

## 1 **Index Supplementary Appendix**

2 Supplementary Methods

3 Supplementary figure 1. Q-Q plot of whole exome-wide association analyses *P* values in IBD cohort and  
4 Lifelines-DEEP cohort respectively.

5 Supplementary figure 2. Box-plot of associations between *SEC16A* (rs10781497) and *WDR78* (rs74609208)  
6 and microbial pathways in IBD cohort and Lifelines-DEEP cohort respectively.

7 Supplementary figure 3. Q-Q plot of gene-based burden test *P* values in IBD cohort and Lifelines-DEEP cohort  
8 respectively.

9 Supplementary figure 4. Distribution of interaction t-statistic based on 999 permutation tests.

10 Supplementary figure 5. Box-plot of association *P* values between *LCT* (rs748841) and *Bifidobacterium*  
11 *adolescentis* in IBD cohort and Lifelines-DEEP cohort respectively.

12 Supplementary table S1 Demographic data for two cohorts.

13 Supplementary table S2.1 Microbial traits in IBD cohort.

14 Supplementary table S2.2 Microbial traits in LifeLines-DEEP cohort.

15 Supplementary table S3 LLDeep cohort-specific mbQTLs used for interaction analyses.

16 Supplementary table S4 Supplementary table S4 IBD cohort-specific mbQTLs used for interaction analyses.

17 Supplementary table S5.1 Summary statistics of genetics case control (IBD vs. control)

18 Supplementary table S5.2 Summary statistics of genetics case control (CD vs. control)

19 Supplementary table S5.3 Summary statistics of genetics case control (UC vs. control)

20 Supplementary table S6 Genetics case-control analysis of 3010 variants located in reported IBD loci (De Lange  
21 et al. Nat. Genet. 2017)

22 Supplementary table S7 Microbial taxa and pathways comparison between IBD/CD/UC and controls

23 Supplementary table S8 mbQTLs in exome-wide approach (adjusted for age, gender, read depth, BMI,  
24 smoking, antibiotics, PPI, laxatives and intestinal disease location (specific to IBD cohort)).

25 Supplementary table S9 Gene function enrichment analysis for genes harboring mbQTLs in whole  
26 exome-wide analysis.

27 Supplementary table S10 mbQTLs in targeted approach.

28 Supplementary table S11 mbQTLs stratified by CD, UC and IBD unclassified (IBDU)

- 29 Supplementary table S12 Interaction Analysis.
- 30 Supplementary table S13 Associations between variants in *LCT* and bacterial taxa.
- 31 Supplementary table S14 Top 20 gene-based associations between CNV and microbial taxa/pathway.
- 32

### 33 **Supplementary Methods**

#### 34 Whole exome sequencing of IBD patients and population controls

35 Whole exome sequencing (WES) was performed on blood samples taken from 939 LifeLines-DEEP  
36 participants and 525 IBD patients. DNA isolation was performed using the AutoPure LS procedure from  
37 Qiagen. Library preparation and sequencing was done at the Broad Institute of MIT and Harvard. On average,  
38 86.06 million high quality reads were obtained for each sample with 98.85% of reads aligned (human  
39 reference genome hg19), resulting in a coverage of 81% of the target region with a read depth of >30X. Next,  
40 the Genome Analysis Toolkit of the Broad Institute was used for calling single-nucleotide polymorphism and  
41 insertions/deletions (<https://software.broadinstitute.org/gatk/>).

42 The following variant/sample filtering parameters were applied to WES data using PLINK tool (v.1.9)<sup>1</sup>: 1)  
43 variants with a call rate <0.99 were removed. Variants with a minor allele frequency (MAF) >5% were used  
44 for microbial quantitative trait loci (mbQTLs) analyses of common variants, and variants with a MAF <5%  
45 were used for low-frequency and rare mbQTL mapping. 2) In the Hardy-Weinberg equilibrium test, we used a  
46 P value <0.0001 as a significance cutoff in the LifeLines-DEEP cohort and discarded those variants in the both  
47 cohorts. 3) To remove related samples, we used PLINK to calculate Identity-by-descent and removed samples  
48 with IBD >0.185<sup>2</sup>. 4) To identify ancestry-based genetic outliers in our dataset, we merged the WES data with

49 genomes of Europeans from publically available 1000 Genome Project (Phase 3) data  
50 (<http://www.internationalgenome.org/>), and performed principal component analysis (PCA) analysis based  
51 on SNPs shared between datasets. Outliers were defined as samples which fall outside of a mean  $\pm$  3 SD  
52 interval in both of the first two PCs, and these samples were removed. 5) determining sex (based on  
53 heterozygosity rates) and identifying mismatched samples were based on the inbreeding coefficient (lower  
54 than 0.4 for females and higher than 0.7 for males). GATK gCNV was used for copy number variation (CNV)  
55 calling. Common CNVs (site frequency >1%) were removed (<https://gatkforums.broadinstitute.org/gatk>).  
56 Finally, we also excluded participants who had their colon removed due to the large effect this procedure has  
57 on the gut microbiome<sup>3</sup>. These filtering steps led to exclusion of 19 samples from LifeLines-DEEP and 90  
58 samples from IBD, and we retained 920 LifeLines-DEEP individuals, 435 IBD individuals, 73,164 common  
59 variants, 98,878 rare variants and 1046 CNVs (site frequency <1%) for downstream analysis.

#### 60 Metagenomic sequencing of gut microbiota and data processing

61 Participants were asked to freeze a stool sample at home within 15 min of production. A research nurse  
62 visited each participant at home shortly after production to collect the sample on dry ice for transport to the  
63 laboratory at -80°C. Microbial DNA was isolated using the Qiagen AllPrep DNA/RNA Mini Kit (Qiagen; cat.  
64 #80204). Metagenomic shotgun sequencing was performed using the Illumina MiSeq platform. An average of

65 3.0 Gb of data (around 32.3 million reads) was obtained per sample. Reads belonging to the human genome  
66 were removed by mapping the data to the human reference genome (version NCBI37) with kneaddata  
67 (v0.5.1, <http://huttenhower.sph.harvard.edu/kneaddata>).

68 Profiling of microbiome taxonomic and functional composition was done using MetaPhlan (v2.6.0)  
69 (<http://huttenhower.sph.harvard.edu/metaphlan>) and HUMAnN2 (v0.6.1)  
70 (<http://huttenhower.sph.harvard.edu/humann2>). Within the IBD cohort, we found 483 microbial metabolic  
71 pathways and 1455 taxa, including 13 phyla, 23 classes, 32 orders, 70 families, 178 genera, 578 species and  
72 561 strains. Within LifeLines-DEEP, we found 468 pathways and 1375 taxa, including 15 phyla, 24 classes, 33  
73 orders, 74 families, 176 genera, 573 species and 480 strains.

74 For each cohort, taxa present in fewer than 10% of samples and pathways present in fewer than 25% of  
75 samples were excluded from the analyses. We removed the redundant taxa by keeping the lowest taxonomic  
76 level that shared identical abundance (for example, species *Odoribacter splanchnicus* and its strain  
77 *GCF\_000190535* had the same detected relative abundance in all samples, so we only kept the lowest level  
78 taxon, *GCF\_000190535*, in this case). Unclassified and unintegrated metabolic pathways were also excluded  
79 as they have little informative biological meaning.

80 Statistical analyses

81 We used a well-established pipeline for our mbQTL analysis<sup>4</sup>  
82 ([https://github.com/alexa-kur/miQTL\\_cookbook](https://github.com/alexa-kur/miQTL_cookbook)). For association tests between microbial features and  
83 common individual variants (MAF > 5%), we normalized the relative abundances of microbial taxa and  
84 metabolic pathways data through inverse rank transformation. Multivariate linear regression was used to  
85 adjust for the effect of confounders using the following model:

$$\begin{aligned} 86 \quad \textit{feature} = & (\textit{intercept}) + \textit{age} + \textit{gender} + \textit{BMI} + \textit{smoking status} + \textit{medication usage} \\ 87 \quad & + \textit{sequence depth} + \textit{disease location (only for IBD)} \end{aligned}$$

88 The corrected microbiome features (residuals from linear model) were considered as quantitative traits. The  
89 rank-based Spearman correlation method was applied to determine the relationship between non-zero  
90 microbiome features and each host genetic variant (where variants were encoded as 0 for homozygote of  
91 major allele, 1 for heterozygotes and 2 for homozygote of minor allele).

92 Restricted by the relatively small sample size in our non-common variants association analysis, we used the  
93 gene-based burden test by adding the variant's score instead of individual genotype into the correlation test,  
94 keeping only PTVs with MAF <5% and, or CNVs with MAF <1%, calculating per-gene scores using the  
95 MetaSKAT packages in R v.3.5.0.

96 For meta-analyses, the metap package (<https://rdrr.io/cran/metap/>) in R was used to perform a

97 weighted-Z-score approach using 'metap' package, considering sample size and separate  $P$  values. All  
98 significance thresholds were calculated by Bonferroni method with respect to the number of tested variants.  
99 For interaction analyses, we added a disease and genotype interaction term to the linear model:

$$100 \quad \text{corrected feature} \sim (\text{intercept}) + \text{genotype} + \text{disease} + \text{genotype} \times \text{disease}$$

101 To ensure the interaction results are not biased by potentially inflated statistics, we reassigned the disease  
102 status across all samples randomly 999 times and retested the interactions. Significance thresholds were set  
103 based on the Bonferroni method corrected by the number of interaction tests. At whole-exome-wide level,  
104 we observed 44 randomly significant interaction signals out of 999 permutation tests (around 0.04 random  
105 signals for each permutation on average), while we observed 14 significant signals from our non-permuted  
106 test, suggesting that the number of observed interactions is enriched and unlikely to have occurred by  
107 chance. The similar tests were also performed for targeted analysis. In addition, the distributions of  
108 interaction  $t$ -statistics and empirical  $P$  values were also assessed (**Supplementary figure 5**).

109 The recessive association of a SNP (rs4988235, G allele) near the gene *LCT* with abundance of  
110 lactose-metabolizing Bifidobacteria is well established in previous research. To evaluate this in our data, we  
111 used a recessive model (labeling the homozygote of minor allele as 1, and the other two genotypes as 0) to  
112 investigate the correlation between *LCT* variants (five protein coding variants linked to rs4988235, including

113 rs748841, rs12988076, rs6719488, rs2236783 and rs309180, linkage disequilibrium  $R^2 > 0.7$ ) and all  
114 taxonomic abundances using Spearman correlation in R (v.3.5.0).

115

#### 116 **References for methods**

- 117 1. Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and  
118 population-based linkage analyses. *The American journal of human genetics*, 2007, 81(3): 559-575.
- 119 2. Anderson, Carl A., et al. "Data quality control in genetic case-control association studies." *Nature*  
120 *protocols* 5.9 (2010): 1564.
- 121 3. Tyler A D, Knox N, Kabakchiev B, et al. Characterization of the gut-associated microbiome in inflammatory  
122 pouch complications following ileal pouch-anal anastomosis. *PLoS one*, 2013, 8(9): e66934.
- 123 4. Bonder, M. J. *et al.* The effect of host genetics on the gut microbiome. *Nat. Genet.* **48**, 1407–1412 (2016).

124

125

#### 126 **Supplementary figures**

127 **Supplementary figure 1.** Q-Q plot of exome-wide association analyses  $P$  values in IBD cohort and  
128 Lifelines-DEEP cohort respectively. X-axis indicates the expected  $-\log_{10} P$  values and Y-axis indicates the  
129 observed  $-\log_{10} P$  values. A) Whole exome-wide approach of 73,164 common variants with microbial  
130 superpathway of glyoxylate bypass and TCA (TCA-GLYOX-BYPASS), B) Whole exome-wide approach of 73,164  
131 common variants with microbial superpathway of acetyl-CoA biosynthesis (PWY-5173).

132 **Supplementary figure 2.** Box-plot of associations between *SEC16A* (rs10781497) and *WDR78* (rs74609208)



133 and microbial pathways in IBD cohort and Lifelines-DEEP cohort respectively. X-axis indicates the genotype of  
134 variants and Y-axis indicates the relative abundance of microbial pathways. A) the associations between  
135 *SEC16A* (rs10781497) and Thiamin diphosphate biosynthesis I (THISYN-PWY), Thiazole biosynthesis I (*E. coli*)  
136 (PWY-6892), B) the associations between *WDR78* (rs74609208) and dTDP-L-rhamnose biosynthesis I  
137 (DTDPRHAMSYN-PWY).

138 *P*, *P* values of associations; *r*, correlation coefficient; IBD, IBD cohort; LLD, LifeLines-DEEP cohort

139 **Supplementary figure 3.** Q-Q plot of burden test *P* values in IBD cohort and Lifelines-DEEP cohort  
140 respectively. X-axis indicates the expected  $-\log_{10} P$  values and Y-axis indicates the observed  $-\log_{10} P$  values.  
141 A) Burden tests of 980 genes with microbial pathway homolactic fermentation (ANAEROFrucat-PWY), B)  
142 Burden tests of 980 genes with microbial pathway glucose and xylose degradation (PWY-6901).

143 **Supplementary figure 4.** Distribution of interaction t-statistic based on 999 permutation tests. We  
144 permuted disease status across all samples 999 times and got the density distribution of interaction  
145 t-statistic. Red line indicates the t-statistic in our non-permutation test. A) An example from whole  
146 exome-wide level interaction analyses. B) An example from targeted level interaction analyses.

147 **Supplementary figure 5.** Box-plot association *P* values between *LCT* (rs748841) and *Bifidobacterium*  
148 *adolescentis* in IBD cohort and Lifelines-DEEP cohort respectively. X-axis indicates the genotype of rs748841

149 and Y-axis indicates the relative abundance of *Bifidobacterium adolescentis*.

150 *P*, *P* values of associations; IBD, IBD cohort; LLD, LifeLines-DEEP cohort