenterology* **137**, 2041–2051
1276 (2009).
- 1277 113.Seiden-Long, I. *et al.* Gab1 but not Grb2 mediates tumor progression in Met
1278 overexpressing colorectal cancer cells. *Carcinogenesis* **29**, 647–655 (2008).
- 1279 114.Liu, D.-Q. *et al.* Increased RIPK4 expression is associated with progression and
1280 poor prognosis in cervical squamous cell carcinoma patients. *Sci. Rep.* **5**, 11955
1281 (2015).
- 1282 115.Huang, X. *et al.* Phosphorylation of Dishevelled by protein kinase RIPK4 regulates
1283 Wnt signaling. *Science* **339**, 1441–1445 (2013).
- 1284 116.Liu, J.-Y. *et al.* RIPK4 promotes bladder urothelial carcinoma cell aggressiveness
1285 by upregulating VEGF-A through the NF- κ B pathway. *Br. J. Cancer* **118**, 1617–
1286 1627 (2018).
- 1287 117.Yang, X., Yue, Y. & Xiong, S. Dpep2 Emerging as a Modulator of Macrophage
1288 Inflammation Confers Protection Against CVB3-Induced Viral Myocarditis. *Front.*
1289 *Cell. Infect. Microbiol.* **9**, 57 (2019).
- 1290 118.Li, S., Ning, L.-G., Lou, X.-H. & Xu, G.-Q. Necroptosis in inflammatory bowel
1291 disease and other intestinal diseases. *World J Clin Cases* **6**, 745–752 (2018).
- 1292 119.Garcia-Carbonell, R., Yao, S.-J., Das, S. & Guma, M. Dysregulation of Intestinal
1293 Epithelial Cell RIPK Pathways Promotes Chronic Inflammation in the IBD Gut.
1294 *Front. Immunol.* **10**, 1094 (2019).
- 1295 120.Liu, Z.-Y. *et al.* RIP3 promotes colitis-associated colorectal cancer by controlling
1296 tumor cell proliferation and CXCL1-induced immune suppression. *Theranostics* **9**,

- 1297 3659–3673 (2019).
- 1298 121. Jayakumar, A. & Bothwell, A. L. M. RIPK3-Induced Inflammation by I-MDSCs
1299 Promotes Intestinal Tumors. *Cancer Res.* **79**, 1587–1599 (2019).
- 1300 122. Conev, N. V. *et al.* RIPK3 expression as a potential predictive and prognostic
1301 marker in metastatic colon cancer. *Clin. Invest. Med.* **42**, E31–E38 (2019).
- 1302 123. Wang, Z. *et al.* Regulation of innate immune responses by DAI (DLM-1/ZBP1) and
1303 other DNA-sensing molecules. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 5477–5482
1304 (2008).
- 1305 124. Kuriakose, T. & Kanneganti, T.-D. ZBP1: Innate Sensor Regulating Cell Death and
1306 Inflammation. *Trends Immunol.* **39**, 123–134 (2018).
- 1307 125. Zhang, T. *et al.* Influenza Virus Z-RNAs Induce ZBP1-Mediated Necroptosis. *Cell*
1308 **180**, 1115–1129.e13 (2020).
- 1309 126. O’Flanagan, C. H. & O’Neill, C. PINK1 signalling in cancer biology. *Biochim.*
1310 *Biophys. Acta* **1846**, 590–598 (2014).
- 1311 127. Chang, J. Y., Yi, H.-S., Kim, H.-W. & Shong, M. Dysregulation of mitophagy in
1312 carcinogenesis and tumor progression. *Biochim. Biophys. Acta Bioenerg.* **1858**,
1313 633–640 (2017).
- 1314 128. Novak, E. A. & Mollen, K. P. Mitochondrial dysfunction in inflammatory bowel
1315 disease. *Front Cell Dev Biol* **3**, 62 (2015).
- 1316 129. Alomair, A. O. *et al.* Colonic Mucosal Microbiota in Colorectal Cancer: A Single-
1317 Center Metagenomic Study in Saudi Arabia. *Gastroenterol. Res. Pract.* **2018**,
1318 5284754 (2018).
- 1319 130. Ai, D. *et al.* Identifying Gut Microbiota Associated With Colorectal Cancer Using a

- 1320 Zero-Inflated Lognormal Model. *Front. Microbiol.* **10**, 826 (2019).
- 1321 131.Chen, J. *et al.* An expansion of rare lineage intestinal microbes characterizes
1322 rheumatoid arthritis. *Genome Med.* **8**, 43 (2016).
- 1323 132.Taverniti, V. & Guglielmetti, S. Methodological issues in the study of intestinal
1324 microbiota in irritable bowel syndrome. *World J. Gastroenterol.* **20**, 8821–8836
1325 (2014).
- 1326 133.Sato, Y. *et al.* Up-regulated Annexin A1 expression in gastrointestinal cancer is
1327 associated with cancer invasion and lymph node metastasis. *Exp. Ther. Med.* **2**,
1328 239–243 (2011).
- 1329 134.Sena, A. *et al.* Dysregulation of anti-inflammatory annexin A1 expression in
1330 progressive Crohns Disease. *PLoS One* **8**, e76969 (2013).
- 1331 135.Babbin, B. A. *et al.* Annexin A1 regulates intestinal mucosal injury, inflammation,
1332 and repair. *J. Immunol.* **181**, 5035–5044 (2008).
- 1333 136.Mazmanian, S. K., Liu, C. H., Tzianabos, A. O. & Kasper, D. L. An
1334 immunomodulatory molecule of symbiotic bacteria directs maturation of the host
1335 immune system. *Cell* **122**, 107–118 (2005).
- 1336 137.Zitomersky, N. L. *et al.* Characterization of adherent bacteroidales from intestinal
1337 biopsies of children and young adults with inflammatory bowel disease. *PLoS One*
1338 **8**, e63686 (2013).
- 1339 138.Park, B. S. & Lee, J.-O. Recognition of lipopolysaccharide pattern by TLR4
1340 complexes. *Exp. Mol. Med.* **45**, e66 (2013).
- 1341 139.Rakoff-Nahoum, S., Paglino, J., Eslami-Varzaneh, F., Edberg, S. & Medzhitov, R.
1342 Recognition of commensal microflora by toll-like receptors is required for intestinal

- 1343 homeostasis. *Cell* **118**, 229–241 (2004).
- 1344 140. Yesudhas, D., Gosu, V., Anwar, M. A. & Choi, S. Multiple roles of toll-like receptor 4
1345 in colorectal cancer. *Front. Immunol.* **5**, 334 (2014).
- 1346 141. Lu, Y., Li, X., Liu, S., Zhang, Y. & Zhang, D. Toll-like Receptors and Inflammatory
1347 Bowel Disease. *Front. Immunol.* **9**, 72 (2018).
- 1348 142. Perera, M., Al-Hebshi, N. N., Speicher, D. J., Perera, I. & Johnson, N. W. Emerging
1349 role of bacteria in oral carcinogenesis: a review with special reference to perio-
1350 pathogenic bacteria. *J. Oral Microbiol.* **8**, 32762 (2016).
- 1351 143. Alam, M. T. *et al.* Microbial imbalance in inflammatory bowel disease patients at
1352 different taxonomic levels. *Gut Pathog.* **12**, 1 (2020).
- 1353 144. Nakatsu, G. *et al.* Gut mucosal microbiome across stages of colorectal
1354 carcinogenesis. *Nat. Commun.* **6**, 8727 (2015).
- 1355 145. Graham, D. B. & Xavier, R. J. Pathway paradigms revealed from the genetics of
1356 inflammatory bowel disease. *Nature* **578**, 527–539 (2020).
- 1357 146. Farooqi, A. A., de la Roche, M., Djamgoz, M. B. A. & Siddik, Z. H. Overview of the
1358 oncogenic signaling pathways in colorectal cancer: Mechanistic insights. *Semin.*
1359 *Cancer Biol.* **58**, 65–79 (2019).
- 1360 147. Lanis, J. M., Kao, D. J., Alexeev, E. E. & Colgan, S. P. Tissue metabolism and the
1361 inflammatory bowel diseases. *J. Mol. Med.* **95**, 905–913 (2017).
- 1362 148. Francescone, R., Hou, V. & Grivennikov, S. I. Cytokines, IBD, and colitis-
1363 associated cancer. *Inflamm. Bowel Dis.* **21**, 409–418 (2015).
- 1364 149. Barbara, G. *et al.* The immune system in irritable bowel syndrome. *J.*
1365 *Neurogastroenterol. Motil.* **17**, 349–359 (2011).

- 1366 150. Luca, F., Kupfer, S. S., Knights, D., Khoruts, A. & Blekhman, R. Functional
1367 Genomics of Host-Microbiome Interactions in Humans. *Trends Genet.* **34**, 30–40
1368 (2018).
- 1369 151. Burns, M. B. *et al.* APOBEC3B is an enzymatic source of mutation in breast cancer.
1370 *Nature* **494**, 366–370 (2013).
- 1371 152. Andrews, S. & Others. FastQC: a quality control tool for high throughput sequence
1372 data. (2010).
- 1373 153. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
1374 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
- 1375 154. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low
1376 memory requirements. *Nat. Methods* **12**, 357–360 (2015).
- 1377 155. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose
1378 program for assigning sequence reads to genomic features. *Bioinformatics* **30**,
1379 923–930 (2014).
- 1380 156. Poyet, M. *et al.* A library of human gut bacterial isolates paired with longitudinal
1381 multiomics data enables mechanistic microbiome research. *Nat. Med.* (2019)
1382 doi:10.1038/s41591-019-0559-3.
- 1383 157. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the
1384 integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat.*
1385 *Protoc.* **4**, 1184–1191 (2009).
- 1386 158. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
1387 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 1388 159. Salter, S. J. *et al.* Reagent and laboratory contamination can critically impact

- 1389 sequence-based microbiome analyses. *BMC Biol.* **12**, 87 (2014).
- 1390 160.Oksanen, J. *et al.* The vegan package. *Community ecology package* **10**, 631–637
1391 (2007).
- 1392 161.Quinn, T. P., Erb, I., Richardson, M. F. & Crowley, T. M. Understanding sequencing
1393 data as compositions: an outlook and review. *Bioinformatics* **34**, 2870–2878 (2018).
- 1394 162.Hotelling, H. Relations Between Two Sets of Variates. *Biometrika* **28**, 321–377
1395 (1936).
- 1396 163.Witten, D., Tibshirani, R., Gross, S. & Narasimhan, B. Pma: Penalized multivariate
1397 analysis. *R package version 1*, (2013).
- 1398 164.Ash, J. T., Darnell, G., Munro, D., Engelhardt, B. E. & Authorship, I. E. Joint
1399 analysis of gene expression levels and histological images identifies genes
1400 associated with tissue morphology. doi:10.1101/458711.
- 1401 165.Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based
1402 approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci.*
1403 *U. S. A.* **102**, 15545–15550 (2005).
- 1404 166.Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**,
1405 1739–1740 (2011).
- 1406 167.Waardenberg, A. J., Basset, S. D., Bouveret, R. & Harvey, R. P. CompGO: an R
1407 package for comparing and visualizing Gene Ontology enrichment differences
1408 between DNA binding experiments. *BMC Bioinformatics* **16**, 275 (2015).
- 1409 168.Morris, J. A. & Gardner, M. J. Calculating confidence intervals for relative risks
1410 (odds ratios) and standardised ratios and rates. *Br. Med. J.* **296**, 1313–1316
1411 (1988).

- 1412 169.Smoot, M. E., Ono, K., Ruscheinski, J., Wang, P.-L. & Ideker, T. Cytoscape 2.8:
1413 new features for data integration and network visualization. *Bioinformatics* **27**, 431–
1414 432 (2011).
- 1415 170.Friedman, J., Hastie, T. & Tibshirani, R. glmnet: Lasso and elastic-net regularized
1416 generalized linear models. *R package version 1*, (2009).
- 1417 171.Dezeure, R., Bühlmann, P., Meier, L. & Meinshausen, N. High-Dimensional
1418 Inference: Confidence Intervals, p-Values and R-Software hdi. *Stat. Sci.* **30**, 533–
1419 558 (2015).
- 1420 172.Zhang, C.-H. & Zhang, S. S. Confidence intervals for low dimensional parameters
1421 in high dimensional linear models. *J. R. Stat. Soc. Series B Stat. Methodol.* **76**,
1422 217–242 (2014).
- 1423 173.Meinshausen, N. & Bühlmann, P. Stability selection. *J. R. Stat. Soc. Series B Stat.*
1424 *Methodol.* **72**, 417–473 (2010).
- 1425 174.Hofner, B., Boccuto, L. & Göker, M. Controlling false discoveries in high-
1426 dimensional situations: boosting with stability selection. *BMC Bioinformatics* **16**,
1427 144 (2015).
- 1428