# Population Structure Discovery in Meta-Analyzed Microbial Communities and Inflammatory Bowel Disease

3  Siyuan Ma[1,2], Dmitry Shungin[2], Himel Mallick[1,2], Melanie Schirmer[2], Long H. Nguyen[3],

4  Raivo Kolde[2], Eric Franzosa[1,2], Hera Vlamakis[2], Ramnik Xavier[2*], Curtis Huttenhower[1,2*]

5 [1] Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA

6 [2] Broad Institute of MIT and Harvard, Cambridge, MA, USA

7 [3] Massachusetts General Hospital, Boston, MA, USA

## Abstract

Microbial community studies in general, and of the human microbiome in inflammatory bowel disease (IBD) in particular, have now achieved a scale at which it is practical to associate features of the microbiome with environmental exposures and health outcomes across multiple large-scale populations. This permits the development of rigorous meta-analysis methods, of particular importance in IBD as a means by which the heterogeneity of disease etiology and treatment response might be explained. We have thus developed MMUPHin (Meta-analysis Methods with a Uniform Pipeline for Heterogeneity in microbiome studies) for joint normalization, meta-analysis, and population structure discovery using microbial community taxonomic and functional profiles. Applying this method to ten IBD cohorts (5,151 total samples), we identified a single consistent axis of microbial associations among studies, including newly associated taxa such as *Acinetobacter* and *Turicibacter* detected due to the sensitivity of meta-analysis. Linear random effects models further revealed associations with medications, disease location, and interaction effects consistent within and between studies. Finally, multiple unsupervised clustering metrics

22    and dissimilarity measures agreed on a lack of discrete microbiome "types" in the IBD gut

23    microbiome. These results thus provide a benchmark for consistent characterization of the IBD

24    gut microbiome and a general framework applicable to meta-analysis of any microbial community

25    types.

## Introduction

27    Meta-analysis for molecular epidemiology in large populations has seen great success in linking

28    high-dimensional 'omic features to complex health-related phenotypes. One example of this is in

29    genome-wide association studies (GWAS[1]), where the appropriate study scale, achieved by

30    rigorous integration of multiple cohorts, has both facilitated reproducible discoveries (in the form

31    of disease-associated loci[2-4]) and addressed confounding due to unobserved population

32    structure[5]. The inflammatory bowel diseases (IBD) represent a particular success story for GWAS

33    meta-analysis[3,4], and environmental and microbial contributors complementing the condition's

34    complex genetic architecture have been detailed by many individual studies[6-8]. However, in the

35    absence of methods appropriate for large-scale microbial meta-analysis, the extent to which these

36    findings reproduce across studies, or can be extended by increased joint sample sizes, remains

37    undetermined. Likewise, it is unclear whether reproducible population structure in the microbiome,

38    such as microbially-driven IBD "subtypes," exists to help explain the clinical heterogeneity of these

39    conditions[9].

40    Meta-analysis of microbial community profiles presents unique quantitative challenges relative to

41    other types of 'omics data such as GWAS[10] or gene expression[11]. These include particularly

42    strong batch, inter-individual, and inter-population differences, and statistical issues including

43    zero-inflation and compositionality[12,13]. Consequently, methods to correct for cohort and batch

44    effects from other 'omics settings[14-17] are not directly appropriate. Two recent studies have

45    suggested quantile normalization[18] and Bayesian Dirichlet-multinomial regression (BDMMA)[19] for

46    microbial profiles, which are applicable to a limited subset of differential abundance tests and do

47    not provide batch-corrected profiles. To date, there are no methods permitting the joint analysis

48    of batch-corrected microbial profiles for most study designs.

49    IBD represents one of the best-studied, microbiome-linked inflammatory phenotypes to date

50    which thus stands to benefit from such approaches[20,21]. Among the inflammatory bowel diseases,

51    Crohn's disease (CD) and ulcerative colitis (UC) have been individually linked with structural and

52    functional changes in the gut microbiome in many individual studies[21]. Each of CD and UC can

53    itself be highly heterogeneous within the IBD population, however, and diversity in disease-

54    associated gut microbial features has not been consistently associated with factors including

55    disease subtype, progression, or treatment response[7,9,22,23]. Of note, two meta-analysis studies

56    included IBD as one of several phenotypes[24,25]. These studies were not IBD-specific, did not have

57    access to appropriate normalization techniques, nor took the aforementioned factors into account.

58    The complexity of microbial involvement in IBD, and the presence of substantial unexplained

59    variation in the manifestation of its symptoms, makes it particularly appropriate for application of

60    meta-analysis techniques.

61    In this work, we introduce and validate a statistical framework for population-scale meta-analysis

62    of microbiome data, and apply it to the largest collection to date of ten published 16S rRNA gene

63    sequencing-based IBD studies (**Table 1**) to identify consistent disease associations and

64    population structure. We found both previously documented and novel microbial links to the

65    disease, with further differentiation among subtypes, phenotypic severity, and treatment effects.

66    We further confidently conclude that there are no apparent, reproducible microbiome-based

67    subtypes within CD or UC, which are instead a population structure gradient from less to more

68    "pro-inflammatory" ecological configurations. Our work thus represents one of the first large-scale

69 efforts to assesses consistency in gut microbial findings for IBD and provides methodology

70 supporting future microbial community meta-analyses.


## Results


**Integrating 10 studies of the IBD stool and mucosal microbiomes**

73 We collected and uniformly processed ten published 16S studies of the IBD gut microbiome

74 (**Table 1, Fig. 1a, Supplemental Table 1**) totaling 2,179 subjects and 5,151 samples. These

75 studies range widely in terms of cohort designs and population characteristics, including recent-

76 onset and established disease patients, cross-sectional and longitudinal sampling, pediatric and

77 adult populations, diseases (CD and UC), treated and treatment-naive patients, biopsy and stool

78 samples, and inclusion of healthy/non-IBD controls. Covariates were manually curated to ensure

79 consistency across studies (**Methods**). Major factors available from all or most studies included

80 demographics (age/sex/race), biogeography, disease location and/or extent, antibiotic usage,

81 immunosuppression, and steroid and/or 5-ASA usage.

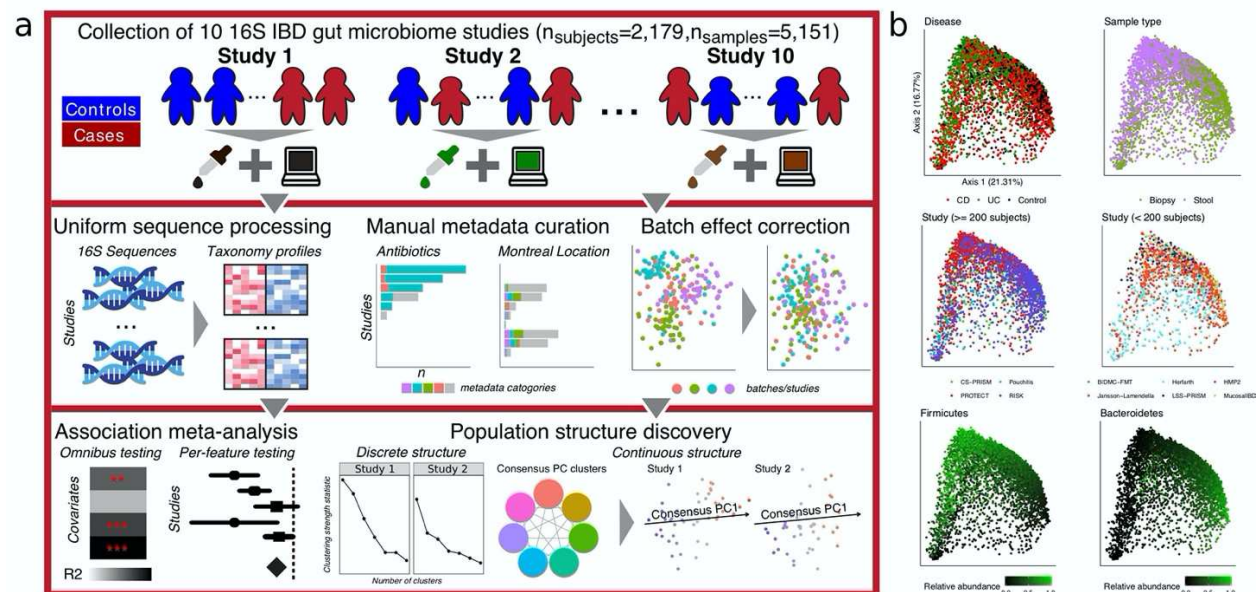

82

83 **Figure 1: A method for large-scale microbial community meta-analysis and its application to inflammatory**

84 **bowel disease. a)** We developed a novel statistical framework, MMUPHin, allowing joint normalization and meta-

85 analysis of large microbial community profile collections with heterogeneous and complex designs (multiple covariates,

86 longitudinal samples, etc.). We applied it to a collection of 10 inflammatory bowel disease studies comprising 2,179

87 subjects and 5,151 total samples (**Table 1**). We uniformly processed the associated sequence data and harmonized

88 metadata across cohorts. Microbial taxonomic profiles were then corrected for batch- and study-effects before

89 downstream analyses for omnibus and per-feature association with disease phenotypes and unsupervised population

90 structure discovery. **b)** MDS ordination of all microbial profiles (Bray-Curtis dissimilarity) before batch correction

91 visualize the strongest associations with gut microbial composition, including disease, sample type (biopsy or stool),

92 cohort (visualized separately for larger and smaller studies), and dominant phyla.

93 Using this joint dataset and upon uniform bioinformatics processing (**Methods**), we first assessed

94 the factors that corresponded to overall variation in microbiome structure, which included disease

95 status, sample type (biopsy versus stool), and dominant phyla (Bacteroidetes and Firmicutes, **Fig.**

96 **1b**). Cohort effects prior to batch correction and meta-analysis were also significant. Microbiome

97 differences associated with disease were notable even without normalization. However, this can

98 be misleading due to the confounding of cohort structure between studies, such as the

99 differentiation between RISK (a predominantly mucosal study of CD) and PROTECT (a

100 predominantly stool study of UC). Inter-individual differences largely independent of population or

101 disease, such as Bacteroidetes versus Firmicutes dominance, were also universal among studies

102 and sample types as expected[9,26]. Many of these factors were of comparable effect size, both

103 visually and as quantified below, emphasizing the need for covariate-adjusted statistical modelling

104 to delineate the biological (disease, treatment) and technical (cohort, batch) effects associated

105 with individual taxa throughout the cohorts (**Supplemental Notes**, **Supplemental Fig. 1-3**).

| Study | Brief description | N subject | N sample | Phenotype(s) | Age | Gender | Sample type(s) |
|---|---|---|---|---|---|---|---|
| PROTECT [23] | Longitudinal cohort of newly diagnosed UC | 405 | 1212 (539) | UC 405 | 12.71 (3.29) | Male 52%/ Female 48% | Biopsy 22%/ Stool 78% |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| RISK[7] | Pediatric cohort of treatment-naïve CD | 631 | 882 | CD 430/ Control 201 | 12.16 (3.22) | Male 59%/ Female 41% | Biopsy 72%/ Stool 28% |
| Herfarth[27] | Densely (daily) sampled longitudinal cohort | 31 | 860 (31) | CD 19/ Control 12 | 36.03 (14.12) | Male 35%/ Female 58%/ Missing 6% | Stool |
| Jansson-Lamendella[22] | Longitudinal follow up with fecal samples | 137 | 683 (137) | CD 49/ UC 60/ Control 28 | | Male 42%/ Female 58% | Stool |
| Pouchitis[28] | Patients recruited underwent IPAA for treatment of UC or FAP prior to enrollment. | 353 | 577 | CD 42/ UC 266/ Control 45 | 46.19 (13.58) | Male 52%/ Female 48% | Biopsy |
| CS-PRISM[29] | Cross sectional cohort nested in PRISM | 397 | 467 | CD 215/ UC 144/ Control 38 | 41.68 (15.22) | Male 47%/ Female 53% | Biopsy 29%/ Stool 71% |
| HMP2[9] | Large cohort of newly diagnosed IBD with multi 'omics measurement. | 81 | 177 (162) | CD 37/ UC 22/ Control 22 | 29.76 (19.63) | Male 51%/ Female 49% | Biopsy |
| MucosalIBD[30] | Pediatric cohort with Paneth cell phenotypes | 83 | 132 | CD 36/ Control 47 | 12.93 (3.65) | Male 58%/ Female 42% | Biopsy |
| LSS-PRISM[31] | Longitudinal cohort nested in PRISM. | 18 | 88 (19) | CD 12/ UC 6 | 30.37 (10.52) | Male 39%/ Female 61% | Stool |
| BIDMC-FMT[32] | FMT Trial design | 8 | 16 | CD 8 | 38.38 (12.73) | Male 62%/ Female 38% | Stool |

106 **Table 1: 10 uniformly processed 16S rRNA gene sequencing studies of the IBD mucosal/stool microbiomes.**

107 For longitudinal cohorts, numbers in parentheses indicate baseline sample size. For age, mean and standard error

108 (parenthesized) are shown. Additional covariates are summarized in **Supplemental Table 1**.


109 **A statistical framework for meta-analysis of microbial community profiles**

110 We developed a collection of novel methods for meta-analysis of environmental exposures,

111 phenotypes, and population structures across microbial community studies, specifically

112 accounting for technical batch effects and interstudy differences (**Methods**, **Fig. 1a**). It consists

113 of three main components: batch and study effect correction, covariate modeling, and population

114 structure discovery. First, we extended methods from the gene expression literature (ComBat[15])

115 to enable batch correction of zero-inflated microbial abundance data. Based on linear modelling,

116 the method can differentiate between technical effects (batch, study) versus covariates of

117 biologically interest (exposure, phenotype). Second, we combined well-validated data

118 transformation and linear modelling combinations for microbial community profiles[33] with fixed and

119 random effect modelling[34] for meta-analytical synthesis of per-feature (taxon, gene, or pathway)

120 differential abundance effects. Lastly, we generalized and formalized approaches from cancer

121 transcriptional subtyping[35] to permit unsupervised discovery and validation of both discrete and

122 continuous population structures in microbial community data (**Supplemental Fig. 4**). Our

123 methods, implemented as Meta-analysis Methods with a Uniform Pipeline for Heterogeneity in

124 microbiome studies (MMUPHin), are available as an R package through Bioconductor[36] and at

125 https://bioconductor.org/packages/release/bioc/html/MMUPHin.html.

126 We validated MMUPHin both in comparison to existing methods and through extensive simulation

127 studies (**Fig. 2**), with simulated realistic microbial abundance profiles at different data

128 dimensionality, biological/technical batch signal strength, and discrete/continuous population

129 structures (**Methods, Supplemental Table 2**, **Supplemental Fig. 5-8**). MMUPHin successfully

130 reduced variability attributable to technical effects in simulated microbial profiles, as first quantified

131     by the PERMANOVA R2 statistic[37] (**Fig. 2a-b, Supplemental Fig. 5**). This was true both in terms

132     of reducing the overall microbial variability attributable to technical artifacts and in terms of the

133     ratio of "biological" versus technical variability (**Fig. 2a**). ComBat correction[15], suited for gene

134     expression data, was capable of reducing batch effects to a lesser degree, but also tended to

135     reduce desirable "biological" variation in the process, likely due to noise introduced by it changing

136     many zero counts to non-zero values. Previously proposed techniques for microbial community

137     data, namely quantile normalization[18] and   BDMMA[19], are only appropriate for differential

138     abundance analysis and do not provide batch-normalized profiles, thus precluding PERMANOVA

139     batch effect quantification; their per-feature testing performance is evaluated together with

140     MMUPHin in the following section. MMUPHin thus provides batch-corrected microbial community

141     profiles that retain biologically meaningful variation more than (or not even possible using) existing

142     methods.

143     For differential abundance testing, MMUPHin successfully corrected for false associations when

144     batch/cohort effects were confounded with variables of interest, which is a common concern for

145     'omics meta-analysis[38], while quantile normalization[18] and BDMMA[19] had either inflated or overly

146     conservative false positive rates (**Fig. 2c-d, Supplemental Fig. 6**). We also validated MMUPHin's

147     support for unsupervised population structure discovery, in addition to these "supervised"

148     differential abundance and statistical association tests. In microbial communities, valid,

149     generalizable population structure can manifest as either discretely clustered subtypes[39] or as

150     continuously variable gradients of community configurations[40], but methods for discovery are

151     particularly susceptible to false positives in the presence of technical artifacts[26,40]. To this end, for

152     discrete structures, MMUPHin utilizes established clustering strength evaluation metrics[41] to a)

153     evaluate the existence of discrete clusters within individual microbiome studies and b) to validate

154     the reproducibility of such structures among studies meta-analytically (**Fig. 2e-f, Supplemental**

155     **Fig. 7**). For continuous structures, our method generalizes single study principal component

156 analysis (PCA[42]) to multiple studies by constructing a network of correlated top PC loadings[35],

157 thus identifying major axes of variation that explain the largest amount of heterogeneity between

158 microbial profiles and are also consistent across studies (**Fig. 2g-h, Supplemental Fig. 8**). As a

159 result, MMUPHin was able to successfully identify discrete clusters (i.e. microbiome "types") when

160 present, as well as significantly consistent continuous patterns of microbiome variation that recur
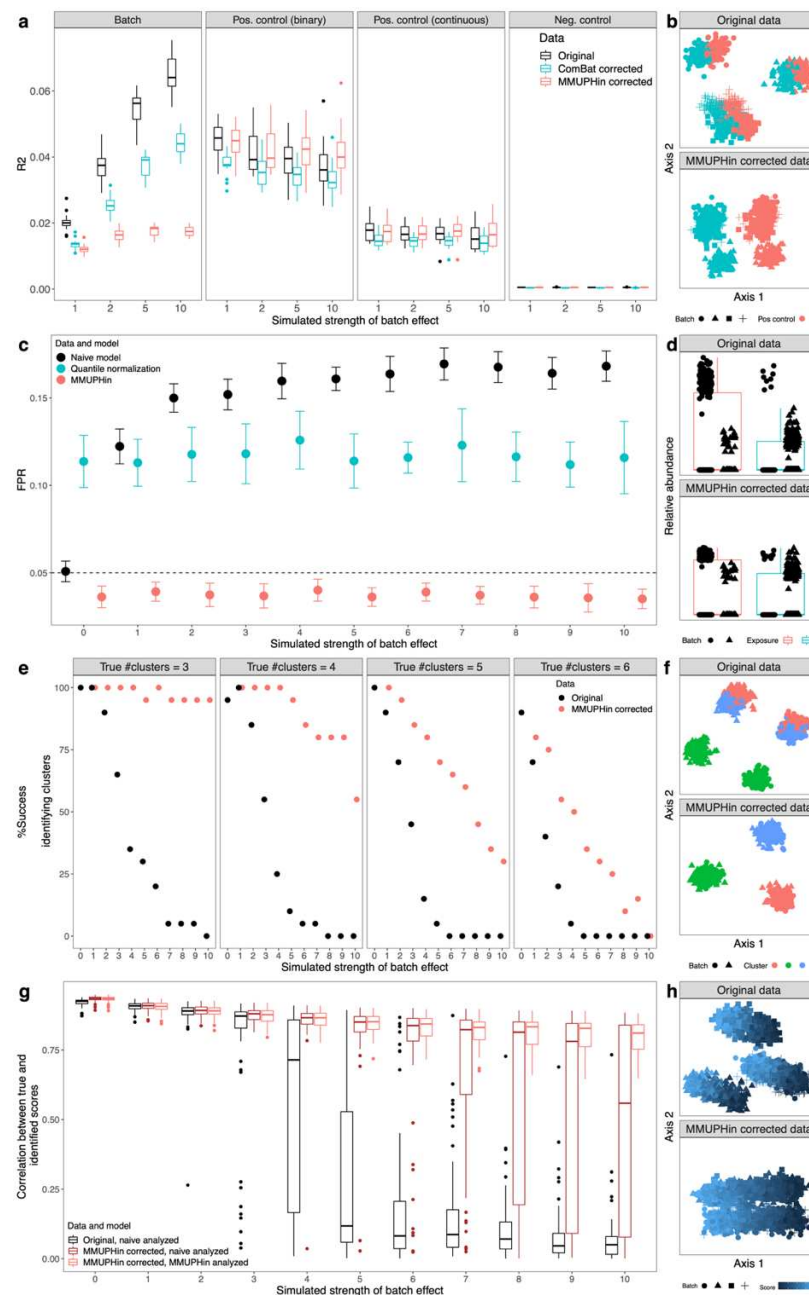
161 among populations (**Supplemental Notes**).

162

163    **Figure 2: Effectiveness of batch correction, association meta-analysis, and unsupervised population structure**

164    **discovery methods.** All evaluations use simulated microbial community profiles as detailed in Methods. Left panels

165    summarize representative subsets of results (full set of simulation cases presented in **Supplemental Table 2** and

166    results in **Supplemental Fig. 5-8**), and right panels show examples of batch-influenced data pre- and post-correction.

167    a, b) MMUPHin is effective for covariate-adjusted batch effect reduction while maintaining the effect of positive control

168    variables. Results shown correspond to the subset of details in **Supplemental Fig. 5** with number of samples per batch

169    = 500, number of batches = 4, and number of features = 1000 with 5% spiked with associations. c, d) Batch correction

170    and meta-analysis reduces false positives when an exposure is spuriously associated with microbiome features due to

171    an imbalanced distribution between batches. Corresponds to **Supplemental Fig. 6** with number of samples per batch

172    = 500, number of features = 1000 with 5% spiked associations, and case proportion difference between batches = 0.8.

173    Evaluations of BDMMA generates low FPRs due to the zero-inflated nature of simulated microbial abundances, and

174    are included only in **Supplemental Fig. 6**. e, f) Batch correction improves correct identification of the true underlying

175    number of clusters during discrete population structure discovery. Corresponds to **Supplemental Fig. 7** with number

176    of batches = 4. g, h) Continuous structure discovery accurately recovers microbiome compositional gradients in a

177    simulated population. Corresponds to **Supplemental Fig. 8** with number of batches = 6.


178    **Meta-analysis of the IBD microbiome**


179    Given these validations of MMUPHin's accuracy in simulated data, we next applied it to the 10-

180    study, 4,789-sample IBD gut amplicon profile meta-analysis introduced above (**Fig. 3**). MMUPHin

181    successfully reduced the effects both of differences among studies, and of batches within studies

182    (study effect correction modelling disease and sample type as covariates, see **Methods**),

183    although these remained among the strongest source of variation among taxonomic profiles as

184    quantified by PERMANOVA R2 (**Fig. 3a**, **Methods**, **Supplemental Table 3**). Among biological

185    variables, sample type (biopsy/stool), biopsy location (multiple, conditional on biopsy samples),

186    disease status (IBD/control), and disease types (CD/UC, conditional on IBD) consistently had the

187    strongest effect on the microbiome among studies. Several relationships between study design

188    and phenotypic effects were apparent. Batches had a particularly strong effect in CS-PRISM and

189    RISK, for example, where biopsy and stool samples were also perfectly separated by batch.

190  Treatment exposures all had small effects on microbiome structure within studies, which typically

191  reached statistical significance only when combined by meta-analysis; antibiotics were an

192  exception with slightly larger effects. Montreal classification did not generally correspond with

193  significant variation, while age (at sample collection as stratified below and above 18, and at

194  diagnosis by Montreal age classification[43]) had small but significant effects. The effects of gender

195  and race were not significant. Lastly, for longitudinal studies, relatively stable differences between

196  subjects over time were large and significant, consistently for both longer-interval (HMP2) as well

197  as densely sampled cohorts (Herfarth, daily samples), in agreement with previous individual
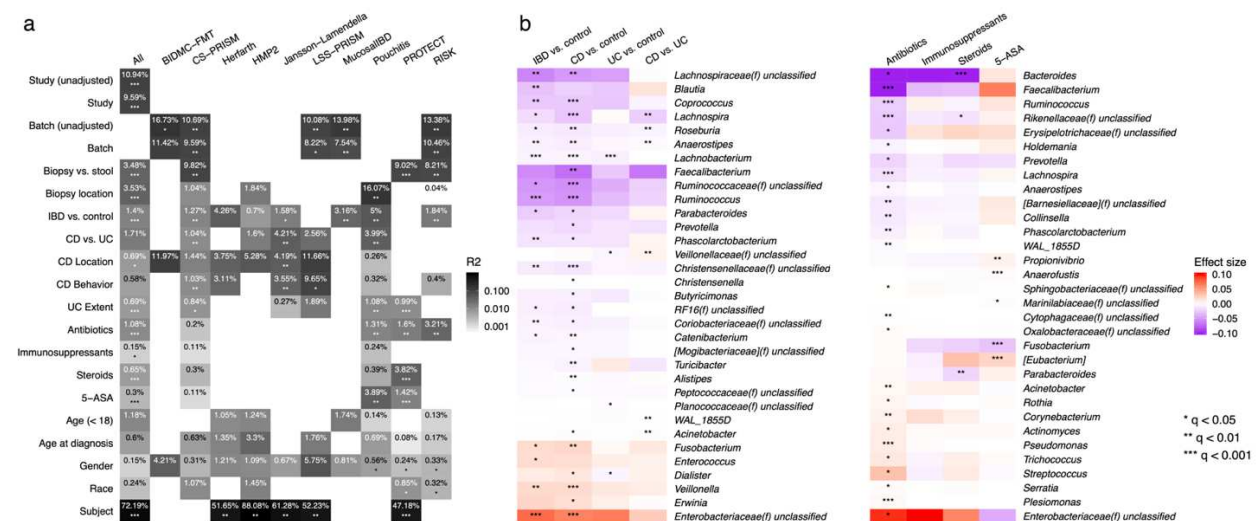
198  studies' observations[9,23].

199



200  **Figure 3: Meta-analytic omnibus and per-feature testing reveal novel and previously documented IBD**

201  **associations. a)** Omnibus testing (PERMANOVA on Bray-Curtis dissimilarities with stratification and covariate control

202  where appropriate, see **Methods** and **Supplemental Table 3**) identified between-subject differences as the greatest

203  source of microbiome variability, with IBD phenotype, disease (CD/UC), and sample type (stool/biopsy) as additional

204  main sources of biological variation. MMUPHin successfully reduced between-cohort and within-study batch effects,

205  although these technical sources also remained significant contributors to variability. **b)** Individual taxa significantly

206  associated with IBD phenotypes or treatments after meta-analysis. Taxa are arranged by family-level median effect

207  size of IBD vs. control for disease results and that of  antibiotic usage for treatment results. Effect sizes are aggregated

208  regression coefficients (across studies with random effects modelling) on arcsin square root-transformed relative

209   abundances. Detailed model information in **Methods** and **Supplemental Table 3**. Individual study results in

210   **Supplemental Table 4**.

211   We identified individual taxonomic features consistently associated with disease and treatment

212   variables (**Fig. 3b**, **Supplemental Table 4**), with meta-analysis multivariate differential

213   abundance analysis, adjusting for common demographics (age, gender, race) and further

214   stratifying for sample type and disease when appropriate (**Methods**, **Supplemental Table 3**). At

215   a very high level, differential abundance patterns between CD and control microbiomes were

216   consistent with, and often more severe than contrasts between UC and control, confirming with

217   increased resolution previous observations that CD patients tend to have more aggravated

218   dysbiosis than UC patients[9]. As expected, our meta-analysis confirms many of the taxa

219   associated with IBD reported by previous individual (**Fig. 3b**, detailed in **Supplemental Notes**);

220   these findings strongly supports the emerging hypotheses of pro-inflammatory aerotolerant

221   clades forming a positive feedback loop in the gut during inflammation, often of oral origin[7], and

222   depleting the gut's typical fastidious anaerobe population as a result.

223   We also identified two taxa not previously associated with IBD, both of modest effect sizes and

224   likely newly detected by the meta-analysis' increased power. The genus *Acinetobacter* was

225   enriched in CD, and *Turicibacter* was depleted. *Turicibater* in particular is poorly represented in

226   reference sequence databases, with only nine genomes for one species (*Turicibacter sanguinis*)

227   currently in the NCBI genome database; this makes it easy to overlook in shotgun metagenomic

228   profiles relative to amplicon sequencing. The genus *Acinetobacter*, conversely, is quite well

229   characterized due to its role in antimicrobial resistant infections[44], and it was previously linked

230   specifically to the primary sclerosing cholangitis phenotype in UC[45], although without follow-up to

231   our knowledge. *Turicibacter* is overall less characterized both in isolation and with respect to

232   disease, although our findings and others' suggest it might be inflammation-sensitive when

233   present; it was one of many clades increased in mice during CD8+ T cell depletion[46] and reduced

234 in a homozygous TNF deletion[47]. As the strains of *Acinetobacter* implicated in gut inflammation

235 are unlikely to be those responsible for e.g. nosocomial infections, further investigation of both

236 clades using more detailed data or IBD-specific isolates is warranted.

237 Among treatment variables (samples or time points during which subjects were receiving

238 antibiotics, immunosuppressants, steroids, and/or 5-ASAs), antibiotics had the strongest effects

239 on individual taxa, as well as the greatest number of significantly associated taxa (**Fig. 3b**). These

240 associations are also broadly in agreement with previous observations for microbiome responses

241 to antibiotics in IBD or generally, e.g. the depletion of *Faecalibacterium*, *Ruminococcus*, and

242 *Bacteroides* in patients treated with antibiotics, and the enrichment of (often stereotypically

243 resistant) taxa such as *Streptococcus*, *Acinetobacter*, and the Enterobacteriaceae, with

244 differential responses to the treatment groups speaking to both administration considerations and

245 their impact on host versus microbial community bioactivities (**Supplemental Notes**).

246 Subsets of IBD-linked taxa were additionally associated with the diseases' phenotypic severity

247 (**Fig. 4a, Supplemental Table 5**). Montreal classification[43] was used as a proxy for disease

248 severity, including Behavior categories for Crohn's disease (B1 non-stricturing, non-penetrating,

249 B2 stricturing, non-penetrating, B3 stricturing and penetrating) and Extent for ulcerative colitis (E1

250 limited to rectum, E2 up to descending colon, E3 pancolitis). We tested for features differentially

251 abundant in the more severe phenotypes when compared against the least severe category (B1

252 CD and E1 UC, **Methods**). Among statistically significant results, many extended those identified

253 above as overall IBD associated (**Fig. 3b**), such as the depletion of *Faecalibacterium* in B3 CD

254 and *Roseburia* in B2 CD, as well as the enrichment of Enterobacteriaceae in E3 UC. In most

255 cases, microbial dysbiosis was also additionally aggravated from the moderate to the most

256 extreme disease manifestations; such differences were statistically significant (**Methods**) in, for

257 example, the progressive depletion of *Bacteroides* in CD and UC, as well as the enrichment of

258 Enterobacteriaceae in UC. This meta-analysis is uniquely powered to detect these subtle

259    differences, which aid in shedding light on the microbiome's response to progressive inflammation

260    and disease subtypes. Pancolitis corresponds with a unique microbial configuration distinct from

261    regional colitis and not generally detectable in smaller studies[6], for example, while more severe

262    CD induces essentially a more extreme form of the same dysbiosis observed in less severe forms

263    of the disease.



264

265    **Figure 4: IBD-associated taxa are aggravated in more severe disease; disease biogeography and CD/UC**

266    **differentially affect some taxa with respect to disease and treatment. a)** Statistically significant genera from meta-

267    analytically synthesized differential abundance effects among severity of CD and UC phenotypes as quantified by

268    Montreal classification. The difference between the most severe phenotype with the least severe one (B3 vs. B1 for

269    CD, E3 vs. E1 for UC) was in most cases more aggravated than that of the intermediate phenotype. Many of the

270    identified features overlap with those associated with IBD vs. control differences, suggesting a consistent gradient of

271    severity effects on the microbiome. Individual study results in **Supplemental Table 5**. **b)** Genus *Dehalobacterium* as

272    an example in which a taxon is uniquely affected in the stool microbiome during CD and not at the mucosa. Likewise,

273    family Enterobacteriaceae as an example in which steroid treatment corresponds with enrichment of the clade in CD

274    samples, but depletion in UC. In all panels, effect sizes are aggregated regression coefficients on arcsin square root-

275    transformed relative abundances. Full sets of statistically significant interactions, with individual study results, are in

276    **Supplemental Table 6**.

277    Additionally, diseases (CD and UC) and their corresponding dysbioses also interacted distinctly

278    with the microbiome under different treatment regimes and in different biogeographical

279    environments (mucosa vs. stool, **Fig. 4b, Supplemental Table 6**). Interaction effects, in the

280    statistical sense, were defined as a main exposure (IBD or treatment) having differential effects

281    on taxon abundance with respect to either sample type (biopsy/stool) or diseases (CD/UC); they

282    were identified via moderator meta-analysis models (**Methods**). Overall, we found elevated

283    effects of both CD (relative to controls) and antibiotic treatment in stool as compared to biopsy-

284    based measurements of the microbiome (**Supplemental Table 6**). An example of this is

285    *Dehalobacterium*, with significantly greater depletion in CD stool relative to biopsies (**Fig. 4b**).

286    *Dehalobacterium*, as with *Turicibacter* above, is underrepresented in reference sequence

287    databases, better-detected by amplicon sequencing, and thus not a common microbial signature

288    of IBD. It has been linked to CD in at least one existing 16S-based stool study[48]. In contrast,

289    several UC-specific microbial disruptions were more prominent at the mucosa (i.e. in biopsies,

290    **Supplemental Table 6**). Coupled with the severity-linked differences above, this suggests CD-

291    induced changes in the entire gut microbial ecosystem largely as a consequence of inflammation,

292    with UC-induced dysbioses both more local and more specific to disease and treatment regime.

293    Additional results include effect of steroids on the Enterobacteriaceae, which tended to be more

294    abundant in CD patients receiving steroids, but less abundant in UC recipients (**Fig. 4b,**

295    **Supplemental Table 6, Supplemental Notes**).

296    **Consistent IBD microbial population structure discovered by unsupervised analysis**

297    The existence of subtypes within gut microbial communities has been a major open question in

298    human microbiome studies, and it is of particular importance within IBD as a potential explanation

299  for heterogeneity in disease etiology and treatment response[6,9]. To systematically characterize

300  population structure in the IBD gut microbiome that was reproducible among studies, we

301  performed both discrete and continuous structure discovery on the 10 cohorts using our meta-

302  analysis framework. To identify potential discrete community types (i.e. clusters), we performed

303  clustering analysis within each cohort's IBD patient population, and evaluated the clustering

304  strength via prediction strength (**Methods**). We found no evidence to support discrete clustering

305  structure within individual cohorts, nor were we able to reproduce each cohort's clustering results

306  externally (**Fig. 5a**). This lack of discrete structure was consistent when we further stratified

307  samples to either CD or UC populations (**Supplemental Fig. 9**), or extended to additional

308  dissimilarity metric and clustering strength measurements (**Supplemental Fig. 9, Methods**). Our

309  observation that the IBD gut microbiome cannot be well characterized by discrete clusters is thus

310  consistent with previous findings on gut microbial heterogeneity for healthy populations[40] and

311  suggests that, at the level powered by this study, such microbiome subtypes are not clearly

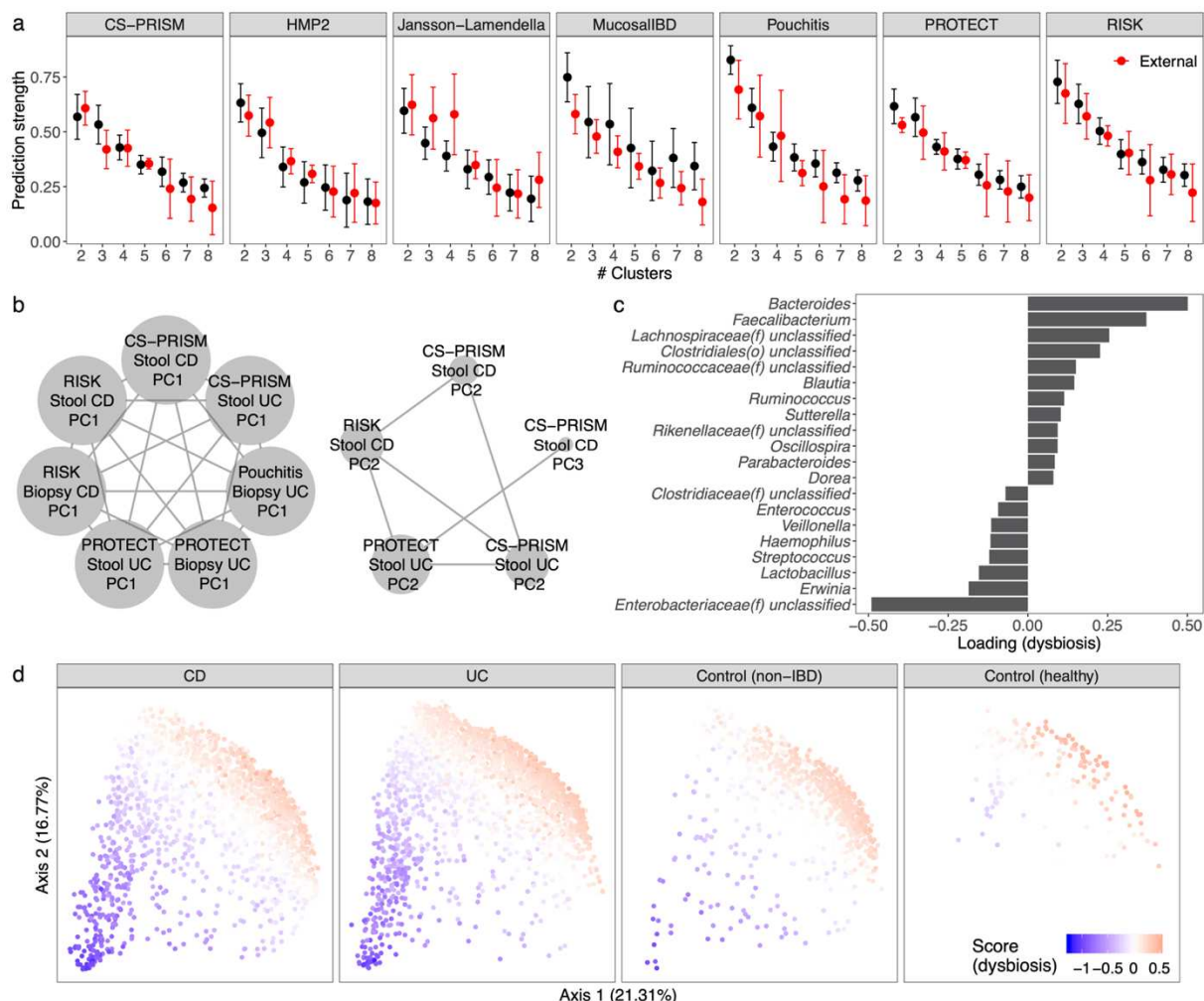312  responsible for clinical heterogeneity.

**Figure 5: Unsupervised population structure discovery finds no evidence of microbiome-based subtypes in the IBD gut, but a reproducible gradient of continuously variable dysbiosis in disease. a)** No support was detected for discrete microbiome subtypes (clusters) within the IBD microbiome, neither within cohort nor when evaluated among studies (red bars) using prediction strength[41]. This remained true during stratification within CD and UC, and for additional dissimilarity metric/clustering strength measurements (**Supplemental Fig. 9**). **b)** Conversely, two reproducible, continuously variable patterns of microbiome population structure were identified using groups of similar principal components (**Methods**)[35]. These patterns were consistent within and between cohorts, disease types, and sample types, as well as under different edge strength cutoffs (**Supplemental Fig. 11**), and their consensus loadings were reproducible among cohorts (**Supplemental Fig. 12**). **c)** Top 20 genera with highest absolute loadings for the disease-associated dysbiosis score corresponding to the first cluster in **b**. Many of these taxa were also IBD-associated (**Fig. 3b**). **d)** Distribution of the dysbiosis pattern across CD, UC, non-IBD control, and healthy populations.

325  Although it was defined in an unsupervised way solely within the IBD population, across which the pattern is highly

326  variable, it also differentiates well between IBD and control populations (**Supplemental Fig. 13**).

327  Conversely, we identified two consistent, continuously varying gradients of microbial community

328  variation in the IBD microbiome (**Fig. 5b-d**, **Supplemental Fig. 10**). These gradients represent

329  patterns of microbes that occur with greater or lesser abundance in tandem, and which covary

330  across subjects in a population; they were identified as principal component (PC) vectors that

331  recur among different cohorts (see **Methods**)[35]. Briefly, we used the four largest IBD cohorts (CS-

332  PRISM, Pouchitis, PROTECT, and RISK) as training datasets to identify two clusters of consistent

333  PCs (**Fig. 5b**), which were confirmed with sensitivity analysis (**Supplemental Fig. 11**) and

334  validated in the remaining cohorts (**Supplemental Fig. 12**). The consensus loadings (i.e. within-

335  cluster average) representing these two clusters (**Fig. 5c**, **Supplemental Fig. 10**, **Supplemental**

336  **Table 7**) were used to assign continuously varying scores to the IBD population that capture

337  gradient changes in the microbiome that occurred consistently within IBD, across diseases,

338  sample types, and cohorts. This disease-linked "type" of microbiome variation corresponded

339  roughly to severity or extent of inflammation, as detailed below.

340  In particular, while the second continuous population structure captured the Firmicutes-

341  Bacteroidetes tradeoff present in most gut microbiome studies (**Supplemental Fig. 10**)[9,26,40], the

342  first continuous score was IBD-specific and corresponded roughly to more extreme disease-

343  associated dysbiosis in CD and UC populations (**Fig. 5d**). This is evidenced by the taxa with

344  highest weights in the scores' consensus loading vector (**Fig. 5c**), which included taxa

345  differentially abundant between IBD and control populations (**Fig. 3**). The score was consistent

346  both within CD and UC while also further differentiating IBD, non-IBD control, and healthy

347  populations (**Fig. 5d**, **Supplemental Fig. 13**), even though it was identified unsupervisedly only

348  from diseased subsets. The composition of the score and its population structure are also

349  consistent with our recent definition of dysbiotic gut microbiome configurations corresponding with

350    multi'omic perturbations during IBD activity[9]. Together with the supervised meta-analysis results

351    above, these unsupervised population structure findings confirm that there are no detectable

352    discrete subtypes of the gut microbiome in IBD even among ~5,000 combined samples, while

353    showing a single continuously variable gradient of microbiome changes reproducibly present

354    during more dysbiotic diseases.


355    **Discussion**


356    Here, we provide a novel framework for microbial community meta-analysis and apply it to the

357    first large-scale integration of over 5,100 amplicon profiles of the stool and mucosal microbiomes

358    in IBD. This identified a significantly reproducible gradient in the gut microbiome indicative of

359    increasing dysbiosis in subsets of patients. The study also showed no evidence of additional

360    population structure, such as microbiome-driven discrete disease subtypes, within CD or UC. The

361    increased power provided by meta-analysis supported many of the taxonomic associations

362    previously ascribed to IBD (e.g. *Faecalibacterium*, *Ruminococcus*, Enterobacteriaceae) while

363    uncovering new associations (*Turicibacter*, *Acinetobacter*) not confidently associated with

364    inflammation by other populations or data types. Almost all effects were exhibited similarly using

365    either stool or mucosal profiling, with a small number of exceptions showing significant

366    differentiation (e.g. *Dehalobacterium*). Novel disease-treatment response interactions were

367    observed (e.g. steroids on Enterobacteriaceae). Finally, the meta-analysis framework developed

368    for the study, MMUPHin, has been extensively evaluated and its performance for batch effect

369    removal, supervised meta-analysis of exposures and covariates, and unsupervised population

370    structure discovery validated on a variety of simulated microbial community types. It is extensible

371    to integration of microbial community taxonomic or functional profiles from other data types (e.g.

372    metagenomic sequencing) or environments.

373    However, all microbial community meta-analyses should be approached with caution, since in

374    many cases unwanted sources of technical variation between studies (i.e. batch effects) are so

375    large as to potentially mask biological signals even after correction[49-51] (**Supplemental Notes**).

376    Reducing inter-study variation in microbial community profiles is challenging relative to other

377    'omics data types due to 1) the extreme heterogeneity of microbes within most communities

378    (exacerbating both technical and biological differences), and 2) feature zero-inflation arising from

379    both biological and technical reasons[13,52]. Notably, despite these challenges, MMUPHin was able

380    to meta-analyze amplicon profiles in this study both to associate microbial shifts with disease

381    outcome, to associate them with treatment-specific differences, and to identify a single pattern of

382    typical microbial variation within IBD. While previous efforts have developed IBD dysbiosis scores

383    by contrasting patients with control groups[7,9], this pattern of microbial variation was present

384    specifically within IBD patients (both CD and UC), and in agreement with supervised methods,

385    captured several classes of microbial functional responses in the gut (**Supplemental Note**).

386    The IBD gut microbiome particularly stands to benefit from meta-analysis, as have other multiply-

387    sampled conditions such as colorectal cancer[53,54], in order to identify ecological and

388    microbiological changes during the disease that are reproducible across populations. We consider

389    this study based on 16S rRNA gene sequencing to be a proof of concept, able to achieve

390    unprecedented power due to the number of amplicon profiled samples available, but with greater

391    precision possible in future work using e.g. metagenomic and other 'omics technologies. This also

392    enabled comparison of responses in the stool versus mucosal microbiomes, the latter of which

393    are not amenable to metagenomic profiling from biopsies; these were in overall good agreement,

394    but the few areas of significantly differential responses to inflammation are likely of particular

395    immunological interest. The large sample and population sizes also provide some confidence in

396    ruling out discrete, microbially-driven population subtypes as an explanation for CD and UCs'

397    clinical heterogeneity. Instead, the work identified a single consistent axis of gradient microbial

398    change corresponding to increasing departures from "normal" microbiome configurations[7,9,55].

399    This pattern of consistent microbial dysbiosis can continue to be explored in further work on its

400    functional, immunological, and clinical consequences. Overall, this study represents one of the

401    first large-scale, methodologically appropriate, targeted meta-analysis of the IBD microbiome, and

402    the corresponding methodology and its implementation are freely available for future meta-

403    analyses of human-associated and environmental microbial populations.


404    **Methods**


405    **MMUPHin: a uniform statistical framework for meta-analysis of microbial community**

406    **studies**


407    We developed MMUPHin (Meta-analysis Methods with a Uniform Pipeline for Heterogeneity in

408    microbiome studies) as a framework for meta-analysis of microbial community studies using

409    taxonomic, functional, or other abundance profiles. It includes components for batch effect

410    adjustment, differential abundance testing, and unsupervised discrete and continuous population

411    structure discovery.


412    <u>Batch adjustment</u>


413    For microbial community batch correction, we extended the batch correction method developed

414    for gene expression data in ComBat[15] with an additional component to allow for the zero inflated

415    nature of microbial abundance data. In our model, sample read count $Y$ was modelled with respect

416    to both batch variable and biologically relevant covariate(s) $X$:


$$Y_{ijp} = exp\{\beta_p X_{ij}{}' + \sigma_p(\gamma_{ip} + \delta_{ip}\epsilon_{ijp})\} \times I_{ijp}$$

418     Where $i$ indicates batch/study, $j$ indicates sample, and $p$ indicates feature. $\gamma_{ip}$ and $\delta_{ip}$ are batch-

419     specific location and scale parameters. $\sigma_p$ is a feature-specific standardization factor. $\beta_p$ are

420     covariate-specific coefficients, and $\epsilon_{ijp}$ is an independent error term following a standard normal

421     distribution. $I_{ijp}$ is a binary (0, 1) zero-count indicator, to allow for zero inflation of features. As in

422     ComBat, $\gamma_{ip}$ and $\delta_{ip}$ are modelled with normal and inverse-gamma priors, respectively.

423     Hyperparameters are estimated with empirical Bayes estimators as in ComBat[15]. The posterior

424     means, $\widehat{\gamma^*}_{ip}$ and $\widehat{\delta^*}_{ip}$, along with standard frequentist estimates $\widehat{\beta_p}$ and $\widehat{\sigma_p}$ are used to provide

425     batch-corrected count data:

$$\widetilde{Y_{ijp}} = exp\{\frac{Y_{ijp} - \widehat{\beta_p}X_{ij}' - \widehat{\gamma^*}_{ip}\widehat{\sigma_p}}{\widehat{\delta^*}_{ip}} + \widehat{\beta_p}X_{ij}'\} \times I_{ijp}$$

426

427     Per-sample feature counts are then re-normalized to keep sample read depth unchanged post-

428     correction. In practice, the user provides sample microbial abundance table ($Y$), batch/study

429     information, and optionally any other covariates $X$ that are potentially confounded with batch but

430     encode important biological information. MMUPHin outputs an adjusted profile $\tilde{Y}$ that is corrected

431     for the effect of batches but retains the effects of $X$ (if provided).


432     <u>Meta-analysis differential abundance testing</u>


433     For meta-analytical differential abundance testing, after batch correction, MMUPHin first performs

434     multivariate linear regression within individual studies using previously validated data

435     transformation and modelling combinations appropriate for microbial community profiles

436     (MaAsLin2[33]). This yields study-specific, per-feature differential abundance effects estimations

437     $\widehat{\beta_{ip}}$, where $i$ indicates study and $p$ indicates feature. These are then aggregated into meta-

438     analysis effect size with fixed/random effects modelling as implemented in the metafor R

439     package[34]:

440
$$\widehat{\beta_{\iota p}} = \beta_p + \epsilon_{ip} + e_{ip}$$

441     $\beta_p$ is the overall differential abundance effect of feature $p$. $\epsilon_{ip}$ is per-study measurement error,

442     and $e_{ip}$ is study-specific random effects term (not present in fixed-effect models). In practice, the

443     user provides a microbial community profile, study design (batch) information, the main exposure

444     variable of interest, and optional additional covariates. If any meta-analyzed studies include

445     repeated measures (e.g. longitudinal designs), then random covariates can also be provided and

446     will be modelled for such studies. MMUPHin then performs MaAsLin2 regression modelling within

447     each study and aggregates effect sizes of the exposure variable $\widehat{\beta_{\iota p}}$ across studies using the

448     resulting random/fixed effects model. The estimated overall effect, $\widehat{\beta_p}$, is reported as the overall

449     differential abundance effect for feature $p$.

450     <u>Unsupervised discrete structure discovery</u>

451     For unsupervised discrete (i.e. cluster) structure discovery of a single study, again after batch

452     correction, MMUPHin uses average prediction strength[41], an established clustering strength

453     metric, to measure the existence of reproducible clusters among meta-analyzed datasets. Briefly,

454     for each individual dataset, the metric randomly and iteratively divides samples into "training" and

455     "validation" subsets. In each iteration, clustering is first performed on the training samples, across

456     a range of cluster numbers $k$, yielding (for a specific $k$) training sample clusters $A_{k1}, A_{k2}, ..., A_{kk}$.

457     Note that $A_{k1}, A_{k2}, ..., A_{kk}$ jointly forms a partition of the testing sample indices. The same

458     clustering analysis is then performed on the validation samples, and the resulting partition of

459     sample space provides classification membership potentially different from clustering

460     memberships $A_{k1}, A_{k2}, ..., A_{kk}$. Prediction strength for $k$ clusters is defined as

461    $ps(k)$

462    $= \min\limits_{1 \le l \le k} \frac{1}{n_{kl}(n_{kl} - 1)} \sum\limits_{j \ne j' \in A_{kl}} I\{\text{validation samples } j \text{ and } j' \text{ are classified to the same group according to training samp}$

463    i.e. the minimum (across validation clusters) proportion of same-cluster sample pairs also being

464    classified as the same group by training samples. $n_{kl} = |A_{kl}|$, or the number of test samples in

465    the $l$th cluster.

466    Average prediction strength is the average of prediction strengths across randomization iterations.

467    Intuitively, it characterizes the degree of agreement between the clustering structures in randomly

468    partitioned validation and training subsets; if $k$ is appropriately describing the true number of

469    discrete clusters in the dataset, then average prediction strength should be close to one (training

470    and validation samples agree most of the time).

471    We additionally generalized this metric to meta-analysis settings, where we aimed to quantify the

472    agreement of clustering structures between studies. In the meta-analytical setting, generalized

473    prediction strength for cluster number $k$ in study $i$ with validation study $i'$ is

474    $gps_{ii'}(k)$

475    $= \min\limits_{1 \le l \le k} \frac{1}{n_{ki;l}(n_{ki'l} - 1)} \sum\limits_{j \ne j' \in A_{ki;l}} I\{\text{validation samples } i'j \text{ and } i'j' \text{ are classified to the same group according to study } i\}$

476    Where $A_{kil}$ indicates the $l$-th cluster membership in study $i$, when cluster number is specified as

477    $k$; $n_{kil} = |A_{kil}|$. The average generalized prediction in study $i$ for cluster number $k$ is then defined

478    as the average of $gps_{ii'}(k)$ across all $i' \ne i$, i.e., all validation studies (instead of iterations of

479    randomized partitions). Similar to the single study prediction strength, it describes the

480    generalizability of clustering structure in study $i$ in external validation studies.

481    Unsupervised continuous structure discovery

482    We extended our previous work in cancer gene expression subtyping[35] to perform unsupervised

483    continuous structure discovery in microbial community profiles. Complementary to discrete cluster

484    discovery, the goal is to identify strong feature covariation signals (gradients) that are reproducible

485    across studies. This is carried out by performing principal component analysis individually in

486    microbiome studies and constructing a network of correlated PCA loading vectors, to identify

487    loadings that are consistently present across studies. In detail, given a collection of training

488    microbial abundance datasets, our method takes the following steps (visualized in **Supplemental**

489    **Fig. 4**):

490    1. For each dataset $i$, PCA is performed on normalized and arcsin square root-transformed

491    microbial abundance data. Given a user-specified threshold on variance explained, we

492    record its top PC loading vectors, $w_{i1}, w_{i2}, \ldots, w_{iJ_i}$, where $J_i$ is the smallest number of top

493    loading vectors that jointly explain percentage of variability in the dataset past a

494    customizable threshold $0 < threshold_v < 1$ (default to 80%).

495    2. For two PC loadings from different datasets $w_{ij}$ and $w_{i'j'}$, similarity is quantified with the

496    absolute value of cosine coefficient[56] $|cos < w_{ij}, w_{i'j'} >|$. This yields a network of PC

497    loading vectors associated by weighted edges $w_{ij}$ and $w_{i'j'}$, retaining edges only if their

498    weight surpasses a customizable similarity threshold ( $|cos < w_{ij}, w_{i'j'} >| >$

499    $threshold_s, 0 < threshold_s < 1$).

500    3. In the resulting network, we perform community detection[57] to identify densely connected

501    modules of PCs. Each module by definition consists of PCs from different datasets that

502    are similar to each other - whether or not they occur in the same order or with similar

503    percent variance explained - and which thu represent strong feature covariation signals

504    that are recurrent in studies.

505    4. For a module $k$ containing PC set $M_k$, its consensus vector $W_k$ is calculated as the

506    average of sign-corrected loading vectors in $M_k$, i.e., $W_k := \frac{\sum_{w_{ij} \in M_k} \widetilde{w_{ij}}}{|M_k|}$. Note that the

507    average is taken not over the original loading vectors $w_{ij}$, but rather their sign-corrected

508    versions $\widetilde{w_{ij}}$. Specifically, the signs of each $w_{ij}$ in $M_k$ are corrected so that all of the

509    loading vectors have positive cosine coefficients.

510    5. The module-wide consensus vectors $W_k$ represent strong, mutually independent, and

511    reproducible covariation signals across the microbial datasets; they are used to identify

512    continuously varying gradients in microbial abundance profiles that represent reproducible

513    population structures. Specifically, given a sample with normalized and transformed

514    microbial abundance measurements $x$, its continuous score for module $k$ is defined as

515    $x'W_k$, as in regular PCA.

516    6. If additional studies are available, the reproducibility of each $W_k$ can be further examined

517    by correlating $W_k$ with the top PC loadings in each such validation study. For each

518    additional study, $W_k$ is considered to be validated in that dataset if its absolute cosine

519    coefficient with at least one of the dataset's top PCs surpasses the coefficient similarity

520    cutoff $threshold_s$; the number of top PCs to consider in the validation dataset loadings is

521    determined with the same cutoff $threshold_v$.


522    **Simulation validation of MMUPHin**

523    We performed extensive simulation studies (**Fig. 2, Supplemental Fig. 5-8, Supplemental Table**

524    **2**) to validate the performance of each component of MMUPHin. In all cases these employed

525    realistic microbial abundance profiles generated using SparseDOSSA

526    (http://huttenhower.sph.harvard.edu/sparsedossa). This is a model of microbial community

527    structure using a set of zero-inflated log-normal distributions fit to selected training data, in this

528    case drawn from the IBD gut microbiome[6]. Controlled microbial associations with simulated

529    covariates can then (optionally) be spiked in. Note that although the assumed null distributions in

530    MMUPHin and SparseDOSSA are the same (zero-inflated log normal), the models of effects for

531    batch and biological variables are substantially different: MMUPHin assumes exponentiated

532    effects, while SparseDOSSA assumes re-standardized linear effects.

533    Specifically, SparseDOSSA models null microbial feature abundances using a zero-inflated log-

534    normal distribution:

535
$$log(Y_{ip}) \sim N(\mu_p, {\sigma^2}_p) \times Bernoulli(\pi_p)$$

536    This is the same initial distributional assumption as the MMUPHin batch correction model, when

537    there are no batch or covariates effects. However, for spiked-in associations with metadata (batch,

538    biological variables, etc.), SparseDOSSA uses a different model. Given a simulated, pre-spiking-

539    in feature count vector $Y_p$ with mean $\mu_p{}^Y$ and standard error $\sigma_p{}^Y$, as well as a metadata variable

540    vector $X$ with mean $\mu^X$ and standard error $\sigma^X$, the post-spiked-in feature count is set to:

541
$$\widetilde{Y_{ip}} = \frac{1}{1+\phi}\{Y_{ip} + \phi \times [\frac{(X_i - \mu^X)\sigma_p^Y}{\sigma^X} + \mu_p{}^Y]\}$$

542    where $\phi$ is a configurable spike-in strength parameter. By this definition, microbial features post-

543    spike-in have the same mean and approximately the same variance as before, the only difference

544    being the added association with the metadata variable(s) used. This is to ensure the counts of

545    the modified feature are not dominated by the values of the target covariate, but instead

546    distributed similarly to real data. The SparseDOSSA association model thus differs from

547    MMUPHin's model in two substantial ways: i) MMUPHin's associations are defined within the

548    exponentiated component and are thus better described as a multiplicative effect, whereas

549    SparseDOSSA's effects are directly applied on untransformed data, and ii) SparseDOSSA

550    additionally ensures realistic data generation with the re-standardization procedure.

551    Thus, the only component of the SparseDOSSA model that requires fitting to training data is the

552    aforementioned zero-inflated log-normal null distribution. In our analysis, this was always PRISM[6],

553    while other parameters were specified across a wide range of combinations to simulate different

554    application scenarios. These include the effect sizes of the associated batch and biological

555    variables (i.e. the $\phi$ parameter), number of batches, sample sizes, as well as dimensionality (both

556    the total number of features and the percentage of features randomized to be associated with

557    batch/biological variables). For each combination of simulation parameters, we performed 20

558    random replications (i.e. running simulation/evaluation with the same parameters but different

559    random seeds). **Supplemental Table 2** presents the full list of parameter combinations.

560    <u>Evaluating batch adjustment</u>

561    For evaluation of MMUPHin's batch effect adjustment component, we simulated metadata that

562    included batch (with varying total batch numbers 2, 4, 6, 8), a binary positive control (simulated

563    "biological" covariate), continuous positive control ("biological"), and negative control (binary, and

564    guaranteed to be non-associated with microbial features) variables. Microbial abundance data

565    was simulated to be associated with the batch and the two positive control variables at varying

566    effect sizes (1, 2, 5, 10 for batch variable and fixed at 10 for positive control variables), but not

567    with the negative control variable. We additionally varied the number of samples per batch (20 to

568    simulate multiple-batches in a single study scenario, 100 to simulate meta-analysis with moderate

569    sized studies and 500 to simulate large meta-analysis), total number of microbial features (n=200

570    and 1000), as well as the percentage of features associated with metadata (5%, 10%, and 20%)

571    (**Supplemental Table 2**).

572    Performance of batch correction methods was quantified by omnibus associations (PERMANOVA

573    $R^2$) between the simulated microbial abundance data with the batch and positive control variables,

574    before and after batch correction. For ComBat[15] and our method, batch correction was performed

575     with both positive control variables as well as the negative control variable as covariates.

576     MMUPHin successfully reduced the confounding batch effect, but retained the effect of positive

577     control variables, and did not inflate the effect of negative control variable (**Fig. 2a**, **Supplemental**

578     **Fig. 5**).

579     <u>Evaluating meta-analytic differential abundance testing</u>

580     We evaluated false positive rates (FPR) in particular for meta-analytic feature association testing,

581     specifically the null case in which there are no associations between microbial features and

582     covariates, but false associations can arise in the presence of batch effects with unbalanced

583     distribution of covariate values across studies (**Fig. 2b**). For simulation, we generated a binary

584     covariate unevenly distributed between two "studies" at varying levels of disparity (**Supplemental**

585     **Table 2**). Microbial abundance data was simulated to be associated only with the two studies and

586     not with the covariate (i.e. study confounded null data), with varying strengths of batch effect (from

587     0 to 10). The number of samples per batch varied between 100 and 500 to, again, simulate

588     moderate- and large-sized meta-analysis. Lastly, we varied total number of microbial features and

589     the percentage of features associated with metadata as above.

590     FPRs were calculated as the percentage of simulated microbial features with nominal p-values <

591     0.05 for associations with the exposure variable. Four data normalization and analysis regimes

592     were evaluated (**Fig. 2c**, **Supplemental Fig. 6**): a) naive MaAsLin2 model on the study effect

593     confounded null data (without explicitly modelling the batches), b) the quantile normalization

594     procedure, paired with two-tailed Wilcoxon tests, as proposed in [18], c) BDMMA as proposed in [19],

595     with the default 1,0000 total MCMC sampling and 5,000 burn-in, d) the complete MMUPHin meta-

596     analysis model for the batch corrected data as described above. Note that due to its computational

597     cost we were only able to evaluate the Dirichlet-multinomial regression model on a subset of

598     parameter combinations, namely number of samples per batch = 100, number of features = 200,

599    and percent of associated microbes = 5%. These parameters roughly agree with those used in

600    the simulation analysis in the method's original publication[19].

601    We also evaluated the computational costs of quantile normalization, BDMMA, and MMUPHin

602    (**Supplemental Fig. 6**). For this, the same subset of 20 replications (batch effect 0, exposure

603    imbalance 0, number of samples per batch 100, and number of features 200) were ran through

604    the three methods under the same computation environment (single core Intel(R) Xeon(R) CPU

605    E5-2680 v2 @ 2.80GHz).

606    <u>Evaluating unsupervised discrete structure discovery</u>

607    To simulate microbial abundance data with known discrete clustering structure, we again used

608    the simulation model above, with microbial feature associations added both with a discrete "batch"

609    variable and a discrete clustering variable, at varying number of batches (2, 4, 6, 8), number of

610    clusters (3, 4, 5, 6), as well as effect size of association (0 to 10 for batch, fixed at 10 for cluster).

611    For the evaluation of MMUPHin's unsupervised methods (both here and during continuous

612    population structure discovery below), we fixed the number of samples per batch at 500, the

613    number of total features at 1,000, and the percent of associated features at 20%. These were

614    guided by the fact that the underlying unsupervised methods (clustering, PCA) require larger

615    sample sizes for good performance even without batch confounding, and are generally only

616    practical with higher feature dimensions (**Supplemental Table 2**).

617    Performance of clustering was evaluated as the percentage of replicates in which the right number

618    of synthetically defined underlying clusters was identified using prediction strength, across

619    technical replicates for a fixed combination of simulation parameters. That is, the number of

620    clusters within a simulation was identified as that which maximized prediction strength. This was

621    compared to the "truth" (i.e. the known simulation parameter) and counted as a success only if

622    the two agreed. The percentage of success for a given parameter combination across the 20

623    random replications was used as the evaluation metric for model performance. We compared the

624    performance of clustering before and after MMUPHin batch correction (**Fig. 2e**, **Supplemental**

625    **Table 7**). Note that batch correction is modelled only using the batch variable and specifically not

626    including the cluster variable as a covariate in the batch correction model above, as the underlying

627    cluster structure is unknown in non-synthetic unsupervised analyses settings.

628    <u>Evaluating unsupervised continuous structure discovery</u>

629    To simulate microbial abundance data with known continuously variable population structure, we

630    spiked in feature associations with both a simulated batch covariate (4, 6, 8) and a continuously

631    varying gradient (uniformly distributed between -1 and 1), at varying number of batches and effect

632    size of both associations (as above). The number of samples per batch, total number of microbial

633    features, and the percentage of features associated were fixed at the same values as above

634    (**Supplemental Table 2**).

635    Performance of continuous structure discovery analysis was evaluated as the Spearman

636    correlation between the known simulated gradient score and the strongest continuously valued

637    population structure as identified by MMUPHin's continuous structure discovery method (above).

638    We again compared the performance of continuous score discovery on the batch confounded and

639    batch corrected data (**Fig. 2g**, **Supplemental Fig. 8**). Note that, as above, batch correction is

640    again modelled only using the batch variable and does not have any access to the synthetic

641    continuous gradient, as any underlying continuous population structure is unknown during

642    unsupervised analyses settings.

**643** **Collection and uniform processing of ten IBD microbiome studies employing 16S rRNA**

**644** **gene sequencing**

**645** <u>Study inclusion and raw sequence data</u>

**646** We curated 10 published 16S rRNA gene sequencing (abbreviated 16S) gut microbiome studies

**647** of IBD for meta-analysis (**Table 1, Supplemental Table 1**). Demultiplexed raw sequences were

**648** either downloaded from EBI (Jansson-Lamendella and Herfarth) or available locally as previously

**649** generated (other eight studies). Metadata were obtained either directly from the sequence

**650** repository/manuscript (Herfarth, Jasson-Lamendella, HMP2, MucosalIBD, PROTECT, RISK), or

**651** from collaborators (BIDMC-FMT, CS-PRISM, LSS-PRISM, Pouchitis). This resulted in a total of

**652** 5,151 samples and 2,179 subjects available prior to processing and quality control.

**653** <u>Metadata curation</u>

**654** We manually curated subject- and sample-specific metadata across studies to ensure

**655** consistency. Variables collected and curated include:

**656** ● Disease (CD, UC, control), universally available.

**657** ● Type of controls (non-IBD, healthy). Control information was available directly for CS-

**658** PRISM, Jansson-Lamendella, and Pouchitis, inferred from study design described in

**659** manuscript for Herfarth, HMP2, MucosalIBD, and RISK (all non-IBD controls), and not

**660** applicable for BIDMC-FMT, LSS-PRISM, and PROTECT (only has IBD subjects).

**661** ● Sample type (biopsy, stool), universally available.

**662** ● Body site of biopsy sample collection (ileum, colon, rectum), with more detailed

**663** classifications recorded separately in case of need. Mappings for the relevant datasets

**664** are:

- CS-PRISM: terminal ileum, neo-ileum, pouch are aggregated as ileum; cecum, ascending/left-sided colon, transverse colon, descending/right-sided colon, sigmoid colon were aggregated as colon; rectum classification was kept unchanged.

- HMP2: ileum classification kept unchanged; cecum, ascending/right-sided colon, transverse colon, descending/left-sided colon, and sigmoid colon were aggregated as colon.

- MucosalIBD: all terminal ileum samples, aggregated to ileum.

- Pouchitis: terminal ileum, pouch, pre-pouch ileum aggregated as ileum; sigmoid colon aggregated to colon.

- PROTECT: all rectum samples, classification kept unchanged.

- RISK: terminal ileum was aggregated to ileum; rectum kept unchanged.

- Montreal classifications:

  - Location for CDs (L1, L2, L3, and possible combinations), available for BIDMC-FMT, CS-PRISM, Herfarth, Jansson-Lamendella, LSS-PRISM, and Pouchitis.

  - Behavior for CDs (B1, B2, and B3), available for CS-PRISM, Herfath, Jansson-Lamendella, LSS-PRISM, Pouchitis, and RISK.

  - Extent for UCs (E1, E2, and E3), available for CS-PRISM, Jansson-Lamendella, LSS-PRISM, Pouchitis, and PROTECT.

- Age at sample collection (in years), available for BIDMC-FMT, CS-PRISM, Herfarth, HMP2, LSS-PRISM, MucosalIBD, Pouchitis, PROTECT, RISK.

- Age at diagnosis (in years). Directly available for CS-PRISM, HMP2, LSS-PRISM, and Pouchitis, inferred as baseline age for PROTECT and RISK as these were new-onset cohorts.

- Race (White, African American, Asian / Pacific Islander, Native American, more than one race, others). Directly available for CS-PRISM, Herfarth, HMP2, PROTECT, and RISK,

691        inferred from manuscript cohort description for Jansson-Lamendella (all Caucasian

692        cohort).

693        ● Gender (male/female). Available for BIDMC-FMT, CS-PRISM, Herfarth, HMP2, Jansson-

694        Lamendella, LSS-PRISM, MucosalIBD, Pouchitis, PROTECT,

695        ● Treatment variables, including antibiotics, immunosuppressants, steroids, and 5-ASA.

696        These variables were encoded as yes/no to indicate, approximately, currently receiving

697        them at the time of sampling. Additional information such as specific medication or delivery

698        method was recorded separately if available in case of need. We note the potentially

699        confounding difference in studies' definitions of treatment: for Pouchitis and PROTECT

700        authors defined antibiotics as receiving the treatment within the past month (30 days for

701        Pouchitis, 27 days for PROTECT), whereas for CS-PRISM, HMP2, LSS-PRISM, and RISK

702        such determination was not possible (antibiotics "yes" was defined as "currently taking").

703        Likewise, we had no additional information to determine the time extent for the other three

704        treatments, beyond that according to metadata/publication, patients were "currently taking"

705        the treatment at sample collection.

706    For a comprehensive list of curation mapping schema, please refer to our metadata curation

707    repository: https://github.com/biobakery/ibd_meta_analysis.

708    <u>16S amplicon sequence bioinformatics and taxonomic profiling</u>

709    Sequences were processed, per-cohort, with the published, standardized bioBakery workflow[58]

710    using the UPARSE protocol[59] (version v9.0.2132-64bit). For all studies, demultiplexed sequences

711    were truncated at 200bp max length and filtered by maximum expected error of one[59]. Operational

712    taxonomic units (OTUs) were clustered at 97% identity and aligned using USEARCH with 97%

713    identity to the Greengenes database 97% reference OTUs (version 13.8)[60] for taxonomy

714    assignment. The resulting Greengenes identifiers for OTUs were used as basis for matching

715    features (taxa) among cohorts.

716    Quality control

717    Across samples, a median of 81.51% reads / sample passed quality control filtering and were

718    successfully assigned to OTUs with Greengenes identifiers (**Supplemental Fig. 1**). These 8,921

719    raw OTUs aggregated to a total of 1,122 genera prior to quality control. We retained taxa that

720    exceeded 5e-5 relative abundance with at least 10% prevalent in at least one study; this criterion

721    generally removes spurious OTU assignments while retaining rare organisms if confidently

722    present in at least one study. Lastly, we also removed low read depth samples with less than

723    3,000 total sequences, which retained 78.34%-100% samples per cohorts (**Supplemental Table**

724    **1**). The final resulting taxonomic profile, used for all further analysis, aggregated into 249 total

725    genera spanning 4,789 samples (OTUs unclassified under a particular taxonomy level were

726    aggregated as "unclassified" feature under that taxon, e.g. "Enterbacteriaceae unclassified"

727    accumulates all OTUs' abundances under the family that could not be classified at the genus level.

728    Data availability

729    Quality controlled (truncated and filtered) sequences, Greengenes mapped OTU count profiles,

730    and curated sample metadata are available at the Human Microbial Bioactives Resource Portal

731    (http://portal.microbiome-bioactives.org).

732    **Applying MMUPHin to IBD gut microbiome meta-analysis**

733    For the resulting collection of microbiome studies, batch and study effects was performed using

734    MMUPHin on both the genus level feature abundance profiles. For either taxonomic rank, batch

735    (i.e., sequencing run) effect correction was first performed within individual studies (when

736    batch/plate information was available, applicable to BIDMC-FMT, CS-PRISM, LSS-PRISM,

737    MucosalIBD, and RISK). Microbial abundance profiles across all studies were then jointly

738    corrected for study effects, while modelling disease status (IBD or control), disease (CD or UC),

739    and sample type (biopsy or stool) as covariates. Reduction of batch and study effects was

740    evaluated by PERMANOVA R2 (**Fig. 3a**).

741    **Association analyses**

742    <u>Omnibus testing of microbial composition associations</u>

743    We used PERMANOVA tests (2,000 permutations) as implemented in the R package vegan[37]

744    using Bray-Curtis dissimilarities for all omnibus association tests of overall microbial community

745    structure with covariates (**Fig. 3a**). Where appropriate, R2s were calculated conditioning on the

746    necessary covariates; specifically, CD/UC Montreal classifications were conditional on CD/UC

747    samples respectively, treatment was conditional on IBD status, biopsy location was conditional

748    on a sample being a biopsy, and all covariates were conditional on being non-missing. Otherwise,

749    variables were tested marginally (that is, each as the sole variable in the model). Importantly, to

750    account for repeated measures within subjects for longitudinal studies, we adopted the blocked

751    permutation strategy as in [9], where per-sample measurements (sample type, biopsy location,

752    treatment) were permuted within subjects, and per-subject measurements (disease,

753    demographics) were permuted along with subjects (but within cohorts, relevant for the all-cohorts

754    evaluation). For a full list of the model and permutation strategies that this resulted in for our

755    analysis, please refer to **Supplemental Table 3**. Finally, per-variable p-values were adjusted with

756    Benjamini-Hochberg false discovery rate control on a per-study basis.

757    <u>Per-feature meta-analysis differential abundance testing</u>

758    To identify microbial features individually significantly associated with one or more covariates, we

759    applied MMUPHin's differential abundance testing model as described above. Cohorts were first

760    stratified by sample type (biopsy or stool) and, where appropriate, diseases (CD or UC) prior to

761    model fitting.  Arcsin square root-transformed genus level taxon abundances were tested for

762     covariate associations in individual cohort strata with multivariate linear modelling (linear random

763     intercept model adopted for longitudinal studies). Covariates used for adjustment include age,

764     gender, and race for disease variables, and additionally disease status for treatment variables.

765     Effect sizes across cohort strata were aggregated with a random effects model with restricted

766     maximum likelihood estimation[34]. P-values were FDR adjusted across features for each variable.

767     For the full list of models adopted as well as cohort stratification strategy, please refer to

768     **Supplemental Table 3**. **Fig. 3b** visualizes the aggregated meta-analysis effects; for individual

769     study results refer to **Supplemental Table 4**.

770     <u>Testing for phenotypic severity within CD and UC patients</u>

771     Meta-analytical testing of features associated with CD behavior and UC extent classifications

772     were performed with similar models (**Supplemental Table 3**). Specifically, within each study's

773     CD patients, the tests for contrasts B2 versus B1 and B3 versus B1 are performed by

774     $\text{Relative abundance} \sim \beta_0 + \beta_1 I\{\text{subject is B2}\} + \text{additional covariates (subsetted to B1, B2 CDs)}$

775     $\text{Relative abundance} \sim \beta_0 + \beta_1 I\{\text{subject is B3}\} + \text{additional covariates (subsetted to B1, B3 CDs)}$

776     The two $\beta_1$ coefficients, once aggregated with meta-analysis, were reported as the effect sizes

777     shown in **Fig. 4a**, along with their FDR corrected q-values (adjusted across features for each

778     test).

779     $\text{Relative abundance} \sim \beta_0 + \beta_1 I\{\text{subject is } B2 \text{ or } B3\} + \beta_2 I\{\text{subject is } B3\} + \text{additional covariates}$

780     $\beta_2$ in this model corresponds to the effect of B3 in addition to the overall contrasts between B23

781     versus B1. The meta-analysis aggregated p-values of these effects were reported as the

782     differentiation between the most severe and "medium" severity phenotypes (vertical bars

783     indicating significance in **Fig. 4a**). Note that FDR adjustment of this effect was performed across

784     the subset of features with at least either B2 versus B1 or B3 versus B1 effect significant (i.e., the

785    subset of features visualized in **Fig. 4a**). Equivalent models were adopted for contrasts between

786    extent categories of UC patients. Individual study results for the aggregated effects in **Fig. 4a** are

787    in **Supplemental Table 5**.

788    Interaction effects testing

789    To test for interaction effects with sample type and diseases, we fit meta-analysis moderator

790    models[34] on the per cohort strata effects:

791    $$\widehat{\beta_{ip}} = \beta_{0p} + \beta_{1p}I\{\text{cohort strata } i \text{ is biopsy}\} + \epsilon_{ip} + e_{ip}$$

792    $$\widehat{\beta_{ip}} = \beta_{0p} + \beta_{1p}I\{\text{cohort strata } i \text{ is CD}\} + \epsilon_{ip} + e_{ip}$$

793    The moderator effects $\beta_{1p}$ correspond to the interaction effect between the exposure under

794    evaluation (disease, treatment, etc.) with the moderator variable. **Fig. 4b** visualizes the two

795    example features, *Dehalobacterium* and Enterobacteriacea; al significant interactions as well as

796    individual study effects are in **Supplemental Table 6**.

797    **Population structure analyses**

798    Discrete structure discovery

799    We performed discrete subtype discovery (i.e. "enterotyping"[61]) in IBD, CD, and UC populations

800    across studies (longitudinal studies subsetted to baseline samples), using MMUPHin's discrete

801    structure discovery component. Only studies with at least 33 samples were considered for

802    clustering analysis, as this was the sample size in the original enterotype paper[26]. Specifically,

803    clustering was performed on Bray-Curtis dissimilarity by the partition-around-medoid method as

804    implemented in R package cluster; the same method was adopted in previous enterotyping efforts

805    including the original enterotype paper[26,40]. Clustering was evaluated with prediction strength and

806    validated externally with MMUPHin's generalized prediction strength as described above. Across

807  studies, we found no evidence to support a particular number of clusters within IBD, CD, or UC

808  populations (**Fig. 5a**, **Supplemental Fig. 9**), suggesting that the IBD microbiome does not have

809  discrete clusters.

810  We additionally extended our clustering evaluation analysis to other dissimilarity metrics (Jaccard,

811  root Jensen-Shannon divergence) and clustering strength measurements (Calinski-Harabasz

812  index, average silhouette width), which were also explored in previous efforts[40], Importantly, the

813  original enterotype paper adopted root Jensen-Shannon divergence and Calinski-Harabasz index

814  for cluster discovery. Across combinations of these additional dissimilarities and clustering

815  strength metrics, we also found no evidence to support discrete clusters (**Supplemental Fig. 9**).

816  <u>Continuous structure discovery</u>

817  Continuous structure discovery was performed with MMUPHin's corresponding component. The

818  four largest studies (CS-PRISM, Pouchitis, PROTECT, RISK) were subsetted to baseline samples

819  (only relevant for PROTECT), stratified by CD/UC and biopsy/stool sample type, and used as the

820  training sets for MMUPHin. The minimum variance explained threshold ($threshold_v$) was set to

821  default (80%), but we varied the PC similarity (evaluated by absolute cosine coefficient)

822  cutoff$threshold_s$ between 0.5 and 0.8 to assess the sensitivity of the two identified PC clusters in

823  **Fig. 5b** (corresponding to $threshold_s$ = 0.65). As we show in **Supplemental Fig. 11**, with a small

824  $threshold_s$(0.5) PC networks become denser, with the two PC clusters in **Fig. 4b** forming key

825  components of two larger clusters; when $threshold_s$ is large (0.8) the network is sparser, with

826  only the most highly similar nodes of the two clusters forming smaller communities. We thus

827  concluded that the two identified clusters in **Fig. 5b** were not sensitive to the cosine coefficient

828  threshold, as they were recurrently identified in both smaller and larger cutoff scenarios.

829 <u>Continuous structure validation</u>

830 We validated the consistency of the two clusters' corresponding continuous scores in all IBD

831 cohorts, non-IBD and healthy control samples, as well as a randomly permuted mock study (as a

832 negative control). The reproducibility of each continuous score within a study was defined as the

833 maximum absolute cosine coefficient between the score's consensus loading (as provided by

834 MMUPHin) and the top three principal component loadings discovered independently within that

835 study. Note that the number of top principal components considered here was set to a fixed value

836 (three) instead of based on a percent variance cutoff as in the MMUPHin continuous structure

837 discovery stage. This is because in the two identified clusters in **Fig. 5c**, the latest included node

838 was PC3. The randomly permuted study consisted of 473 samples (median validation data sets

839 sample size) randomly selected from the entire meta-analysis collection, but each sample's

840 microbial abundance was independently permuted across features. This was to simulate a

841 "negative control" dataset where there should be no continuous population structures.

842 As we show in **Supplemental Fig. 12**, the dysbiosis score was well validated across studies,

843 except for healthy control samples and the negative control dataset. The Firmicutes-versus-

844 Bacteroidetes trade-off score, on the other hand, was reasonably well reproduced in all studies

845 and particularly well-established in healthy samples, but, again, was not significantly detected in

846 the negative control dataset.

847 <u>Continuous score assignment</u>

848 Assignment of continuous scores was straightforward given the two consensus loading vectors

849 provided by MMUPHin. Within each study, arcsin square root-transformed relative abundances

850 were centered per-feature, the transformed abundance matrix was then multiplied by each

851 consensus loading via dot product to generate per-sample continuous scores. These scores were

852 used for visualization as in **Fig. 4d** and **Supplemental Fig. 10**, as well as for testing the difference

853    between CD, UC, non-IBD, and healthy control populations as in **Supplemental Fig. 13** We

854    provide the two consensus loadings in **Supplemental Table 7;** interested researchers can follow

855    these steps to assign the two continuous scores in other datasets.

856

# References

1. Hirschhorn, J. N. & Daly, M. J. Genome-wide association studies for common diseases and complex traits. *Nat Rev Genet* **6**, 95-108, doi:10.1038/nrg1521 (2005).

2. Lambert, J. C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. *Nat Genet* **45**, 1452-1458, doi:10.1038/ng.2802 (2013).

3. Franke, A. *et al.* Genome-wide meta-analysis increases to 71 the number of confirmed Crohn's disease susceptibility loci. *Nat Genet* **42**, 1118-1125, doi:10.1038/ng.717 (2010).

4. Anderson, C. A. *et al.* Meta-analysis identifies 29 additional ulcerative colitis risk loci, increasing the number of confirmed associations to 47. *Nat Genet* **43**, 246-252, doi:10.1038/ng.764 (2011).

5. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet* **38**, 904-909, doi:10.1038/ng1847 (2006).

6. Morgan, X. C. *et al.* Dysfunction of the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biol* **13**, R79, doi:10.1186/gb-2012-13-9-r79 (2012).

7. Gevers, D. *et al.* The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe* **15**, 382-392, doi:10.1016/j.chom.2014.02.005 (2014).

8. Ananthakrishnan, A. N. Environmental risk factors for inflammatory bowel diseases: a review. *Dig Dis Sci* **60**, 290-298, doi:10.1007/s10620-014-3350-9 (2015).

9. Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature* **569**, 655-662, doi:10.1038/s41586-019-1237-9 (2019).

10. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190-2191, doi:10.1093/bioinformatics/btq340 (2010).

11. Rhodes, D. R. *et al.* Large-scale meta-analysis of cancer microarray data identifies common transcriptional profiles of neoplastic transformation and progression. *Proc Natl Acad Sci U S A* **101**, 9309-9314, doi:10.1073/pnas.0401994101 (2004).

12. Li, H. Microbiome, Metagenomics, and High-Dimensional Compositional Data Analysis. *Annual Review of Statistics and Its Application* **2**, 73-94, doi:10.1146/annurev-statistics-010814-020351 (2015).

13. Mallick, H. *et al.* Experimental design and quantitative analysis of microbial community multiomics. *Genome Biol* **18**, 228, doi:10.1186/s13059-017-1359-z (2017).

14. Li, Y., Willer, C., Sanna, S. & Abecasis, G. Genotype imputation. *Annu Rev Genomics Hum Genet* **10**, 387-406, doi:10.1146/annurev.genom.9.081307.164242 (2009).

15. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* **8**, 118-127, doi:10.1093/biostatistics/kxj037 (2007).

893   16   Leek, J. T. & Storey, J. D. Capturing heterogeneity in gene expression studies by
894        surrogate variable analysis. *PLoS Genet* **3**, 1724-1735,
895        doi:10.1371/journal.pgen.0030161 (2007).

896   17   Leek, J. T. *et al.* Tackling the widespread and critical impact of batch effects in high-
897        throughput data. *Nat Rev Genet* **11**, 733-739, doi:10.1038/nrg2825 (2010).

898   18   Gibbons, S. M., Duvallet, C. & Alm, E. J. Correcting for batch effects in case-control
899        microbiome studies. *PLoS Comput Biol* **14**, e1006102, doi:10.1371/journal.pcbi.1006102
900        (2018).

901   19   Dai, Z., Wong, S. H., Yu, J. & Wei, Y. Batch effects correction for microbiome data with
902        Dirichlet-multinomial regression. *Bioinformatics* **35**, 807-814,
903        doi:10.1093/bioinformatics/bty729 (2019).

904   20   Manichanh, C., Borruel, N., Casellas, F. & Guarner, F. The gut microbiota in IBD. *Nat Rev*
905        *Gastroenterol Hepatol* **9**, 599-608, doi:10.1038/nrgastro.2012.152 (2012).

906   21   Kostic, A. D., Xavier, R. J. & Gevers, D. The microbiome in inflammatory bowel disease:
907        current status and the future ahead. *Gastroenterology* **146**, 1489-1499,
908        doi:10.1053/j.gastro.2014.02.009 (2014).

909   22   Halfvarson, J. *et al.* Dynamics of the human gut microbiome in inflammatory bowel disease.
910        *Nat Microbiol* **2**, 17004, doi:10.1038/nmicrobiol.2017.4 (2017).

911   23   Schirmer, M. *et al.* Compositional and Temporal Changes in the Gut Microbiome of
912        Pediatric Ulcerative Colitis Patients Are Linked to Disease Course. *Cell Host Microbe* **24**,
913        600-610 e604, doi:10.1016/j.chom.2018.09.009 (2018).

914   24   Pasolli, E., Truong, D. T., Malik, F., Waldron, L. & Segata, N. Machine Learning Meta-
915        analysis of Large Metagenomic Datasets: Tools and Biological Insights. *PLoS Comput*
916        *Biol* **12**, e1004977, doi:10.1371/journal.pcbi.1004977 (2016).

917   25   Duvallet, C., Gibbons, S. M., Gurry, T., Irizarry, R. A. & Alm, E. J. Meta-analysis of gut
918        microbiome studies identifies disease-specific and shared responses. *Nat Commun* **8**,
919        1784, doi:10.1038/s41467-017-01973-8 (2017).

920   26   Arumugam, M. *et al.* Enterotypes of the human gut microbiome. *Nature* **473**, 174-180,
921        doi:10.1038/nature09944 (2011).

922   27   Vazquez-Baeza, Y. *et al.* Guiding longitudinal sampling in IBD cohorts. *Gut* **67**, 1743-1745,
923        doi:10.1136/gutjnl-2017-315352 (2018).

924   28   Morgan, X. C. *et al.* Associations between host gene expression, the mucosal microbiome,
925        and clinical outcome in the pelvic pouch of patients with inflammatory bowel disease.
926        *Genome Biol* **16**, 67, doi:10.1186/s13059-015-0637-x (2015).

927   29   Franzosa, E. A. *et al.* Gut microbiome structure and metabolic activity in inflammatory
928        bowel disease. *Nat Microbiol* **4**, 293-305, doi:10.1038/s41564-018-0306-4 (2019).

929   30   Liu, T. C. *et al.* Paneth cell defects in Crohn's disease patients promote dysbiosis. *JCI*
930        *Insight* **1**, e86907, doi:10.1172/jci.insight.86907 (2016).

31    Hall, A. B. *et al.* A novel Ruminococcus gnavus clade enriched in inflammatory bowel disease patients. *Genome Med* **9**, 103, doi:10.1186/s13073-017-0490-5 (2017).

32    Vaughn, B. P. *et al.* Increased Intestinal Microbial Diversity Following Fecal Microbiota Transplant for Active Crohn's Disease. *Inflamm Bowel Dis* **22**, 2182-2190, doi:10.1097/MIB.0000000000000893 (2016).

33    Mallick, H., Rahnavard, A. & McIver, L. Maaslin2: Maaslin2. R package version 1.2.0, http://huttenhower.sph.harvard.edu/maaslin2. *Bioconductor*, doi:10.18129/B9.bioc.Maaslin2 (2019).

34    Viechtbauer, W. Conducting Meta-Analyses in R with the metafor Package. *2010* **36**, 48, doi:10.18637/jss.v036.i03 (2010).

35    Ma, S. *et al.* Continuity of transcriptomes among colorectal cancer subtypes based on meta-analysis. *Genome Biol* **19**, 142, doi:10.1186/s13059-018-1511-4 (2018).

36    Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Methods* **12**, 115-121, doi:10.1038/nmeth.3252 (2015).

37    Oksanen, J. *et al.*    (2019).

38    Vilhjalmsson, B. J. & Nordborg, M. The nature of confounding in genome-wide association studies. *Nat Rev Genet* **14**, 1-2, doi:10.1038/nrg3382 (2013).

39    Ravel, J. *et al.* Vaginal microbiome of reproductive-age women. *Proc Natl Acad Sci U S A* **108 Suppl 1**, 4680-4687, doi:10.1073/pnas.1002611107 (2011).

40    Koren, O. *et al.* A guide to enterotypes across the human body: meta-analysis of microbial community structures in human microbiome datasets. *PLoS Comput Biol* **9**, e1002863, doi:10.1371/journal.pcbi.1002863 (2013).

41    Tibshirani, R. & Walther, G. Cluster Validation by Prediction Strength. *Journal of Computational and Graphical Statistics* **14**, 511-528, doi:10.1198/106186005X59243 (2005).

42    Wold, S., Esbensen, K. & Geladi, P. Principal component analysis. *Chemometrics and intelligent laboratory systems* **2**, 37-52, doi:10.1016/0169-7439(87)80084-9 (1987).

43    Satsangi, J., Silverberg, M. S., Vermeire, S. & Colombel, J. F. The Montreal classification of inflammatory bowel disease: controversies, consensus, and implications. *Gut* **55**, 749-753, doi:10.1136/gut.2005.082909 (2006).

44    Lin, M. F. & Lan, C. Y. Antimicrobial resistance in Acinetobacter baumannii: From bench to bedside. *World J Clin Cases* **2**, 787-814, doi:10.12998/wjcc.v2.i12.787 (2014).

45    Kevans, D. *et al.* Characterization of Intestinal Microbiota in Ulcerative Colitis Patients with and without Primary Sclerosing Cholangitis. *J Crohns Colitis* **10**, 330-337, doi:10.1093/ecco-jcc/jjv204 (2016).

46    Presley, L. L., Wei, B., Braun, J. & Borneman, J. Bacteria associated with immunoregulatory cells in mice. *Appl Environ Microbiol* **76**, 936-941, doi:10.1128/AEM.01561-09 (2010).

969  47  Jones-Hall, Y. L., Kozik, A. & Nakatsu, C. Ablation of tumor necrosis factor is associated
970      with decreased inflammation and alterations of the microbiota in a mouse model of
971      inflammatory bowel disease. *PLoS One* **10**, e0119441, doi:10.1371/journal.pone.0119441
972      (2015).

973  48  Imhann, F. *et al.* Interplay of host genetics and gut microbiota underlying the onset and
974      clinical presentation of inflammatory bowel disease. *Gut* **67**, 108-119, doi:10.1136/gutjnl-
975      2016-312135 (2018).

976  49  Sinha, R. *et al.* Assessment of variation in microbial community amplicon sequencing by
977      the Microbiome Quality Control (MBQC) project consortium. *Nat Biotechnol* **35**, 1077-1086,
978      doi:10.1038/nbt.3981 (2017).

979  50  Brooks, J. P. *et al.* The truth about metagenomics: quantifying and counteracting bias in
980      16S rRNA studies. *BMC Microbiol* **15**, 66, doi:10.1186/s12866-015-0351-6 (2015).

981  51  Schloss, P. D. Identifying and Overcoming Threats to Reproducibility, Replicability,
982      Robustness, and Generalizability in Microbiome Research. *MBio* **9**,
983      doi:10.1128/mBio.00525-18 (2018).

984  52  Tsilimigras, M. C. & Fodor, A. A. Compositional data analysis of the microbiome:
985      fundamentals, tools, and challenges. *Ann Epidemiol* **26**, 330-335,
986      doi:10.1016/j.annepidem.2016.03.002 (2016).

987  53  Thomas, A. M. *et al.* Metagenomic analysis of colorectal cancer datasets identifies cross-
988      cohort microbial diagnostic signatures and a link with choline degradation. *Nat Med* **25**,
989      667-678, doi:10.1038/s41591-019-0405-7 (2019).

990  54  Wirbel, J. *et al.* Meta-analysis of fecal metagenomes reveals global microbial signatures
991      that are specific for colorectal cancer. *Nat Med* **25**, 679-689, doi:10.1038/s41591-019-
992      0406-6 (2019).

993  55  Lloyd-Price, J. *et al.* Strains, functions and dynamics in the expanded Human Microbiome
994      Project. *Nature* **550**, 61-66, doi:10.1038/nature23889 (2017).

995  56  Jaskowiak, P. A., Campello, R. J. & Costa, I. G. On the selection of appropriate distances
996      for gene expression data clustering. *BMC Bioinformatics* **15 Suppl 2**, S2,
997      doi:10.1186/1471-2105-15-S2-S2 (2014).

998  57  Csardi, G. & Nepusz, T. The igraph software package for complex network research.
999      *InterJournal* **Complex Systems**, 1695-1695 (2006).

1000 58  McIver, L. J. *et al.* bioBakery: a meta'omic analysis environment. *Bioinformatics* **34**, 1235-
1001     1237, doi:10.1093/bioinformatics/btx754 (2018).

1002 59  Edgar, R. C. UPARSE: highly accurate OTU sequences from microbial amplicon reads.
1003     *Nat Methods* **10**, 996-998, doi:10.1038/nmeth.2604 (2013).

1004 60  McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological
1005     and evolutionary analyses of bacteria and archaea. *ISME J* **6**, 610-618,
1006     doi:10.1038/ismej.2011.139 (2012).

1007 61  Costea, P. I. *et al.* Enterotypes in the landscape of gut microbial community composition.
1008     *Nat Microbiol* **3**, 8-16, doi:10.1038/s41564-017-0072-8 (2018).

1009