# Multivariable Association Discovery in Population-scale Meta-omics Studies

3    Himel Mallick[1,2], Ali Rahnavard[3], Lauren J. McIver[1,2], Siyuan Ma[1,2], Yancong Zhang[1,2], Long H.

4    Nguyen[1,4,5], Timothy L. Tickle[2], George Weingart[1,2], Boyu Ren[1,2], Emma H. Schwager[1,2], Suvo

5    Chatterjee[6], Kelsey N. Thompson[1], Jeremy E. Wilkinson[1], Ayshwarya Subramanian[1,2], Yiren Lu[1], Levi

6    Waldron[7], Joseph N. Paulson[8], Eric A. Franzosa[1,2], Hector Corrada Bravo[9], Curtis Huttenhower[1,2]*

7    [1]Biostatistics Department, Harvard T. H. Chan School of Public Health, Boston, MA 02115, USA

8    [2]The Broad Institute, 415 Main Street, Cambridge, MA 02142, USA

9    [3]Computational Biology Institute, Department of Biostatistics and Bioinformatics, Milken Institute School of Public

10    Health, George Washington University, Washington, DC 20052, USA

11    [4]Clinical and Translational Epidemiology Unit, Massachusetts General Hospital and Harvard Medical School, Boston,

12    MA 02144, USA

13    [5]Division of Gastroenterology, Massachusetts General Hospital and Harvard Medical School, Boston, MA 02144,

14    USA

15    [6]Epidemiology Branch, Division of Intramural Population Health Research, Eunice Kennedy Shriver National Institute

16    of Child Health and Human Development, National Institutes of Health, Bethesda, MD 20892, USA

17    [7]Department of Epidemiology and Biostatistics, CUNY School of Public Health, NY 10027, USA

18    [8]Department of Biostatistics, Product Development, Genentech, Inc., South San Francisco, CA 94080, USA

19    [9]Center for Bioinformatics and Computational Biology, University of Maryland, College Park, MD 20742, USA

20    *Correspondence to chuttenh@hsph.harvard.edu

## Abstract

22    It is challenging to associate features such as human health outcomes, diet, environmental

23    conditions, or other metadata to microbial community measurements, due in part to their

24    quantitative properties. Microbiome multi-omics are typically noisy, sparse (zero-inflated), high-

25    dimensional, extremely non-normal, and often in the form of count or compositional

26    measurements. Here we introduce an optimized combination of novel and established

27    methodology to assess multivariable association of microbial community features with complex

28    metadata in population-scale observational studies. Our approach, MaAsLin 2 (Microbiome

29    Multivariable Associations with Linear Models), uses general linear models to accommodate a

30    wide variety of modern epidemiological studies, including cross-sectional and longitudinal

31    designs, as well as a variety of data types (e.g. counts and relative abundances) with or without

32    covariates and repeated measurements. To construct this method, we conducted a large-scale

33    evaluation of a broad range of scenarios under which straightforward identification of meta-omics

34    associations can be challenging. These simulation studies reveal that MaAsLin 2's linear model

35    preserves statistical power in the presence of repeated measures and multiple covariates, while

36    accounting for the nuances of meta-omics features and controlling false discovery. We also

37    applied MaAsLin 2 to a microbial multi-omics dataset from the Integrative Human Microbiome

38    (HMP2) project which, in addition to reproducing established results, revealed a unique,

39    integrated landscape of inflammatory bowel disease (IBD) across multiple time points and omics

40    profiles.

41    **Keywords**: Human Microbiome, Metagenomics, Differential Abundance Analysis, Multivariable

42    Association, Microbiome Epidemiology, Longitudinal Analysis


## Introduction

43

44    Human-associated microbiota has been convincingly linked to the development and progression

45    of a wide range of complex, chronic conditions including inflammatory bowel diseases (IBD),

46    obesity, diabetes, cancer, and cardiovascular disorders[1,2]. Due to recent advances in multiple

47    high-throughput functional profiling technologies, research has expanded well beyond bacteria-

48    specific 16S rRNA gene amplicon profiles to multi-omics surveys, i.e. non-bacterial,

49    metagenomic, metatranscriptomic, metabolomic, and metaproteomic measurements assessed

50    simultaneously in the same biological specimens[3,4]. Additionally, due to diminishing sequencing

51    costs, longitudinal, within-subject study designs are becoming increasingly common, especially

52    when assessing the microbiome in population health[5,6]. These large, complex data contain

53    abundant information to enable microbe-, gene-, and compound-specific hypothesis generation

54    at scale. However, robust quantitative methods to do so at scale can still be challenging to

55    implement without excessive false positives - one of the main hurdles in accurate translational

56    applications of the microbiome to human health.

57    One of the primary limitations of leveraging such population-wide multi-omics surveys is thus

58    computational, in part due to the complexity and heterogeneity of microbial community data that

59    have made reliable application of statistical methods difficult. In particular, best practices to guard

60    against spurious discoveries in meta-omics datasets remain scarce[7-14]. High-throughput meta-

61    omics datasets have specific characteristics that complicate their analyses: high-dimensionality,

62    count and compositional data structure, sparsity (zero-inflation), over-dispersion, and hierarchical,

63    spatial, and temporal dependence, among others. To combat these challenges, specialized

64    methods implemented in usable, reproducible software are needed to accurately characterize

65    microbial communities within large human population studies, while maintaining both sensitivity

66    and specificity.

67    Early advances in microbiome epidemiology focused on omnibus testing for identifying overall

68    associations between aggregate microbiome structure and host or environmental phenotypes and

69    covariates (e.g. disease status, diet, antibiotics or medication usage, age, etc.). Many of these

70    rely on permutation-based procedures for moderated significance testing[11]. These methods

71    assess whether overall community patterns of variation are associated with the covariates of

72    interest, but fail to provide feature-level inference to enable follow-up characterization (where a

73    feature can be any profiled omics abundance, e.g. taxa, genes, pathways, chemicals, etc.) To

74    facilitate actionable outcomes, it is important to discern feature-specific associations at the highest

75    possible resolution. This has led to the development of a variety of per-feature (or feature-wise)

76      association testing methods, most of which are based on similar statistical frameworks, differing

77      primarily in (i) the choice of normalization or transformation, (ii) observation model or likelihood,

78      and (iii) the associated statistical inference[11]. As compared to omnibus testing approaches, per-

79      feature methods (i) identify associations for each individual feature-metadata pair, (ii) facilitate

80      feature-wise covariate adjustment, and (iii) call out specific features (as opposed to complex

81      combinations of features implicated in associations in omnibus testing), leading to increased

82      interpretability for translational and basic biological applications.

83      Despite a rich literature on feature-wise association testing for microbial communities, methods

84      that can accommodate a wide variety of modern epidemiological study designs remain scarce.

85      For instance, many early methods do not explicitly account for the sparsity observed in microbial

86      meta-omics observations, and only a few scale beyond routine univariate (differential abundance)

87      analyses without becoming overly susceptible to false positive or false negative results[7,11].

88      Furthermore, most methods for microbiome data do not explicitly adjust for repeated measures

89      and multiple covariates in a unified statistical framework, a lack of which can have a profound

90      (and typically anti-conservative) impact on subsequent epidemiological inference.

91      Here, we address these issues by providing a flexible approach to identify multivariable

92      associations in large, heterogeneous meta-omics datasets. We have implemented this method

93      as MaAsLin 2 (Microbiome Multivariable Associations with Linear Models, with software version

94      2.0 released with this study), a successor to MaAsLin 1[15,16]. Unlike MaAsLin 1's single-model

95      framework based on applications of arcsine square root-transformed linear model following Total

96      Sum Scaling (TSS) normalization[15,16], MaAsLin 2 has evaluated and combined the best set of

97      analysis steps to facilitate high-precision association discovery in microbiome epidemiology

98      studies. It provides a coherent paradigm through a multi-model framework with arbitrary

99      coefficients (phenotypes and covariates) and contrasts of interest, along with support for data

100     exploration, normalization, and transformation options to aid in the selection of appropriate data-
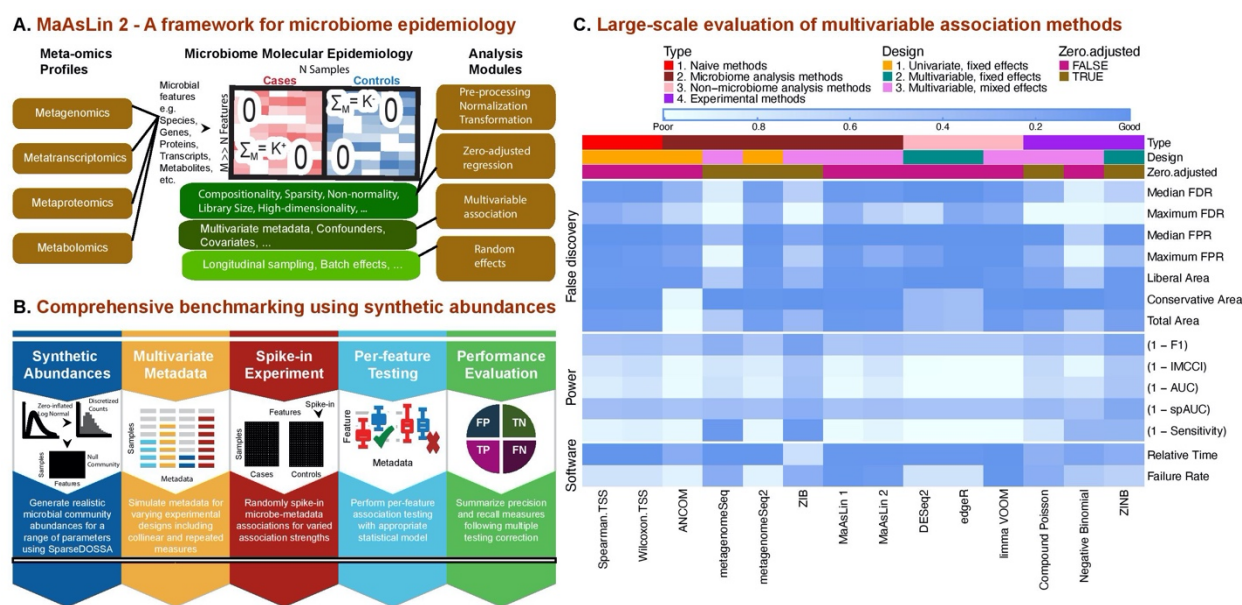
101   and design-driven statistical techniques for analyzing microbial multi-omics data. In this study, we

102   also conducted a large-scale synthetic evaluation of a broad range of circumstances under which

103   straightforward identification of meta-omics features can be challenging. These simulation studies

104   revealed that MaAsLin 2 preserves statistical power in the presence of repeated measurements

105   and multiple covariates while accounting for the nuances of meta-omics features and, critically,

106   controlling false discovery rates. We concluded with an application to novel biomarker discovery

107   in multiple omics datasets from the Integrative Human Microbiome Project (iHMP or HMP2[6]). The

108   implementation of MaAsLin 2, associated documentation and tutorial, and example data sets are

109   freely    available    in    the    MaAsLin    2    R/Bioconductor    software    package    at

110   https://huttenhower.sph.harvard.edu/maaslin2.


111   **Results**


112   **MaAsLin 2 methodology and validation**

113   MaAsLin 2 provides a comprehensive multi-model system for performing multivariable

114   association testing in microbiome profiles - taxonomic, functional, or metabolomic - with analysis

115   modules for preprocessing, normalization, transformation, and data-driven statistical modeling to

116   tackle the challenges of microbial multi-omics (compositionality, overdispersion, zero-inflation,

117   variable library size, high-dimensionality, etc.; **Fig. 1A**). The MaAsLin 2 implementation requires

118   two inputs: (i) microbial feature abundances (e.g. taxa, genes, transcripts, or metabolites) across

119   samples, in either counts or relative counts; and (ii) environmental, clinical, or epidemiological

120   phenotypes or covariates of interest (together "metadata"). Both metadata and microbial features

121   are first processed for missing values, unknown data values, and outliers. If indicated, microbial

122   measurements   are   then   normalized   and   transformed   to   address   variable   depth   of   coverage

123   across   samples.   Feature   standardization   is   optionally   performed,   and   a   subset   or   the   full

124   complement of metadata are used to model the resulting quality-controlled microbial features and

125    define p-values for each metadata association per feature using one of a wide range of possible

126    multivariable models. After all features are evaluated, p-values are adjusted for multiple

127    hypothesis testing and a table summarizing statistically significant associations is reported. While

128    the default MaAsLin 2 implementation uses a log-transformed linear model on TSS-normalized

129    quality-controlled data, the software supports several other statistical models including GLM (e.g.

130    Negative Binomial[17]), zero-adjusted models (e.g. Compound Poisson[18,19], ZINB[20]), and multiple

131    normalization/transformation schemes under one estimation umbrella. In the presence of

132    repeated measures, MaAsLin 2 additionally identifies covariate-associated microbial features by

133    appropriately modeling the within-subject (or environment) correlations in a mixed model

134    paradigm, while also accounting for inter-individual variability by specifying between-subject

135    random effects in the model. A variety of summary and diagnostic plots are also provided to

136    visualize the top results.



137

**Figure 1: MaAsLin 2 for feature-wise association of microbial communities with phenotypes. A)** MaAsLin 2 is a statistical method for association analysis of microbial community meta-omics profiles. It comprises several steps, including data transformation, multivariable inference, multiple hypothesis test correction, and visualization. These are based on a set of flexible and computationally efficient linear models, while accounting for the nuances of microbiome data, repeated measures, and multiple covariates. **B)** Comprehensive benchmarking of multivariable methods for microbiome epidemiology. To identify appropriate methods for associating microbiome features with health outcomes and other covariates, we assessed up to 84 combinations of normalization/transformation, zero-inflation, and regression models (**Supplementary Fig. S1A**). These were applied to synthetic data using a hierarchical model (SparseDOSSA, http://huttenhower.sph.harvard.edu/sparsedossa) to generate realistic, model-agnostic datasets with varying scopes and effect sizes of microbiome associations. Individual per-feature association methods were performed repeatedly to evaluate

147　method-specific recall and precision measures. **C)** Association method performance summary across major evaluation criteria. Three
148　aspects of performance were considered: (i) false discovery, (ii) sensitivity, and (iii) computational efficiency. Evaluation metrics are
149　shown (in rows) for the resulting microbial multivariable association methods (both state-of-the-art and novel), averaged over all
150　simulation parameters (**Supplementary Fig. S1B**). The top-performing methods (as measured by average F1 score) from each class
151　of models (**Supplementary Fig. S1C**) are shown (in columns). With the exception of Spearman and Wilcoxon that maintained best
152　performance on TSS-normalized data, all methods exhibited superior performance with no/default normalization (ANCOM,
153　metagenomeSeq, metagenomeSeq2, DESeq2, edgeR, MaAsLin 1, MaAsLin 2, limma VOOM, ZIB) or library size normalization in
154　which log-transformed library size is included as an offset in the associated GLM likelihood (Compound Poisson, Negative Binomial,
155　ZINB). Top colored boxes represent method characteristics including the capability to handle zero-inflation and random effects. Based
156　on synthetic evaluations, MaAsLin 2 includes optimized default models for epidemiological testing in microbial multi-omics data.

157　To identify model components appropriate for MaAsLin 2's microbiome feature association testing

158　and to objectively benchmark current association methods, we assessed realistic synthetic

159　datasets generated by SparseDOSSA (full details of individual association methods, as well as

160　simulation parameters, are described in **Methods** and are available online at

161　https://github.com/biobakery/maaslin2_benchmark). Briefly, SparseDOSSA is a synthetic data

162　generation routine that models biologically plausible synthetic data from diverse template

163　microbiome profiles by taking into account (i) feature-feature, (ii) feature-metadata, and (iii)

164　metadata-metadata correlations, superseding previous efforts by including multiple covariates

165　and longitudinal designs (**Methods**). As compared to previous simulation schemes,

166　SparseDOSSA allows multivariable spike-in both in the presence and absence of repeated

167　measures, as well as arbitrary covariance structure in the metadata design matrix.

168　For this study, we carried out several spike-in experiments to induce and test controlled

169　associations, as governed by configurable simulation parameters (**Supplementary Fig. S1**).

170　When used for this purpose, SparseDOSSA first generates null microbial community features

171　containing no significant association patterns using a Bayesian hierarchical model independently

172　of metadata features (**Fig. 1B**, **Methods**). In addition to varying sample size and feature

173　dimension, a broad range of metadata and experimental designs are then considered, including

174　repeated measures and univariate and multivariate covariates (both continuous and binary) of

175　varying dimension and effect size (**Supplementary Fig. S1A**). Specifically, in each instance, we

176　varied sample sizes from small (10) to large (200) for a fixed feature size (up to 500), and within

177　each sample size, the effect size parameter was again varied from modest (e.g. <2-fold

178 differences) to strong (10-fold). In each simulation, 10% of features (and 20% of metadata for

179 multivariable scenarios) were modified as an in-silico spike-in. Precision and recall measures

180 were averaged over 100 simulation runs (**Supplementary Fig. S1B, Methods)**. All methods were

181 corrected for multiple hypothesis testing using standard approaches for FDR control, declaring

182 significant associations at a target of FDR 0.05. For a fair comparison, a basic, model-free filtering

183 step to remove low-abundance features was performed before statistical modeling for all

184 methods. Methods unable to process specific simulation configurations due to high computational

185 overhead or slow convergence were omitted for those cases.

186 To compare the detection power of various methods in identifying true positive feature

187 associations, we first comprehensively evaluated published metagenomic tools and

188 representative methods for bulk RNA-seq analysis within each simulation scenario. These

189 methods were combined with several microbiome-appropriate normalization, transformation, and

190 linkage models (**Supplementary Fig. S1C, Methods**). In particular, we considered six distinct

191 categories of methods in our evaluations: (i) published methods specifically designed for microbial

192 community, such as metagenomeSeq[21], ANCOM[14,22], and ZIB[23,24], (ii) published bulk RNA-seq

193 differential expression methods, such as DESeq2[25], edgeR[26], and limma VOOM[27,28]; (iii) existing

194 generalized linear model (GLM) count models, such as the negative binomial[17], (iv) methods

195 based on linear models, such as limma[29] and "pure" linear models (LMs); (v) representative zero-

196 adjusted methods from the microbiome and scRNA-seq literature such as the compound

197 Poisson[18,19] and the zero-inflated negative binomial (ZINB[20,30]); and finally (vi) traditional,

198 simplistic non-parametric methods, such as Spearman correlation and Wilcoxon tests. Of note,

199 many of these methods can only compare two groups (i.e. a single binary metadatum), and not

200 all are compatible with continuous and multivariate metadata, resulting in a distinct set of

201 comparable methods for each experimental design.

202   Our first consideration in designing MaAsLin 2 for microbiome epidemiology was to ensure that

203   both current statistical theory and practical issues were considered during the analysis of

204   microbiome multi-omics data. To this end, we rigorously characterized various finite-sample

205   properties of different association methods focusing on three broadly defined aspects: (i) false

206   discovery, (ii) detection power, and (iii) software implementation, with multiple performance

207   indicators for each category (**Fig. 1C**). Rather than focusing on a single evaluation metric like the

208   Area Under the Curve (AUC) or the False Positive Rate (FPR), we ranked methods based on a

209   combination of metrics (**Methods**), many not considered in previous benchmarking. To

210   summarize each evaluation criteria, a normalized continuous score ranging between 0 and 1 was

211   assigned (**Methods**). Methods were then eliminated based on the presence of 'red flags' with

212   respect to at least one evaluation criteria, i.e. extreme departures from the best possible value.

213   Finally, metrics that are mainly descriptive rather than quantitative were also evaluated (e.g. the

214   ability to handle complex metadata designs, zero-inflation, or repeated measures) to achieve a

215   final consensus. This tiered strategy not only allowed us to select a set of "best" methods based

216   on the fewest 'red flags' across all scenarios, but also to identify a method that is (i) sufficiently

217   robust to false discovery control and detection power, (ii) scalable to large multi-omics datasets,

218   and (iii) accommodating of most modern epidemiological designs and microbial data types.

219   Notably, previous benchmarking in this area has only focused on differential abundance testing

220   without the simultaneous consideration of multiple covariates and repeated measures[7-9].

221   Additionally, with the exception of Hawinkel et al.[7], these benchmarking efforts lacked important

222   considerations to the extent that they either (i) did not consider FDR as a metric of evaluation[9,31,32]

223   or (ii) misreported false positive rate as FDR[8] (**Methods**). While a majority of these studies made

224   a final recommendation based on the traditional AUC metric or a combination of sensitivity and

225   specificity, we argue that without considering the FDR-controlling behavior of a method, the AUC

226   values alone are too optimistic to draw any meaningful conclusions about its practical utility. In

227   other words, particularly for biological follow-up, high precision among the most confident (lowest

228   recall) predictions is essential. To this end, our large-scale benchmarking enables a progressive

229   unification of traditional and practically important evaluation metrics by providing a comprehensive

230   connected view of microbiome multivariable association methods, especially in the context of

231   modern study designs, multiple covariates, and repeated measures.

232   Overall, our simulation study revealed that virtually all high-sensitivity methods with an

233   overoptimistic median AUC, especially those targeted to microbial communities, exhibited a highly

234   inflated average FDR (**Fig. 1C**). A similar pattern was observed for other AUC-like measures such

235   as F1 score and Matthew's correlation coefficient (MCC). On the other end of the spectrum,

236   compositionality-corrected methods such as ANCOM exhibited an extreme departure from 'good'

237   performance with respect to several criteria including sensitivity and p-value calibration, as

238   measured by both Conservative and Total Area[7] (**Methods**). Overall, these simulations reveal

239   that while there is no single method that outperforms others in all scenarios, MaAsLin 2 was the

240   only multivariable method tested that controlled FDR with the fewest 'red flags' across scenarios

241   (**Fig. 1C).**

242   This initial phase of our study thus expands the understanding of statistical association methods

243   appropriate for microbial community data under varying study designs, and it especially highlights

244   the inability of many common methods to control false discoveries. This is of critical importance

245   to past and present microbiome association methods, as failure to control the FDR causes

246   uncertainty in both scientific reproducibility and interpretability. Based on these evaluations, a

247   linear model with TSS normalization and log transformation was adopted as the default model in

248   MaAsLin 2, and the software provides several flexible options to apply a combination of other

249   normalization, transformation, and statistical methods to customize specific analysis tasks. The

250   implementation of MaAsLin 2, associated documentation, and example data sets are freely

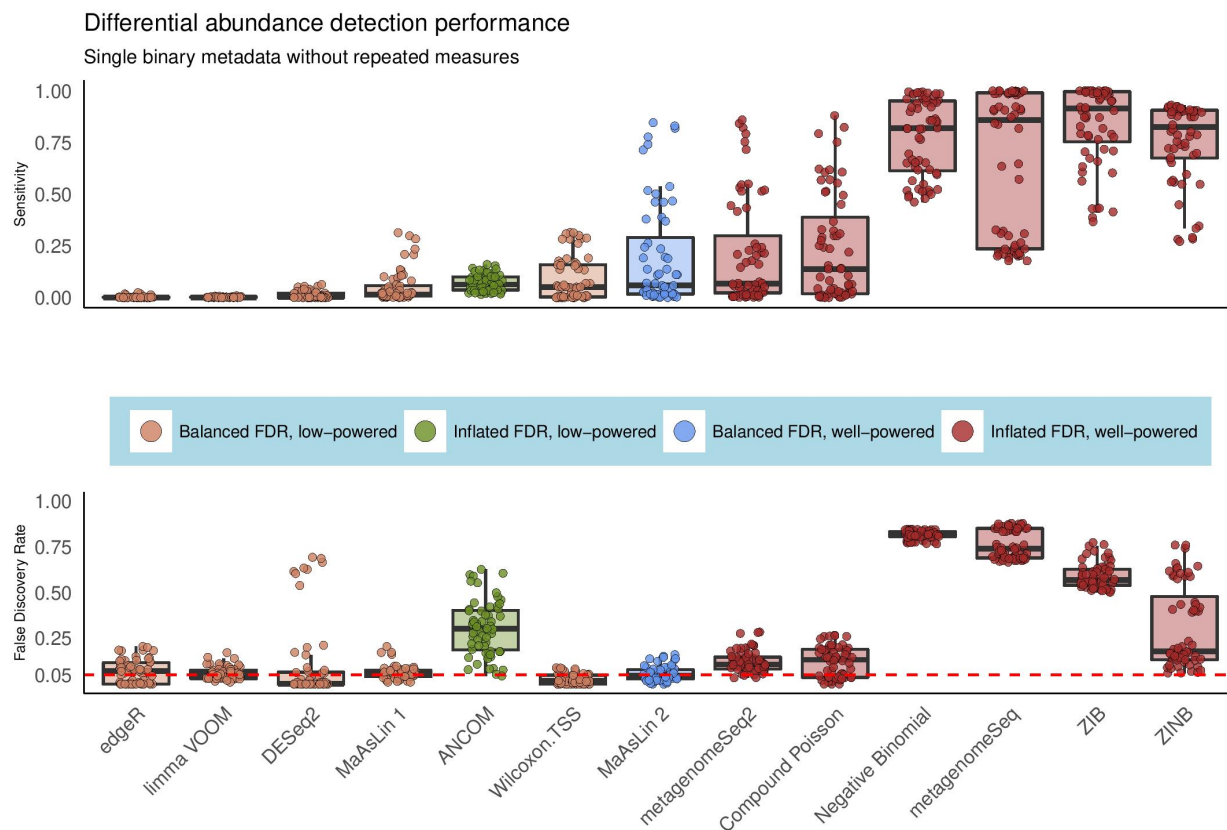251   available in R/Bioconductor and at https://huttenhower.sph.harvard.edu/maaslin2.

252 **MaAsLin 2 controls false discovery rate while maintaining power in differential**

253 **abundance analysis**

254 Differential abundance testing for microbial community features (taxa, pathways, etc.) is one of

255 the most commonly used strategies to identify features that differ between sample categories

256 such as cases and controls. Despite a large number of developments in the area, a lack of

257 consensus on the most appropriate statistical method has been a major concern[11]. To model

258 experimental designs of this type, we used synthetic count data with spiked-in features

259 differentially abundant between two defined groups of samples. In particular, we multiplied the

260 mean relative abundance of a randomly sampled fraction of 10% of the features with a given

261 effect size (fold change) in one of the groups and renormalized the ensemble of relative

262 abundances to a unit sum to create features differentially abundant between groups. We repeated

263 this procedure for each unique combination of sample size (10, 20, 50, 100, 200), feature

264 dimension (100, 200, 500), and fold change (1, 2, 5, and 10), each time summarizing performance

265 over 100 simulation runs (**Methods**). Before model fitting, features with a low prevalence (<10%)

266 were trimmed from the generated data sets.

267 As in our overall evaluation (**Fig. 1C**), we observed marked differences between the FDR-

268 controlling behavior of different methods in the simple case of single binary metadata and non-

269 longitudinal design, in some cases exceeding 75% (**Fig. 2**). Among the methods with good, robust

270 FDR control, only those based on linear models achieved moderate power, whereas, for methods

271 such as DESeq2 and edgeR, the FDR control came at the cost of reduced power. Among other

272 methods, practically all count and zero-inflated models, as well as newer methods based on log-

273 ratios such as ANCOM, struggled to correctly control the FDR at the intended (nominal) level, and

274 the best performance in this class of methods was obtained by metagenomeSeq2, Compound

275 Poisson, and ZINB (as measured by the F1 score). Many of the remaining methods were too

276 liberal, with metagenomeSeq and Negative Binomial standing out with a large number of false

277    positive findings. Overall, linear models (LMs) remained critically the only class of methods tested

278    that has good control of FDR across study designs, and they all exhibited a boost in statistical

279    power with increased sample size and association strength (**Supplementary Fig. S2).**

280    We also evaluated the average FPR of these methods by recording the fraction of tested

281    unassociated (negative) features that were deemed significant following FDR correction. Nearly

282    all methods controlled the FPR well below the imposed level (**Supplementary Fig. S3**). Relatedly,

283    we employed a previously proposed metric called "departure from uniformity" (i.e. departure from

284    uniformity of p-value under the null), which, complementary to FPR, quantifies the liberal or

285    conservative area between observed and theoretical quantiles of a uniform distribution[7]

286    (**Methods**). As expected, methods with high average false discovery rates, including zero-inflated

287    and count models, showed extreme departures from uniformity in the liberal direction, whereas

288    conservative methods such as DESeq2 and edgeR showed the same in the opposite direction,

289    suggesting extreme violation of uniformly distributed p-values under the null hypothesis

290    (**Supplementary Fig. S4**). While these results raise potential concerns about the FDR-controlling

291    behaviors of most existing methods, LM-based approaches did not exhibit this trend. In general,

292    tools based on linear models (such as limma) performed very similarly when calibrated with

293    MaAsLin 2's default model parameters, as expected, but not with their recommended default

294    parameters (**Supplementary Fig. S3**). Additionally, their options for handling sparsity and

295    compositionality were generally not appropriate for microbiome data. Amplicon, metagenomic

296    taxonomic, and functional profiles each show distinct count and zero-inflation properties, for

297    example, that are best handled by a multi-model system. As such, in addition to the binary

298    metadata design, we repeated the above simulation experiments for univariate continuous

299    metadata as well, which led to similar conclusions (**Supplementary Fig. S5**), further supporting

300    MaAsLin 2's default model's performance across metadata types and experimental designs.

Differential abundance detection performance

301

**Figure 2: MaAsLin 2 controls false discovery rate while maintaining power in differential abundance analysis of microbial communities**. To assess models' behaviors during differential abundance analysis, we simulated 100 independent datasets per parameter combination, each containing a single binary metadatum and a fixed number of true positive features (10% of features differentially abundant) for varying association strengths and sample sizes (**Supplementary Fig. S1C**). We then evaluated the ability of different microbiome association methods to recover these associations using a variety of performance metrics and summarized the results across runs (**Methods**). Both sensitivity and false discovery rates (FDR) are shown for the best-performing method from each class of models (as measured by average F1 score; **Methods**; full results in **Supplementary Figs. S2-S5**). Compared to zero-inflated and count-based approaches, MaAsLin 2's linear model formulation consistently controlled false discovery rate at the intended nominal level while maintaining moderate sensitivity. Red line parallel to the x-axis is the target threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score across all simulation parameters in this setting.

312    As a final evaluation, we assessed the impact of various normalization schemes on the associated

313    statistical modeling, evaluating all combinations of normalizations appropriate for each applicable

314    method (**Supplementary Fig. S1C, Methods**). Focusing on the best-performing linear models,

315    we found that model-based normalization schemes such as relative log expression (RLE[33]) as

316    well as data-driven normalization methods such as the trimmed mean of *M*-values (TMM[34]) and

317    cumulative sum scaling (CSS[21]) led to good control of FDR, but they also led to a dramatic

318    reduction in statistical power (**Supplementary Figs. S3, S5).** In contrast, TSS showed the best

319    balance of performance among all tested normalization procedures, leading to more powerful

13

320    detection of differentially abundant features. These results have potential implications for other

321    analyses in addition to differential abundance testing, as normalization is usually the first critical

322    step before any analysis of microbiome data, and an inappropriate normalization method may

323    severely impact post-analysis inference. In summary, our synthetic evaluation indicates that TSS

324    normalization, although simplistic in nature, may be superior to other normalization schemes

325    especially in the context of feature-wise differential abundance testing (and more generally for

326    multivariable association testing, as described later), in addition to community-level comparisons
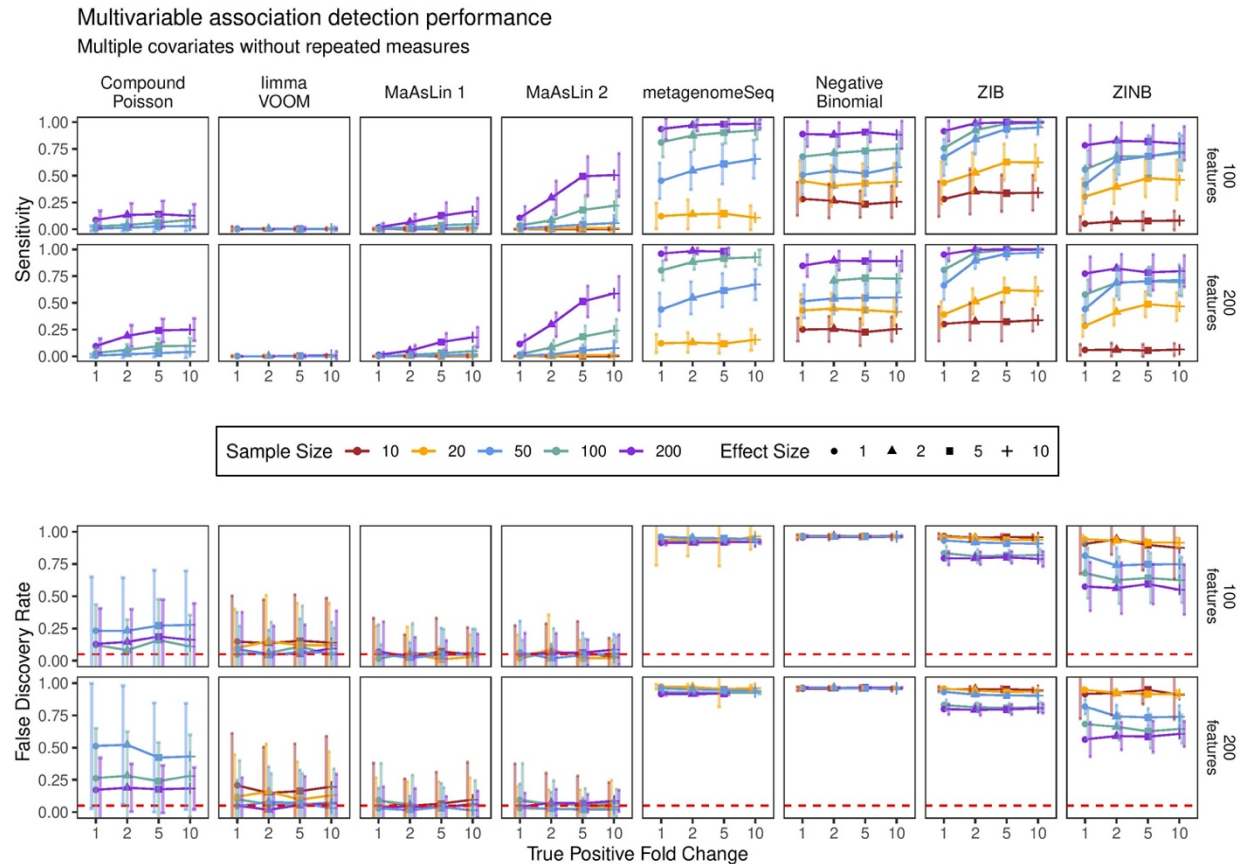
327    as previously described[35].

**MaAsLin 2 facilitates multivariable association discovery in population-scale**

**epidemiological studies**

330    Moving beyond univariate comparisons, we next assessed MaAsLin 2's performance in

331    multivariable association testing in comparison to other methods. Although widespread in

332    microarray and gene expression literature, multivariable analysis methods have remained

333    underdeveloped in microbial community studies. From an epidemiological point of view,

334    coefficients from a covariate-adjusted regression model are arguably more interpretable than its

335    individual, unadjusted counterparts. As a result, major conclusions from existing benchmarking

336    studies geared towards univariate associations are not generalizable to this broader setting,

337    where challenges such as zero-inflation and multiple testing are likely to be exacerbated,

338    especially in relation to multiple rounds of independently conducted univariate analyses as

339    commonly practiced.

340    To introduce multivariable associations into synthetically generated microbial feature profiles, we

341    supplemented each "sample" with multiple covariates consisting of both binary and continuous

342    metadata, either independent or correlated (**Supplementary Fig. S1A, Methods**). In each of

343    these datasets, 10% randomly selected features were modified ("spiked") to be associated with

344    randomly chosen 20% metadata features with a given magnitude (effect size). After spiking in,

345    samples were rescaled to their original (simulated) sequencing depth. As before, we repeated this

346    procedure for each unique combination of sample size (10, 20, 50, 100, 200), feature dimension

347    (100, 200, 500), and effect size (1, 2, 5, and 10), each time summarizing performance over 100

348    simulation runs.

349    The results from this set of simulations revealed that MaAsLin 2's default linear model had the

350    highest sensitivity among the methods that controlled the FDR at the target level, which also

351    remained consistent at larger sample sizes and stronger effect sizes (**Fig. 3**). We also observed

352    an improvement in performance when TSS normalization was employed (as compared to no

353    normalization) but did not observe similar improvement for other normalization methods

354    (**Supplementary Fig. S6)**. As before, zero-inflated and count models failed to control the FDR at

355    the nominal level, in the sense that the actual FDR was always above the nominal threshold used

356    for identifying significant features - a phenomenon that was surprisingly consistent regardless of

357    the metadata covariance structure (**Supplementary Fig. S7**). Taken together, these findings

358    further confirm that MaAsLin 2's default linear model is able to detect relevant associations across

359    a broad range of metadata designs, facilitating population-level analyses of microbial
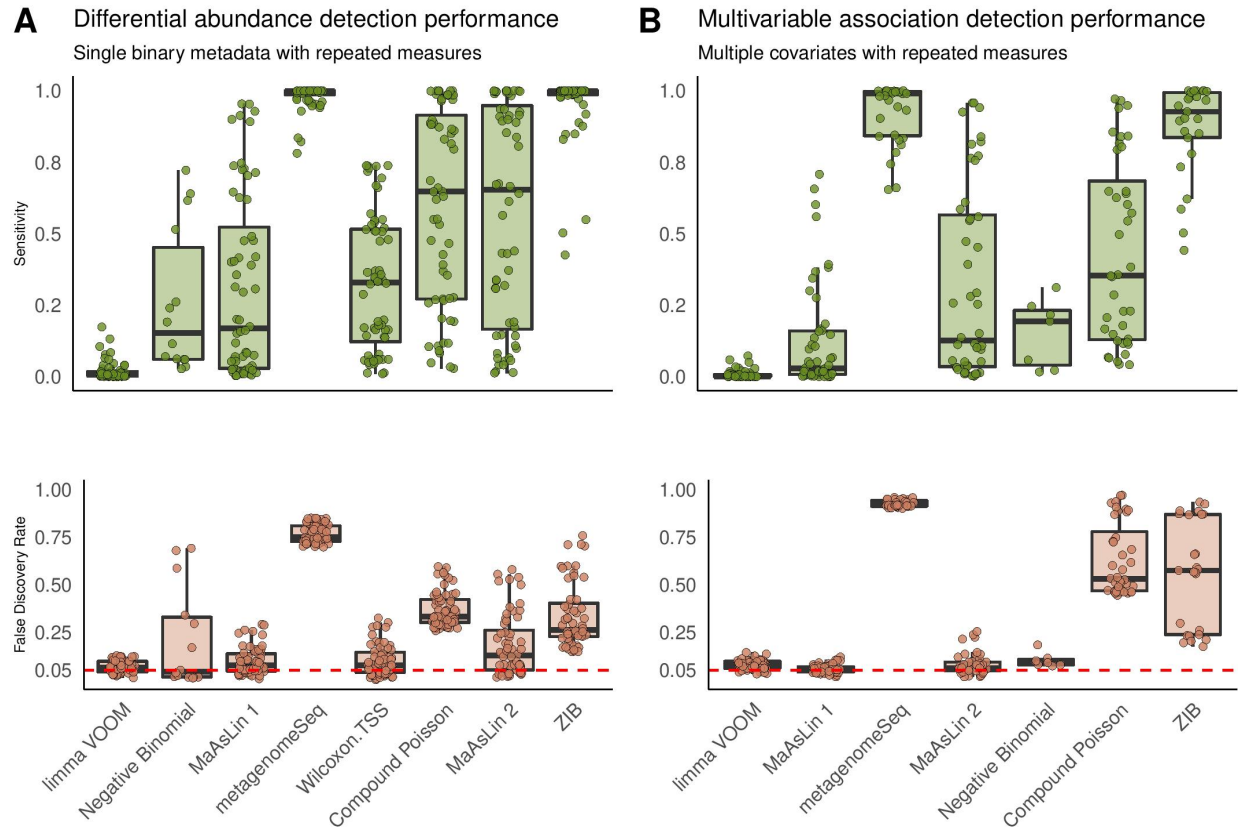
360    communities.

361

**Figure 3: MaAsLin 2 facilitates multivariable association discovery in large-scale human epidemiological and other microbial community studies**. Synthetic datasets containing five "metadata" with varying types of induced feature associations were analyzed using a variety of multivariable approaches (**Supplementary Fig. S1C**). As measured by power (recall) and false discovery rate (FDR), MaAsLin 2's default linear model outperformed other methods in controlling FDR while maintaining power across true-positive fold-change values, regardless of the total number of features. As expected, MaAsLin 2 has better power for stronger effect sizes, eventually attaining the highest power among all FDR-controlling methods. Red line parallel to the x-axis is the nominal FDR. Values are averages over 100 iterations for each parameter combination. The x-axis (effect size) within each panel represents the linear effect size parameter; a higher effect size represents a stronger association. For visualization purposes, the best-performing methods from each class of models (as measured by average F1 score; **Methods**; full results in **Supplementary Figs. S6-S7**) are shown. Methods are sorted by increasing order of average F1 score across all simulation parameters in this setting.

## MaAsLin 2 enables targeted microbiome hypothesis testing in the presence of repeated measures

To further validate MaAsLin 2 for longitudinal (or other repeated measures) microbiome data, we modified our simulation scheme to introduce subject-specific random effects - a key feature of modern microbiome population studies[36]. To this end, we tested MaAsLin 2 and related methods on two types of study designs. The first comprised univariate binary metadata designed to be challenging by the inclusion of non-independence of the data across time points. Second, we also simulated more realistic datasets using multiple independent covariates specific to longitudinal

16

380     microbiome studies. In both these regimes, realistic data were generated using SparseDOSSA

381     each with five time points[24], as in previous studies, but after introducing within-subject correlations

382     and between-subject random effects drawn from a multivariate normal distribution (**Methods**). It

383     is to be noted that the set of evaluable models is greatly reduced from the previous set of cross-

384     sectional association tests, as methods not capable of assessing repeated measures were

385     discarded.

386     Using these longitudinal synthetic "microbial communities," we compared the estimation and

387     inference from MaAsLin 2 with those of the existing methods, which revealed that MaAsLin 2 had

388     much lower false discovery rates than alternatives including ZIB (**Fig. 4, Supplementary Figs.**

389     **S8-S11**), a method specifically designed for microbiome longitudinal data. Both ZIB and MaAsLin

390     2's linear mixed-effects models are capable of identifying covariate-associated features by jointly

391     modeling all time points. However, the computational overhead of ZIB is significantly higher than

392     that of MaAsLin 2, which is prominent even for small datasets (**Supplementary Fig. S12)**.

393     Notably, although not nearly as severe as count-based and zero-inflated models, MaAsLin 2 had

394     a slightly inflated FDR in the univariate repeated measures scenario (**Fig. 4A**) but not in the

395     multivariable scenario (**Fig. 4B**). Among other methods, GLM-based methods such as Negative

396     Binomial and Compound Poisson performed similarly to their non-longitudinal counterparts for

397     both normalized and non-normalized counts (**Supplementary Figs. S8-S9**). This remained

398     consistent for both univariate continuous metadata (**Supplementary Fig. S10**) as well as multiple,

399     correlated covariates (**Supplementary Fig. S11**). Overall, these results suggest that MaAsLin 2's

400     linear mixed effect model consistently provides lower false discovery rates across metadata

401     designs and can effectively aid in testing differential abundance and multivariable association of

402     longitudinal microbial communities.

**Figure 4: MaAsLin 2 enables targeted microbial feature testing in the presence of repeated measures**. Results on simulated data comprising SparseDOSSA-derived compositions with five repeated measures per sample. The FDR is close to the target 0.05 level for MaAsLin 2's default linear model but not for zero-inflated and count models. As before, MaAsLin 2's linear model is consistently better powered than both negative binomial and limma VOOM at comparable FDR values, which remains consistent for both univariate continuous metadata (**A**) and multivariable mixed metadata designs (**B**) (a combination of continuous and binary covariates with five metadata features; **Methods,** full results in **Supplementary Figs. S8-S11**). The red line parallel to the x-axis is the given threshold for FDR in multiple testing. Within each panel, methods are sorted by increasing order of average F1 score across all associated simulation parameters in each setting.

## Multi-omics associations from the Integrative Human Microbiome Project

We applied MaAsLin 2 to identify relevant microbial features associated with the inflammatory bowel diseases (IBD) using longitudinal multi-omics data from the Integrative Human Microbiome Project (iHMP or HMP2[36]). The HMP2 Inflammatory Bowel Disease Multi-omics (IBDMDB) dataset included 132 individuals recruited in five US medical centers with Crohn's disease (CD), ulcerative colitis (UC), and non-IBD controls, followed longitudinally for one year with up to 24 time points each (**Methods**).

18

419   Integrated multi-omics profiling of the resulting 1,785 stool samples generated a variety of data

420   types including metagenome-based taxonomic profiles as well as metagenomic and

421   metatranscriptomic functional profiles, producing one of the largest publicly available microbial

422   multi-omics datasets. Metagenomes and metatranscriptomes were functionally profiled using

423   HUMAnN 2[37] to quantify MetaCyc pathways[38], and taxonomic profiles from metagenomes were

424   determined using MetaPhlAn 2[39] (**Methods**). For each of these data modalities (i.e. taxonomic

425   profiles and DNA/RNA pathways), independent filtering was performed before downstream

426   testing to reduce the effect of zero-inflation on the subsequent inference. In particular, features

427   for which the variance across all samples was very low (below half the median of all feature-wise

428   variances) or with >90% zeros were removed[36]. To further remove the effect of redundancy in

429   pathway abundances (explainable by at most a single taxon), only features (both DNA and RNA)

430   with low correlation with individual microbial abundances (Spearman correlation coefficient <0.5)

431   were retained.

432   We first used the IBDMDB to perform an additional semi-synthetic evaluation of association

433   methods' performance in "real" data, specifically when attempting to associate randomized, null

434   microbial taxonomic profiles to covariates. To this end, we permuted all samples 1,000 times to

435   construct shuffled "negative control" datasets, each time assessing the number of significant

436   associations (unadjusted $p$ <0.05) for each applicable method. These were averaged across

437   iterations to derive the expected number of null associations per method (which should remain

438   near-zero under usual circumstances). In particular, we fit (i) a baseline model as a function of

439   IBD diagnosis (a categorical variable with non-IBD as the reference group) while adjusting for

440   enrollment age (as a continuous covariate) and antibiotic use (as a binary covariate), and (ii) a

441   mixed effects model (with subject as random effects) with IBD dysbiosis state as an additional

442   time-varying covariate in addition to the time-invariant covariates considered in the baseline

443   model. Consistent with prior simulations, we found that several methods produced inflated

444 empirical type I error rates (**Supplementary Fig. S13)**. This conclusion remained unchanged

445 across varying significance thresholds, and as a result, we did not further apply these methods to

446 the non-permuted data. Relevantly and importantly, linear models did not suffer from this problem,

447 providing additional support for MaAsLin 2's robustness to false positive findings.

448 To dissect dysbiotic changes in IBD at greater resolution, we applied MaAsLin 2 to each individual

449 microbial feature type (i.e. species and DNA/RNA pathways) to test association with IBD

450 phenotype while controlling for IBD dysbiosis state, diagnosis, age, and antibiotic use (**Fig. 5**;

451 **Methods**). Nominal p-values for UC- and CD-specific effects were subjected to multiple

452 hypothesis testing correction using the Benjamini-Hochberg method[40] with an FDR threshold of

453 0.25. MaAsLin 2 identified a comparable number of significant associations with those initially

454 reported by the IBDMDB[36]. Among microbial species, MaAsLin 2's default linear model identified

455 206 significant associations, among which 150 (72.8%) overlapped with the original study

456 (**Supplementary Fig. S14**). MaAsLin 2 also reported many significant associations that were not

457 discovered in the original study (**Supplementary Dataset S1)**. For instance, we observed a

458 significant increase in *Bacteroides ovatus* in both UC and CD dysbiotic patients that was not

459 previously captured, as well as detecting (with MaAsLin 2's increased power) specific depleted

460 *Roseburia* species (*R. inulinivorans* and *R. hominis*) not captured by the previous analysis.

461 Notably, top hits from both MaAsLin 2 and the original study yielded nearly identical rankings

462 across data types, which broadly manifested as a characteristic increase in facultative anaerobes

463 at the expense of obligate anaerobes, in agreement with the previously observed depletion of

464 butyrate producers such as *Faecalibacterium prausnitzii* in IBD (**Fig. 5A**).

465 As an additional validation, we next re-analyzed the HMP2 taxonomic and functional profiles using

466 a zero-adjusted model (implemented in MaAsLin 2 as the compound Poisson). While this

467 maintained type-I error control in our shuffled data validation (as did the default linear model,

468 **Supplementary Fig. S13),** it was generally less desirable due to FDR inflation in simulations
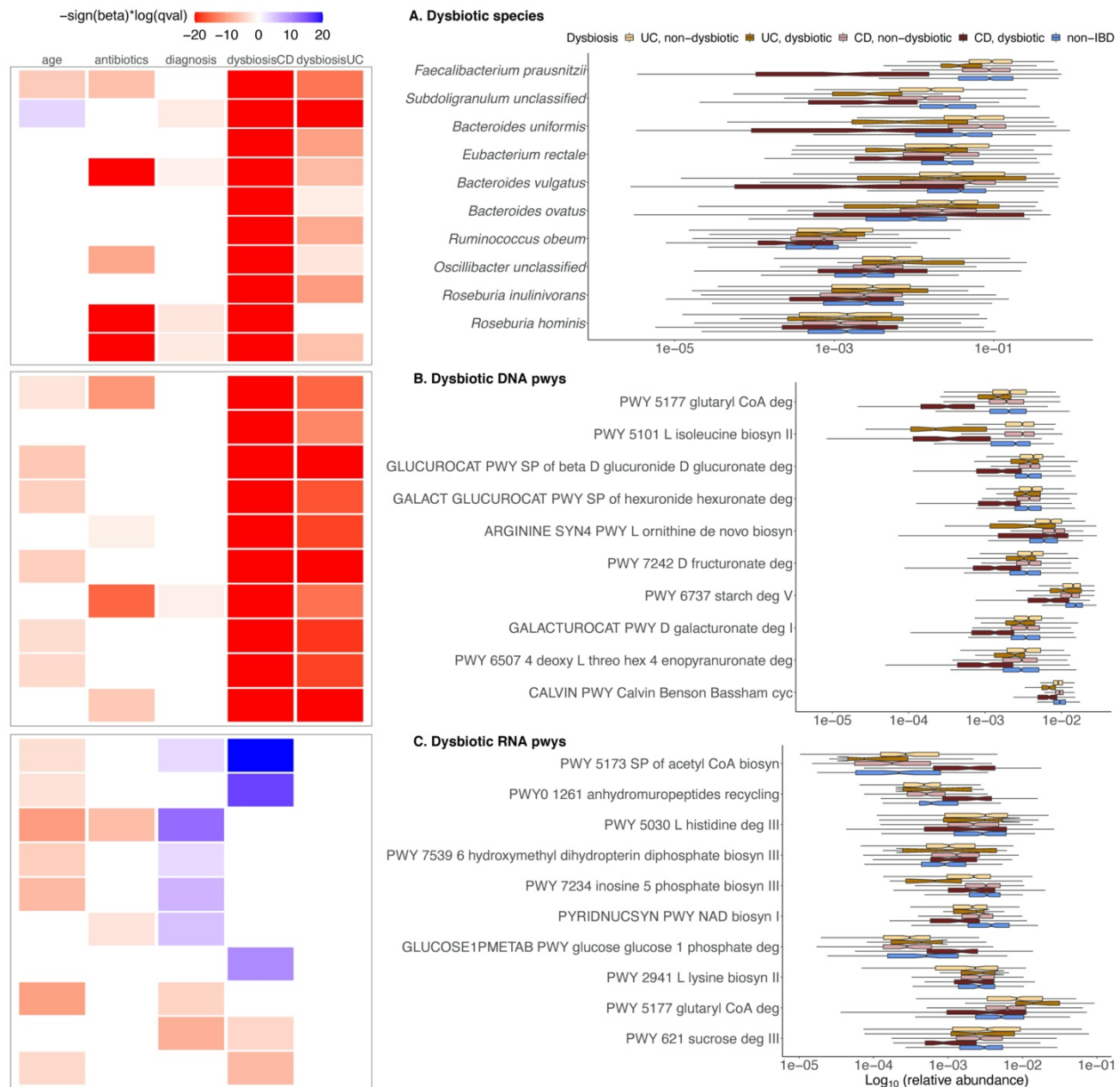
469    (**Figs. 2-4**). In terms of the number of differentially abundant features detected, both the default

470    linear model and the compound Poisson model performed similarly, with a significant overlap

471    between the top hits identified by each method (**Supplementary Fig. S15**). Among other

472    methods, ZIB and limma VOOM also maintained good Type I error control in these experiments,

473    although again both underperformed along other axes in our simulation studies. These results

474    further strengthen the finding that a combination of controlled parametric simulations and

475    'negative control' experiments based on data shuffling are useful together in identifying methods

476    for real-world applications, as the lack of either can lead to misleading (and irreproducible)

477    conclusions across independent evaluations[7]. This also highlights the flexibility of MaAsLin 2's

478    multi-analysis framework, wherein researchers are well-served with multiple (i) normalization

479    schemes, (ii) statistical models, (iii) multiplicity adjustments, (iv) multiple fixed and random effects

480    specifications, and (v) in-built visualization and pre-processing options, facilitating seamless

481    application of methods across diverse experimental designs under a single estimation umbrella.

482    Finally, in addition to taxonomic associations, MaAsLin 2 also detected 492 and 58 significant

483    functional associations for metagenomic (DNA) and metatranscriptomic (RNA) pathways,

484    respectively (**Supplementary Datasets S2-S3),** among which 358 (72.7%) and 39 (67.2%)

485    overlapped with the original study (**Supplementary Fig. S14**). While the original analysis of these

486    data included only community-wide functional profiles, we extended this by considering

487    metagenomic and metatranscriptomic functional profiles at both whole-community and species-

488    stratified levels in order to quantify overall dysbiotic functions while simultaneously assigning them

489    to specific taxonomic contributors. In particular, this considers a per-feature DNA covariate model,

490    in which per-feature normalized transcript abundance is treated as a dependent variable,

491    regressed on per-feature normalized DNA abundances along with other regressors in the model

492    (**Methods)**. Surprisingly, bioinformatics and statistics for metatranscriptomics are not yet

493    standardized, and our results indicate that subtle model variations can produce substantially

494    different results, due to the interactions between two compositions (DNA and RNA relative

495    abundances, **Supplementary Dataset S4).** This novel modeling strategy thus led to the discovery

496    of several novel transcript associations relative to the original study.

497    In a majority of these pathways, functional perturbations were driven by shifts in their

498    characteristic contributing taxa (**Fig. 5B**). For example, the most significant DNA pathways

499    enriched in CD patients were characteristic of facultative anaerobes such as *Escherichia coli*,

500    which are broadly more abundant during inflammation. These included pathways such as

501    synthesis of the enterobactin siderophore, lipid A, and sulfate reduction. A second set of enriched

502    pathways was depleted due to the loss of microbes such as *F. prausnitzii*, a particularly prevalent

503    organism that, when abundant, tended to contribute the majority of all enriched pathways it

504    encodes in this cohort (e.g. synthesis of short-chain fatty acids and various amino acids).

505    With the increased sensitivity of this analysis for species-stratified pathways, the overwhelming

506    majority of significant metagenomic differences were attributable solely to the most differential

507    individual organisms, as expected (**Supplementary Datasets 5-6**). Essentially every pathway

508    reliably detectable in *E. coli* was enriched during CD, UC, or both, and most *F. prausnitzii*

509    pathways depleted, along with many pathways from other gut microbes common in "health"

510    (*Bacteroides vulgatus*, *B. ovatus*, *B. xylanisolvens*, *B. caccae*, *Parabacteroides* spp., *Eubacterium*

511    *rectale*, several *Roseburia* spp., and others). Interestingly, since both more potentially causal

512    "driver" pathways, along with all other "passenger" pathways encoded by an affected microbe,

513    are detected by this more sensitive stratified analysis, it can be in many ways more difficult to

514    interpret than the non-stratified, community-wide, cross-taxon metagenomic responses to broad

515    ecological conditions such as immune activity, gastrointestinal bleeding, or oxygen availability.

**Figure 5: Multi-omics associations from the Integrative Human Microbiome Project. A)** Top 10 significant associations (FDR < 0.25) detected by MaAsLin 2's default linear model (significance and coefficients in **Supplementary Datasets S1-S3.**). All detected associations are adjusted for subjects and sites as random effects and for other fixed-effects metadata including the subject's age, diagnosis status (CD, UC, or non-IBD), disease activity (defined as median Bray-Curtis dissimilarity from a reference set of non-IBD samples), and antibiotic usage. Representative significant associations with dysbiosis state from each 'omics profile are shown: species (**B**), metagenomic (DNA) pathways (**C**), and metatranscriptomic (RNA) pathways (**D**). Values are log-transformed relative abundances with half the minimum relative abundance as pseudo count; full results in **Supplementary Datasets S1-S3.**

Conversely, differentially abundant microbe- and pathway-specific transcript levels highlighted a much more specific and dramatic shift toward oxidative metabolism, away from anaerobic fermentation, and towards Gram-negative (often *E. coli*) growth during inflammation (**Fig. 5C**)[41]. Many of these processes were either more extreme during (e.g. gluconeogenesis) or unique to

23

528    (e.g. glutathione utilization) active CD, as compared to UC. CD and UC responses were opposed

529    in a small minority of cases (e.g. glutaryl-CoA degradation). When stratified among contributing

530    taxa, these differences were almost universally attributable to a few key species, particularly an

531    increase in *E. coli* activity during inflammation and decreases of *F. prausnitzii* transcript

532    representation. Condition-specific transcriptional changes were occasionally contributed (or not)

533    by "passenger" *Bacteroides* spp. (*B. fragilis*, *B. xylanisolvens*, *B. dorei*) instead. Note that these

534    differences include pathways more likely to be "causal" in some sense, as significant

535    transcriptional changes were generally a subset of those detected due to whole-taxon shifts in

536    DNA content (including housekeeping pathways such as general amino acid or nucleotide

537    biosynthesis). These findings further support the importance of disease-specific transcriptional

538    microbial signatures in the inflamed gut relative to metagenomic profiles of functional potential,

539    suggesting that a potential loss of species exhibiting altered expression profiles in disease may

540    have more far-reaching consequences than suggested by their genomic abundances alone.


## 541    Discussion

542    A longstanding goal of microbial community studies, be they for human epidemiology or

543    environmental microbiomes, is to identify microbial features associated with phenotypes,

544    exposures, health outcomes, and other important covariates in large, complex experimental

545    designs. This parallels other methods for high-throughput molecular biology, but microbial

546    community multi-omics must account for properties such as variable sequencing depth, zero-

547    inflation, overdispersion, mean-variance dependency, measurement error, and the importance of

548    repeated measures and multiple covariates. To this end, we have developed and validated a

549    highly flexible, integrated framework utilizing an optimized combination of novel and well-

550    established methodology, MaAsLin 2. This accommodates a wide variety of modern study

551    designs ranging from within-subject, longitudinal to between-subject, cross-sectional, diverse

552    covariates, and a range of quality control and statistical analysis modules to identify statistically

553    significant as well as biologically relevant associations in a reproducible framework. The

554    embedding of these strategies in the paradigm of generalized linear and mixed models enables

555    the treatment of both simple and quite complex designs in a unified setting, improving the power

556    of microbial association testing while controlling false discoveries. To validate this framework, we

557    have extensively evaluated its performance alongside a set of plausible methods for differential

558    abundance analysis in a wide range of scenarios spanning simple univariate to complex

559    multivariable with varying scopes and effect sizes of microbiome associations. Finally, we applied

560    MaAsLin 2 to identify disease-associated features by leveraging the HMP2's multi-omics profile

561    of the IBD microbiome, confirming known associations and suggesting novel ones for future

562    validation.

563    A unique aspect of our synthetic evaluation of microbial community feature-wise association

564    methods while developing MaAsLin 2 is their comprehensive assessment in the presence of

565    multiple covariates and repeated measures, an increasingly common characteristic of modern

566    study designs. To identify covariate-associated microbial features from longitudinal, non-

567    independent measurements, it is necessary to jointly model data from all time points and

568    appropriately account for the within-subject correlations while allowing for multiple covariates.

569    This is particularly critical in the human microbiome, where baseline between-subject differences

570    can be far greater than those within-subjects over time, or of the effects of phenotypes of interest.

571    To the best of our knowledge, the synthetic evaluation presented here is the first to consider such

572    aspects of large-scale microbiome epidemiology in statistical benchmarking. This enabled us to

573    investigate key aspects of published methods that would be difficult to generalize from univariate

574    comparisons alone[7-9]. Note that the resulting conclusion is largely independent of the association

575    models being evaluated, as the synthetic data were generated from an additional, completely

576    external model (i.e. the zero-inflated log-normal, **Methods**), which is fundamentally different from

25

577    any of the evaluated parametric methods. Our simulation results thus complement the findings of

578    previous studies in several important aspects. Consistent with previous reports, nearly all zero-

579    inflated models suffer from poor performance (i.e. inflated false positives and higher computation

580    costs), here in both univariate and multivariable scenarios with or without repeated measures.

581    This calls for methodological advancements in statistical modeling of zero-inflated data, as

582    existing theory seems to differ very surprisingly from practice when implemented by established

583    optimization algorithms and applied to noisy data.

584    One noteworthy finding of our evaluation is that a random effect implementation of the same

585    underlying statistical model can lead to different substantive conclusions than its fixed-effects

586    counterpart. This was particularly evident for the negative binomial case, where a substantially

587    better control of FDR (albeit inflated) was observed for the random effect analog. Interestingly,

588    the negative binomial model (with or without zero inflation) is in many ways considered the most

589    "appropriate" model for count-based microbial community profiles, but we observed extremely

590    inconsistent behavior for negative binomial and ZINB implementations during our evaluation, as

591    also observed in previous findings[42]. In particular, our negative binomial evaluation used the

592    *glm.nb()* function from the *MASS* R package[43] for fixed-effects and the *glmer()* function from the

593    R package *lme4*[44] for random effects, whereas the ZINB evaluation used the *zeroinfl()* function

594    from the R package *pscl*[45]. This additionally highlights the potential reproducibility concerns

595    induced by differences in algorithms, implementations, and computational environments even for

596    the same underlying model, suggesting that great caution should be taken when interpreting

597    multiple implementations of the same statistical model for challenging microbial community

598    settings in the absence of an experimentally validated gold standard.

599    In agreement with previous studies, we confirmed that most RNA-seq differential expression

600    analysis tools tend to provide suboptimal performance when applied unmodified to zero-inflated

601    microbial community profiles. In particular, count-based models, due to their strong parametric

602    assumptions on the distributions or parametric specifications of the mean-variance dependency,

603    tend to have inflated FDR when the assumptions are violated. In sharp contrast to previous claims,

604    however, compositionality-corrected methods such as ANCOM[14,22] as well as specialized

605    normalization and transformation methods such as CLR[46] did not improve performance over non-

606    compositional approaches[8,47], consistent with recent findings that compositional methods may not

607    always outperform non-compositional methods[32]. Importantly, these conclusions hold regardless

608    of the nature of the modeling paradigm (i.e. univariate vs. multivariable), thus providing a

609    generalizable benchmark for future evaluation studies of applied microbiome association

610    methods. Though we primarily focused on data generated in microbial community surveys, many

611    of our conclusions are extendible to similar zero-inflated count data arising in other applications

612    such as single-cell RNA-seq. Taken together, these simulation results revealed that further

613    investigation into the causes of the failure of FDR correction and development of specialized false

614    positive-controlling methods are important upcoming challenges in microbiome statistical

615    research.

616    Limitations of the current MaAsLin 2 method include, first, its restriction to associating one feature

617    at a time. While this strategy has the advantage of being straightforward to interpret, implement,

618    and parallelize, it sacrifices inferential accuracy by ignoring any correlation structure among

619    features. This can certainly exist due both to compositionality and to biology and will differ e.g.

620    between taxonomic features (related by phylogeny) vs. functional ones (such as pathways). A

621    potential extension would be to adopt an additional multivariate framework that allows modeling

622    simultaneously rather than sequentially, thus improving power by borrowing strength across non-

623    independent features. Second, as revealed by our synthetic evaluation, not surprisingly, linear

624    models remain underpowered in detecting weak effects among microbial communities, especially

625    when accompanied by a small sample size. This is in some ways a necessary consequence of

626    the restrictions of current microbiome measurement technologies, and it emphasizes the

627    importance of an informed power analysis before study planning to ensure an optimal sample size

628    with adequate detection power. Finally, and relatedly, it is not straightforward to incorporate any

629    type of graph structure knowledge such as phylogeny or pathway-based functional roles into the

630    per-feature linear model framework. Bayesian linear models can potentially improve on this by

631    including such information through a suitable prior distribution.

632    Several aspects of microbiome epidemiology remain to be investigated both biologically and

633    computationally, in addition to the challenges addressed here. For example, while it is possible to

634    obtain strain-level resolution from metagenomic sequencing data, strain variants are generally

635    unique to particular individuals rather than broadly carried by many people, presenting unique

636    challenges for strain-level multi-omics. From a computational point of view, this calls for further

637    refinements to MaAsLin 2's methodology when applied to strain-resolved community profiles. In

638    addition, the modeling framework adopted here can only inform undirected associations, and

639    hence cannot be used to infer causation. Advanced methods from other molecular epidemiology

640    fields such as causal modeling and mediation analysis methods can help overcome these

641    issues[48]. Another opportunity for future extension of our method is the integration of established

642    missing data imputation methods across features and metadata, a common pitfall in many

643    molecular epidemiology studies[36]. Combined, such extensions will lead to further improvement in

644    downstream inference, allowing researchers to investigate a range of hypotheses related to

645    differential abundance and multivariable association.

646    As currently implemented, MaAsLin 2 is designed to be applicable to most human and

647    environmental microbiome study designs, including cross-sectional and longitudinal. Clearly,

648    these can also be extended to additional designs, such as nested case-control and case-cohorts.

649    It is to be noted that MaAsLin 2's capability extends well beyond association analysis. For

650    instance, MaAsLin 2's multi-analysis framework has been used in the context of meta-analysis[49],

651    and the extracted residuals and random effects from a MaAsLin 2 fit can be used for further

28

652   downstream analysis (e.g. as has been done in the original HMP2 study for cross-measurement

653   correlation analysis[36]). By adhering to a flexible mixed effects framework, MaAsLin 2 is able to

654   analyze multiple groups and time points jointly with other associated covariates, which allows

655   formulation of both fixed effects (for cross-sectional associations) and random effects (for within-

656   subject correlations) in a single unified framework. This is particularly appropriate for non-

657   longitudinal studies (those with a small number of repeated measures e.g. multiple tissues or

658   families), or from sparse and irregular longitudinal data from many subjects (e.g. with unequal

659   number of repeated measurements per subject, as commonly encountered in population-scale

660   epidemiology). This aspect could also be extended in the future, based on the increasing

661   availability of dense time-series profiles appropriate for non-linear trajectory-based methods from

662   Bayesian nonparametrics, such as Gaussian processes, particularly in the presence of multiple

663   covariates[5,50]. Finally, methods need to be developed to accommodate the increasing availability

664   of microbiome-host interactomics and electronic health records in population-scale microbiome-

665   wide epidemiology[6], moving beyond observational discovery toward translational applications of

666   the human microbiome. In summary, the methodology presented here provides a starting point

667   for more efficient identification of microbial associations from large microbial community studies,

668   offering practitioners a wide set of analysis strategies with state-of-the-art inferential power for the

669   human microbiome and other complex microbial environments.

670   **Methods**

671   **Data for differential feature model evaluations**

672   <u>Synthetic null community abundances</u>

673   Realistic null community data were generated using the SparseDOSSA[51] (Sparse Data

674   Observations   for   the   Simulation   of   Synthetic   Abundances)   hierarchical   model

675  (http://huttenhower.sph.harvard.edu/sparsedossa). SparseDOSSA is a newly developed

676  simulator designed to model the fundamental characteristics of real microbial communities (e.g.

677  zero-inflation, compositionality, etc.) and to simulate new, realistic metagenomic count data with

678  known feature-feature and feature-metadata correlations and provide a gold standard to enable

679  benchmarking of statistical metagenomics methods, superseding previous efforts by including

680  multiple covariates and longitudinal designs.

681  Briefly, SparseDOSSA's Bayesian model captures microbial features (taxon, gene, or pathway

682  abundances) as truncated, zero-inflated log-normal distributions, the parameters of which are

683  hierarchically derived from a parent log-normal distribution. SparseDOSSA estimates feature-

684  specific parameters by fitting to a real-world template dataset, and generates synthetic features

685  from zero-inflated, truncated log-normal distributions based on both fitted and user-defined

686  parameters on a per feature basis (**Supplementary Fig. S1A**). All feature-specific parameters,

687  namely the log-mean, zero-inflation proportion, truncation point, and log-variance are empirically

688  determined to resemble the template dataset's properties. After sampling, the samples are

689  rounded to the nearest integer to mimic count data. A combined dataset of the RISK[52], PRISM[15],

690  pouchitis[16], and NLIBD[53] gut microbiomes, totaling several thousand samples, was used as

691  empirical microbiome template data for the simulations reported in this study. To mimic realistic

692  variation in library size, sequencing depth was generated from a lognormal distribution with

693  average sequencing depth 50,000, resulting in approximately 30-fold to 100-fold variation in

694  sequencing depth.

695  <u>Synthetic metadata generation</u>

696  Simulated metadata matrices in simple univariate cases (UVA, UVB) were generated with

697  continuous values from a standard normal distribution. For the univariate binary case (UVB), we

698  additionally dichotomized the continuous variable by coding samples in the bottom and top half

699  of the distribution as 0 and 1, respectively. For multivariate cases (MVA, MVB), we repeated the

700    above discretization for multiple metadata by first generating from a multivariate normal

701    distribution, and concurrently binarizing half of the metadata features at random. We considered

702    two frequently encountered correlation structures for the multivariate cases: independent and AR

703    (1) with coefficient 0.5, which correspond to MVA and MVB, respectively. Additionally, we

704    considered repeated measures by incorporating random effects in these cross-sectional design

705    matrices. To that end, we generated a simple blocking variable that is normally distributed (with

706    mean 0 and variance 1) across subjects but invariant within subjects, representing a single

707    random effect factor such as subject or time point (block size determined by the simulation

708    parameters as reported in **Supplementary Fig. S1A**). Subsequently, we added this as an

709    additional covariate to the fixed-effects metadata to impose correlations within the blocks,

710    mimicking longitudinal studies. For multivariable cases (MVA, MVB), the number of covariates is

711    fixed to 5. Similarly, for the repeated measures settings, $T = 5$ time points per subject is

712    considered.

713    <u>Multivariable spike-ins of synthetic feature-metadata associations</u>

714    To introduce associations between features and metadata, we used SparseDOSSA's default

715    additive spike-in procedure. Following Weiss et al.[8], we implement the spike-ins in a balanced

716    way across all metadata to avoid compositional bias. Briefly, SparseDOSSA standardizes both

717    (microbial) features and metadata and randomly chooses (microbial) null features and metadata

718    without replacement. The standardization procedure ensures that the spiked-in features are not

719    dominated by the values of the target metadata but rather distributed similarly to the real data.

720    Next, the standardized non-zero abundances of the selected features are modified by adding a

721    linear combination of all spiked-in standardized metadata, in which a real-valued effect size

722    parameter (**Supplementary Fig. S1A**) governs the strength of association for each associated

723    feature-metadata pair. To create differentially abundant features, a randomly sampled fraction of

724    10% of the features are spiked-in to be associated with the metadata. In the multivariable case,

725      20% of the metadata are randomly selected to be associated with the 10% 'differentially abundant'

726      features.

727      **Multivariable association test evaluation**

728      <u>Preprocessing, normalizations, and transformations</u>

729      We considered several commonly used normalization procedures including Total Sum Scaling

730      (TSS), Trimmed Mean of *M*-values (TMM[34]), Relative Log Expression (RLE[33]), and Cumulative

731      Sum Scaling (CSS[21]) (**Supplementary Fig. S1C**). For TSS normalization, raw counts were

732      converted into relative abundances by scaling each sample by the total sum (across features).

733      For the remainder, we used the default settings of the edgeR[26], DESeq2[25], and metagenomeSeq[21]

734      R packages, respectively.

735      In addition to the above normalization procedures, several parametric transformations were also

736      considered. When appropriate, these variance-stabilizing transformations aim at improving

737      parametric estimation models in the presence of violated data assumptions (such as normality

738      and homoscedasticity). These include logit and arcsine square root (AST) for TSS-normalized or

739      proportional relative abundance data, and log transformation (**Supplementary Fig. S1C**). For

740      both log and logit transformations, undefined values were replaced with zeroes (equivalent to

741      adding a small pseudo count of 1 to the zero observations before applying the log transformation).

742      Among other normalization/transformation methods, a 'Default/None' category was also

743      considered which represents either (i) default normalization/transformation employed by the

744      associated software or (ii) no normalization/transformation, or (iii) library size normalization in the

745      form of a GLM offset modeling. Prior to applying any normalization and transformation, a basic

746      filtering was performed to prune features absent in >90% of samples. As in previous

747      benchmarking[7,8,10], correction for multiple testing was performed using the Benjamini-Hochberg[40]

748      FDR threshold of 0.05.

749    Statistical methods

750    We selected several commonly used methods for differential abundance and multivariable

751    association testing along with a set of experimental methods to apply on the synthetic datasets,

752    using a combination of statistical model and normalization/transformation schemes for each

753    method as appropriate (**Supplementary Fig. S1C**). All tests were conducted using the statistical

754    software R and parallelized using custom bash scripts in a high-performance computing

755    environment (full source code available at: https://github.com/biobakery/maaslin2_benchmark).

756    The selected statistical models (abbreviations in parentheses) are as follows:

757    • ANCOM: following Weiss et al.[8], we used the default implementation of ANCOM[14] using

758        the *ANCOM()* function call with default settings. Unlike other methods, ANCOM does not

759        report p-values but instead returns logical indicators of whether a feature is differentially

760        abundant based on a test statistic W. It is to be noted that in the presence of multiple

761        covariates, ANCOM does not return statistically significant feature-metadata pairs with

762        respect to every covariate in the model, making it infeasible for our multivariable setting.

763        Also, we did not test the ANCOM method in repeated measures settings as it was too slow

764        and unstable for assessment, as noted elsewhere[32].

765    • metagenomeSeq: for the fixed effects, counts were first CSS-normalized with the default

766        quantile supplied by the *cumNormStat()* function and the (log-transformed) CSS-

767        normalized counts were subjected to final testing using *fitZig*()[21]. For random effects,

768        *useMixedModel was* set to TRUE in the *fitZig*() function call.

769    • metagenomeSeq2: same as metagenomeSeq[21], except that the final testing was done

770        using the *fitFeatureModel*() function.

771    • DESeq2: for fixed effects, following Thorsen et al.[9], geometric means were first calculated

772        manually from the raw counts and supplied to the *estimateSizeFactors()* function before

33

773      calling the *DESeq()* function for final testing[25]. Random effects modeling compatible with

774      our setting is currently not supported by the DESeq2 software[54].

775      • edgeR: for fixed effects, following Thorsen et al.[9], normalization factors were calculated

776      with TMM, which was followed by common and tagwise dispersion estimation steps,

777      before invoking the standard test with the *exactTest()* function[26]. Random effects modeling

778      compatible with our setting is currently not supported by the edgeR software[54].

779      • limma: the default functionality of *lmFit()* was applied to the feature counts[29]. Repeated

780      measures were handled using the *duplicateCorrelation()* function before calling *lmFit(),* in

781      combination with appropriate normalization/transformation (**Supplementary Fig. S1C**).

782      • limma VOOM: same as limma, except that features were subjected to a voom

783      transformation before applying limma[27,28].

784      • limma2: same as limma, except that library size or scale factor is included as an additional

785      covariate, in combination with appropriate normalization/transformation (**Supplementary**

786      **Fig. S1C**).

787      • Wilcoxon: the built-in R function *wilcox.test()* using default parameters was applied to the

788      features in combination with appropriate normalization/transformation (**Supplementary**

789      **Fig. S1C**).

790      • Spearman: the built-in R function cor.*test()* was applied to the features in combination with

791      appropriate normalization/transformation (**Supplementary Fig. S1C**).

792      • Linear model (LM): the built-in R function *lm()* with default settings was used in

793      combination with appropriate normalization/transformation (**Supplementary Fig. S1C**).

794      • Linear model (LM2): same as LM, except that library size or scale factor is included as an

795      additional covariate in the model, in combination with appropriate

796      normalization/transformation (**Supplementary Fig. S1C**).

797    • Negative binomial (negbin): we used the *glm.nb()* function from the *MASS* package[43] and

798       the *glmer.nb()* function from the *lme4* package[44] for fixed and random effects respectively.

799       In both cases, we used the logarithm of library size (for no normalization) or scaling factor

800       (for other normalization schemes such as CSS, RLE, and TMM) as offset.

801    • Zero-inflated Negative Binomial (ZINB): for fixed effects, we used the *zeroinfl()* function

802       from the  *pscl* package[45] with the logarithm of library size (for no normalization) or scaling

803       factor (for other normalization schemes such as CSS, RLE, and TMM) as offset. In the

804       absence of a robust random effect implementation of the same, the ZINB method was not

805       considered in the repeated measures settings.

806    • Zero-inflated Beta (ZIB): following Peng et al.[23], we used the *gamlss()* function from the  R

807       package *gamlss*[55] for fixed effects and the *ZIBR()* function from the *ZIBR* R package for

808       random effects[24]. In both cases, the features are TSS-normalized before statistical testing.

809    • Compound Poisson (CPLM): we used the *cpglm()* function from the *cplm* package[56] and

810       the *glmmPQL()* function from the *MASS* package[56] for fixed and random effects

811       respectively. In both cases, we used the logarithm of library size (for no normalization) or

812       scaling factor (for other normalization schemes such as CSS, RLE, and TMM) as offset.

813       No offset was used when combined with the TSS-normalized relative counts.

814    • MaAsLin 1: we used the default TSS-normalized, arcsine square root-transformed linear

815       model without gradient boosting[15,16].

816    • MaAsLin 2: we used the default TSS-normalized, log-transformed linear model with half

817       the minimum relative abundance as pseudo count.

818    Naming convention

819    The nomenclature for the model/normalization/transformation combinations for each method is

820    described in the following set of rules:

35

821    1. For published methods with default parameters, there is no additional identifier following

822         the name of the method, indicating default or no normalization/transformation. These

823         include ANCOM, metagenomeSeq, metagenomeSeq2, limma, limma2, limma VOOM,

824         DESeq2, edgeR, and ZIB.

825    2. Similarly, for experimental methods with custom normalization/transformation schemes,

826         no additional identifier simply indicates either no normalization (for non-GLM methods

827         such as LM) or library size normalization (for specific GLMs such as Negative Binomial,

828         Compound Poisson, and ZINB).

829    3. Finally, for methods with additional identifiers, method name is always accompanied by a

830         normalization scheme (after the first dot) which is followed by a transformation (after the

831         second dot) except in cases where either no normalization or no transformation is applied.

832         As an example, limma.CSS.LOG denotes a default limma model followed by CSS

833         normalization and log transformation. Similarly, LM.CLR denotes a vanilla linear model

834         followed by a CLR transformation and no normalization, whereas, ZINB.TMM denotes a

835         zero-inflated negative binomial model with TMM normalization and no transformation, and

836         so on and so forth.


837    Performance evaluation

838    Several performance metrics were considered for evaluation, all derived from some combination

839    of the elements from the confusion matrix: false positives (FPs), true positives (TPs), true

840    negatives (TNs), and false negatives (FNs). These include measures related to (i) statistical

841    power, (ii) false discovery, and (iii) software implementation and scope, all as averages over 100

842    simulation runs (**Supplementary Fig. S1B**). Several measures were considered for statistical

843    power - Sensitivity, Area Under the Curve (AUC), and scaled partial AUC (spAUC). The AUC was

844    calculated as the area under the ROC curve, obtained by plotting the sensitivity versus 1-

845    specificity for the varying p-value threshold. spAUC was calculated as the partial area over the

846  high specificity range (0, 0.20), rescaled to mimic the interpretation of AUC (i.e. 0.5 for a random

847  guess and 1 for a perfect classifier using p-values to discriminate between spiked and non-spiked

848  features). The R package *ROCR*[57] was used to calculate both these AUC measures. We also

849  considered Matthew's correlation coefficient as well as F1 scores as alternate accuracy measures

850  of performance.

851  Among false discovery metrics, maximum and average of several commonly used metrics

852  including False Discovery Rate (FDR) and False Positive Rate (FPR) were considered. When no

853  features were declared significant (i.e. TP = FP = 0), the false discovery rate (FDR) was set to 0.

854  Notably, Weiss et al.[8] misreported false positive rate as FDR, as evident from the supplemental

855  R code of that paper (Additional files 9 and 10 of Weiss et al[8]). In order to avoid any ambiguity,

856  we provide the analytical expressions of the above-mentioned measures (except AUC and

857  spAUC) as follows:

858
$$\text{FDR (1} - \text{ Precision)} = \frac{\text{FP}}{\text{FP} + \text{TP}}$$

859
$$\text{FPR (1} - \text{ Specificity)} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

860
$$\text{Sensitivity (Power or Recall)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

861
$$\text{F1 score} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

862
$$\text{Matthew's correlation coefficient (MCC)} = \frac{\text{TP} * \text{TN} - \text{FP} * \text{FN}}{\sqrt{(\text{TP} + \text{FP}) * (\text{TP} + \text{FN}) * (\text{TN} + \text{FP}) * (\text{TN} + \text{FN})}}$$

863  Following Hawinkel et al.[7], an alternative measure based on the p-value distribution under the

864  null, 'Departure from Uniformity', was also considered. Briefly, to quantify the departures from

865  uniformity into liberal (or conservative) direction, twice the mean distance between the diagonal

866  line and the points in the QQ plot below (or above) the diagonal was computed. We called these

37

867    measures 'Liberal Area' and 'Conservative Area', respectively. Both calculated areas are

868    averages over all features, and they both range from 0 to 1. A combined metric called 'Total Area'

869    that defines departure in either direction (defined as Total Area = Liberal Area + Conservative

870    Area) was also computed.

871    Finally, we calculated computational time and convergence aspects of different methods based

872    on their available implementation. Following Soneson and Robinson[58], we record the actual time

873    required to run each method using a single core and normalize all times for a given data set

874    instance so that the maximal value across all methods is 1 (as reported in **Fig. 1C**). Thus, a

875    'relative' computational time of 1 for a given method and a given data set instance means that this

876    method was the slowest one for that particular instance, and a value of, for example, 0.1 means

877    that the time requirement was 10% of that for the slowest method. Similarly, we estimated the

878    'relative' convergence failure rates for each method, as before, with the worst method as a

879    reference.

880    **Analysis of the iHMP (HMP2) IBDMDB multi-omics dataset**

881    <u>Study design, data, and quality control</u>

882    Data were obtained from the Integrative Human Microbiome Project (HMP2 or iHMP), which is

883    described in detail in Lloyd-Price et al.[36] and available through the Inflammatory Bowel Disease

884    Multi-omics Database (IBDMDB, http://ibdmdb.org). Briefly, subjects included in this cohort were

885    recruited from five academic medical centers across the US: three pediatric sub-cohorts including

886    Cincinnati Children's Hospital, Massachusetts General Hospital (MGH) Pediatrics, and Emory

887    University Hospital, and two adult sub-cohorts including MGH and Cedars-Sinai Medical Center.

888    132 subjects were followed for one year each to generate integrated longitudinal molecular

889    profiles of host and microbial activity during disease (up to 24 time points each; in total 2,965

890     stool, biopsy, and blood specimens). Self-collected stool samples were transported in ethanol

891     fixative before storage at -80 C until DNA extraction.

892     Multiple measurement types were generated from many individual stool specimens, including 305

893     samples that contain all stool-derived measurements and 791 metagenome-metatranscriptome

894     pairs. Metagenomic data generation and processing were performed at the Broad Institute. After

895     standard sequence- and sample-level quality control as described in Lloyd-Price et al.[36], species-

896     level taxonomic abundances were inferred for all samples using MetaPhlAn 2[39] and functional

897     profiling was performed by using HUMAnN 2[37]. The resulting data types including metagenome-

898     based taxonomic abundances and pathway abundance profiles for both metagenomics and

899     metatranscriptomics (summarized as structured pathways from MetaCyc[59]) were used as inputs

900     for MaAsLin 2 analysis.

901     <u>Significance testing with shuffled data</u>

902     In order to quantify whether MaAsLin 2 and other multivariable association methods identified

903     more significant associations than expected by chance (i.e. when all the shared signal between

904     features and metadata are broken), we repeatedly shuffled the metadata sample labels, applied

905     multivariable association methods on the randomized data to link features to metadata, and

906     compared the number of statistically significant associations obtained with these randomized data

907     to the number of statistically significant associations obtained with the original data based on the

908     unadjusted p-values. For a comprehensive comparison of both count and noncount models in this

909     experiment, prior to data shuffling, we multiplied the species-level taxonomic abundances by 5%

910     of the filtered read counts as a "proxy" for the underlying raw sequencing count data. The

911     procedure was repeated 1,000 times to estimate the null distribution of the detection performance

912     in both baseline and longitudinal models (with the exception of Compound Poisson mixed effects

913     model which was repeated 100 times to save computation time). While the baseline model

914     included all time-invariant covariates (age, antibiotic use, IBD diagnosis), the longitudinal model

915 also included subjects as random effects with an additional time-variant fixed effect i.e. IBD

916 dysbiosis state, as stated below.

917 <u>Statistical analysis of species, DNA pathways, and RNA pathways</u>

918 For both species and DNA pathways (whole-community and species-stratified), we regressed the

919 log-transformed relative abundances (with half the minimum relative abundance as pseudo count,

920 the default in MaAsLin 2) using the following per-feature linear mixed-effects model:

921 *feature ~ (intercept) + diagnosis + diagnosis/dysbiosis + antibiotic use + consent age + (1 |*

922 *recruitment site) + (1 | subject).*

923 Additionally, we modeled the log-transformed relative abundances of the whole-community and

924 species-stratified RNA pathways (with half the minimum relative abundance per feature as

925 pseudo count) using the similar linear mixed-effects model as before, while additionally adjusting

926 the corresponding DNA pathways abundance as a continuous covariate to filter out the influence

927 from gene copies:

928 *RNA feature ~ (intercept) + diagnosis + diagnosis/dysbiosis + antibiotic use + consent age +*

929 *DNA feature + (1 | recruitment site) + (1 | subject)*

930 That is, in each per-feature multivariable model, recruitment sites and subjects were included as

931 random effects to account for the correlations in the repeated measures (denoted by (1 |

932 recruitment site) and (1 | subject) respectively) and the transformed abundances of each feature

933 was modeled as a function of diagnosis (a categorical variable with non-IBD as the reference

934 group) and dysbiosis state as a nested binary variable (with non-dysbiotic as reference) within

935 each IBD phenotype (UC, CD, and non-IBD), while adjusting for consent age as a continuous

936 covariate, and antibiotics as a binary covariate. Nominal p-values were adjusted for multiple

937 hypothesis testing with a target false discovery rate of 0.25 with this FDR chosen to match the

938 original study.

## Data Availability

940  The iHMP dataset is publicly available at the IBDMDB website (https://ibdmdb.org) and the HMP

941  DACC web portal (https://www.hmpdacc.org/ihmp/). The processed HMP2 datasets analyzed in

942  this manuscript are also available as **Supplementary Datasets S1-S6**.

## Implementation and Software Availability

946  The implementation of MaAsLin 2 is publicly available with source code, documentation, tutorial

947  data, and as an R/Bioconductor package at http://huttenhower.sph.harvard.edu/maaslin2. The

948  software packages used in this work are free and open source, including bioBakery[60] methods

949  available via http://huttenhower.sph.harvard.edu/biobakery as source code, cloud-compatible

950  images, and installable packages. Analysis scripts using these packages to generate figures and

951  results from this manuscript (and associated usage notes) are available from

952  https://github.com/biobakery/maaslin2_benchmark. The following R packages were used to

953  generate the manuscript figures: ComplexHeatmap[61], ggalluvial[62], ggplot2[63], UpSetR[64], and

954  cowplot[65].

## Acknowledgements

## Supplementary Materials

962  **Supplementary Figure S1: Details of simulation parameters, evaluation metrics, and**

963  **benchmarking methods. A.** Four broad metadata designs commonly encountered in

964  microbiome epidemiology for varying sample size, effect size, and feature dimensions are

965  considered: UVA (Single continuous metadata), UVB (Single binary metadata), MVA (Multiple

966   independent metadata), and MVB (Multiple correlated metadata). For each of this broad metadata

967   design, both cross-sectional and longitudinal cases are evaluated (**Methods**). **B.** Three aspects

968   of performance are considered: (i) false discovery, (ii) sensitivity, and (iii) scope and

969   computational efficiency of the associated software, each comprising multiple evaluation metrics

970   (**Methods**). **C.** A combination of statistical models, normalization, and transformation schemes

971   are employed to the synthetic datasets for a variety of association methods, leading up to 84

972   combinations of normalization/transformation, zero-inflation, and regression models.

973

974   **Supplementary Figure S2: Full summary of detection performance for varying effect size,**

975   **sample size, and feature dimensions in the simple case of univariate binary metadata**

976   **without repeated measures.** Both sensitivity and false discovery rates (FDR) are shown for the

977   best-performing methods from each class of methods (as measured by average F1 score). Values

978   are averages over 100 iterations for each parameter combination. The x-axis (effect size) within

979   each panel represents the linear effect size parameter; a higher effect size represents a stronger

980   association. For visualization purposes, the best-performing methods from each class of models

981   (as measured by average F1 score) are shown. Red line parallel to the x-axis is the target

982   threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

983   across all simulation parameters in this setting. All methods were parallelized using custom bash

984   scripts in a high-performance computing environment and methods unable to process specific

985   simulation configurations due to high computational overhead or slow convergence were omitted

986   for those cases.

987

988   **Supplementary Figures S3: Meta-summary of detection performance in the simple case of**

989   **univariate binary metadata without repeated measures**. Detection performance measures

990   (Sensitivity, FPR, FDR) for all methods are provided. Values are averages over all parameter

991   combinations each summarized over 100 iterations. Red line parallel to the x-axis is the target

992   threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

993   across all simulation parameters in this setting.

994

995   **Supplementary Figures S4: Meta-summary of p-value calibration performance in the**

996   **simple case of univariate binary metadata without repeated measures**. P-value calibration

997   measures as measured by 'departure from uniformity' (Liberal Area, Conservative Area, Total

998   Area) for all methods are provided. Values are averages over all parameter combinations. Values

999   are averages over all parameter combinations each summarized over 100 iterations. Redd line

1000  parallel to the x-axis is the target threshold for FDR in multiple testing. Methods are sorted by

1001  increasing order of average F1 score across all simulation parameters in this setting.

1002

1003  **Supplementary Figure S5: Full summary of detection performance for varying effect size,**

1004  **sample size, and feature dimensions in the simple case of univariate continuous metadata**

1005  **without repeated measures.** Both sensitivity and false discovery rates (FDR) are shown for the

1006  best-performing methods from each class of methods (as measured by average F1 score). Values

1007  are averages over 100 iterations for each parameter combination. The x-axis (effect size) within

1008  each panel represents the linear effect size parameter; a higher effect size represents a stronger

1009  association. For visualization purposes, the best-performing methods from each class of models

1010  (as measured by average F1 score) are shown. Red line parallel to the x-axis is the target

1011  threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

1012  across all simulation parameters in this setting. All methods were parallelized using custom bash

1013  scripts in a high-performance computing environment and methods unable to process specific

1014  simulation configurations due to high computational overhead or slow convergence were omitted

1015  for those cases.

1016

1017 **Supplementary Figures S6: Meta-summary of detection performance in the presence of**

1018 **multiple independent covariates without repeated measures**. Detection performance

1019 measures (F1 score, Matthew's correlation coefficient, FDR) for all methods are provided. Values

1020 are averages over all parameter combinations each summarized over 100 iterations. Red line

1021 parallel to the x-axis is the target threshold for FDR in multiple testing. Methods are sorted by

1022 increasing order of average F1 score across all simulation parameters in this setting.

1023

1024 **Supplementary Figure S7: Full summary of detection performance for varying effect size,**

1025 **sample size, and feature dimensions in the presence of multiple independent covariates**

1026 **without repeated measures.** Both sensitivity and false discovery rates (FDR) are shown for the

1027 best-performing methods from each class of methods (as measured by average F1 score). Values

1028 are averages over 100 iterations for each parameter combination. The x-axis (effect size) within

1029 each panel represents the linear effect size parameter; a higher effect size represents a stronger

1030 association. For visualization purposes, the best-performing methods from each class of models

1031 (as measured by average F1 score) are shown. Red line parallel to the x-axis is the target

1032 threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

1033 across all simulation parameters in this setting. All methods were parallelized using custom bash

1034 scripts in a high-performance computing environment and methods unable to process specific

1035 simulation configurations due to high computational overhead or slow convergence were omitted

1036 for those cases.

1037

1038 **Supplementary Figures S8: Meta-summary of detection performance in the presence of**

1039 **repeated measures and univariate binary metadata.** Detection performance measures

1040 (Sensitivity, FPR, FDR) for all methods are provided. Values are averages over all parameter

1041 combinations each summarized over 100 iterations. Red line parallel to the x-axis is the target

1042    threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

1043    across all simulation parameters in this setting.

1044

1045    **Supplementary Figures S9: Meta-summary of detection performance in the presence of**

1046    **repeated measures and multiple independent covariates.** Detection performance measures

1047    (Sensitivity, FPR, FDR) for all methods are provided. Values are averages over all parameter

1048    combinations each summarized over 100 iterations. Red line parallel to the x-axis is the target

1049    threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

1050    across all simulation parameters in this setting.

1051

1052    **Supplementary Figures S10: Meta-summary of detection performance in the presence of**

1053    **repeated measures and univariate continuous metadata.** Detection performance measures

1054    (Sensitivity, FPR, FDR) for all methods are provided. Values are averages over all parameter

1055    combinations each summarized over 100 iterations. Red line parallel to the x-axis is the target

1056    threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

1057    across all simulation parameters in this setting.

1058

1059    **Supplementary Figures S11: Meta-summary of detection performance in the presence of**

1060    **repeated measures and multiple correlated covariates.** Detection performance measures

1061    (Sensitivity, FPR, FDR) for all methods are provided. Values are averages over all parameter

1062    combinations each summarized over 100 iterations. Red line parallel to the x-axis is the target

1063    threshold for FDR in multiple testing. Methods are sorted by increasing order of average F1 score

1064    across all simulation parameters in this setting.

1065

1066    **Supplementary Figure S12**. **Runtime of association methods**. CPU time (in minutes) is shown

1067    for all models faceted by feature dimension (100, 200, 500) and colored by metadata design (i.e.

1068    univariate and multivariable) in both cross-sectional (top) and longitudinal (bottom) settings.

1069    Values are averages over 100 iterations for each parameter combination. All methods were

1070    parallelized using custom bash scripts in a high-performance computing environment and

1071    methods unable to process specific simulation configurations due to high computational overhead

1072    or slow convergence were omitted for those cases.

1073

1074    **Supplementary Figure S13**. **Performance of multivariable association methods on negative**

1075    **training data across a range of significance levels**. MaAsLin 2's default linear model produced

1076    a consistently lower proportion of significant associations in negative training data (or repeatedly

1077    shuffled training set) (averaged over 1,000 permutations) than the positive training (unshuffled)

1078    counterpart in both baseline and longitudinal models (**Methods**).

1079

1080    **Supplementary Figure S14**: **Statistically significant overlap of detected features by**

1081    **MaAsLin 2 and those found in the original study**. Contingency tables describing the

1082    intersection of detected features between MaAsLin 2 and Lloyd-Price et al.[36] for various data

1083    modalities in the IBDMDB dataset are shown.

1084

1085    **Supplementary Figure S15**: **Overlap of detected taxonomic features by various MaAsLin**

1086    **models**. Upset plot describing the intersection of detected taxonomic features between various

1087    MaAsLin 2 models in the IBDMDB dataset reveal significant overlap across methods. A similar

1088    pattern was observed for functional profiles.

1089

1090    **Supplementary Datasets S1-S6: MaAsLin 2 associations between HMP2 multi-omics**

1091    **features and covariates**. List of statistically significant associations (FDR<0.25) between IBD

1092    disease phenotype (with non-IBD as reference), IBD dysbiosis state (with non-dysbiotic as

1093    reference), age, and antibiotic use with multiple data modalities (**S1**: species, **S2**: unstratified DNA

1094    pathways, **S3:** unstratified RNA pathways, **S4:** pathway RNA/DNA ratios, **S5**: stratified DNA

1095    pathways, **S6**: stratified RNA pathways) using a multivariable linear mixed effects model

1096    (**Methods**). Features are sorted by minimum FDR-adjusted p-values. For each feature, coefficient

1097    estimates and test statistics and the associated two-tailed p-values are also reported. For each

1098    data modality, input features and metadata are also provided.

1099

1100    **References**

1101

1102    1    Lynch, S. V. & Pedersen, O. The Human Intestinal Microbiome in Health and Disease. *N*
1103         *Engl J Med* **375**, 2369-2379, doi:10.1056/NEJMra1600266 (2016).
1104    2    Shreiner, A. B., Kao, J. Y. & Young, V. B. The gut microbiome in health and in disease.
1105         *Curr Opin Gastroenterol* **31**, 69-75, doi:10.1097/mog.0000000000000139 (2015).
1106    3    Franzosa, E. A. *et al.* Sequencing and beyond: integrating molecular 'omics' for microbial
1107         community profiling. *Nat Rev Microbiol* **13**, 360-372, doi:10.1038/nrmicro3451 (2015).
1108    4    Hasin, Y., Seldin, M. & Lusis, A. Multi-omics approaches to disease. *Genome Biol* **18**, 83,
1109         doi:10.1186/s13059-017-1215-1 (2017).
1110    5    Lloyd-Price, J. *et al.* Strains, functions and dynamics in the expanded Human Microbiome
1111         Project. *Nature* **550**, 61-66, doi:10.1038/nature23889 (2017).
1112    6    iHMPConsortium. The Integrative Human Microbiome Project. *Nature* **569**, 641-648,
1113         doi:10.1038/s41586-019-1238-8 (2019).
1114    7    Hawinkel, S., Mattiello, F., Bijnens, L. & Thas, O. A broken promise: microbiome
1115         differential abundance methods do not control the false discovery rate. *Brief Bioinform* **20**,
1116         210-221, doi:10.1093/bib/bbx104 (2019).
1117    8    Weiss, S. *et al.* Normalization and microbial differential abundance strategies depend
1118         upon data characteristics. *Microbiome* **5**, 27, doi:10.1186/s40168-017-0237-y (2017).
1119    9    Thorsen, J. *et al.* Large-scale benchmarking reveals false discoveries and count
1120         transformation sensitivity in 16S rRNA gene amplicon data analysis methods used in
1121         microbiome studies. *Microbiome* **4**, 62, doi:10.1186/s40168-016-0208-8 (2016).
1122    10    McMurdie, P. J. & Holmes, S. Waste not, want not: why rarefying microbiome data is
1123         inadmissible. *PLoS Comput Biol* **10**, e1003531, doi:10.1371/journal.pcbi.1003531 (2014).
1124    11    Mallick, H. *et al.* Experimental design and quantitative analysis of microbial community
1125         multiomics. *Genome Biol* **18**, 228, doi:10.1186/s13059-017-1359-z (2017).
1126    12    Jonsson, V., Österlund, T., Nerman, O. & Kristiansson, E. Statistical evaluation of methods
1127         for identification of differentially abundant genes in comparative metagenomics. *BMC*
1128         *Genomics* **17**, 78, doi:10.1186/s12864-016-2386-y (2016).
1129    13    Jonsson, V., Österlund, T., Nerman, O. & Kristiansson, E. Variability in Metagenomic
1130         Count Data and Its Influence on the Identification of Differentially Abundant Genes. *J*
1131         *Comput Biol* **24**, 311-326, doi:10.1089/cmb.2016.0180 (2017).
1132    14    Mandal, S. *et al.* Analysis of composition of microbiomes: a novel method for studying
1133         microbial composition. *Microb Ecol Health Dis* **26**, 27663, doi:10.3402/mehd.v26.27663
1134         (2015).
1135    15    Morgan, X. C. *et al.* Dysfunction of the intestinal microbiome in inflammatory bowel
1136         disease and treatment. *Genome Biol* **13**, R79, doi:10.1186/gb-2012-13-9-r79 (2012).

1137 16    Morgan, X. C. *et al.* Associations between host gene expression, the mucosal microbiome,
1138          and clinical outcome in the pelvic pouch of patients with inflammatory bowel disease.
1139          *Genome Biol* **16**, 67, doi:10.1186/s13059-015-0637-x (2015).
1140 17    Zhang, X. *et al.* Negative binomial mixed models for analyzing microbiome count data.
1141          *BMC Bioinformatics* **18**, 4, doi:10.1186/s12859-016-1441-7 (2017).
1142 18    Sharpton, T. *et al.* Development of Inflammatory Bowel Disease Is Linked to a Longitudinal
1143          Restructuring of the Gut Metagenome in Mice. *mSystems* **2**,
1144          doi:10.1128/mSystems.00036-17 (2017).
1145 19    Armour, C. R., Nayfach, S., Pollard, K. S. & Sharpton, T. J. A Metagenomic Meta-analysis
1146          Reveals Functional Signatures of Health and Disease in the Human Gut Microbiome.
1147          *mSystems* **4**, doi:10.1128/mSystems.00332-18 (2019).
1148 20    Xinyan, Z., Himel, M. & Nengjun, Y. Zero-inflated negative binomial regression for
1149          differential abundance testing in microbiome studies. *Journal of Bioinformatics and*
1150          *Genomics*, 1-1 (2016).
1151 21    Paulson, J. N., Stine, O. C., Bravo, H. C. & Pop, M. Differential abundance analysis for
1152          microbial marker-gene surveys. *Nat Methods* **10**, 1200-1202, doi:10.1038/nmeth.2658
1153          (2013).
1154 22    Kaul, A., Mandal, S., Davidov, O. & Peddada, S. D. Analysis of Microbiome Data in the
1155          Presence of Excess Zeros. *Front Microbiol* **8**, 2114, doi:10.3389/fmicb.2017.02114
1156          (2017).
1157 23    Peng, X., Li, G. & Liu, Z. Zero-Inflated Beta Regression for Differential Abundance
1158          Analysis with Metagenomics Data. *J Comput Biol* **23**, 102-110,
1159          doi:10.1089/cmb.2015.0157 (2016).
1160 24    Chen, E. Z. & Li, H. A two-part mixed-effects model for analyzing longitudinal microbiome
1161          compositional data. *Bioinformatics* **32**, 2611-2617, doi:10.1093/bioinformatics/btw308
1162          (2016).
1163 25    Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion
1164          for RNA-seq data with DESeq2. *Genome Biol* **15**, 550, doi:10.1186/s13059-014-0550-8
1165          (2014).
1166 26    Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for
1167          differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-
1168          140, doi:10.1093/bioinformatics/btp616 (2010).
1169 27    Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing
1170          and microarray studies. *Nucleic Acids Res* **43**, e47, doi:10.1093/nar/gkv007 (2015).
1171 28    Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model
1172          analysis tools for RNA-seq read counts. *Genome Biol* **15**, R29, doi:10.1186/gb-2014-15-
1173          2-r29 (2014).
1174 29    Smyth, G. K. Linear models and empirical bayes methods for assessing differential
1175          expression in microarray experiments. *Stat Appl Genet Mol Biol* **3**, Article3,
1176          doi:10.2202/1544-6115.1027 (2004).
1177 30    Van den Berge, K. *et al.* Observation weights unlock bulk RNA-seq tools for zero inflation
1178          and single-cell applications. *Genome Biol* **19**, 24, doi:10.1186/s13059-018-1406-4 (2018).
1179 31    McMurdie, P. J. & Holmes, S. phyloseq: an R package for reproducible interactive analysis
1180          and graphics of microbiome census data. *PLoS One* **8**, e61217,
1181          doi:10.1371/journal.pone.0061217 (2013).
1182 32    Calgaro, M., Romualdi, C., Waldron, L., Risso, D. & Vitulo, N. Assessment of statistical
1183          methods from single cell, bulk RNA-seq, and metagenomics applied to microbiome data.
1184          *Genome Biol* **21**, 191, doi:10.1186/s13059-020-02104-1 (2020).
1185 33    Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome*
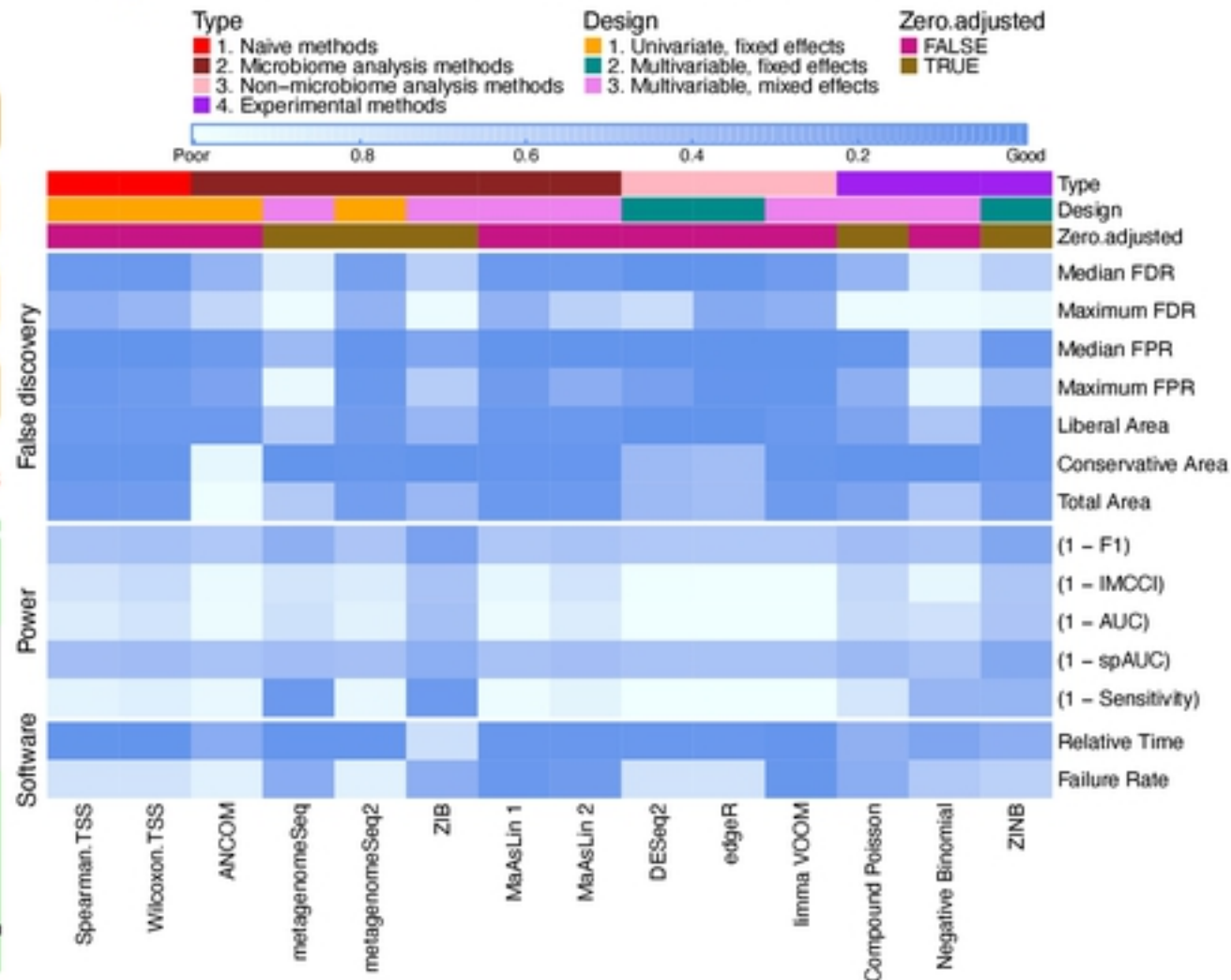1186          *Biol* **11**, R106, doi:10.1186/gb-2010-11-10-r106 (2010).

1187  34   Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression
1188       analysis of RNA-seq data. *Genome Biol* **11**, R25, doi:10.1186/gb-2010-11-3-r25 (2010).
1189  35   McKnight, D. T. *et al.* Methods for normalizing microbiome data: an ecological perspective.
1190       *Methods in Ecology and Evolution* **10**, 389-400 (2019).
1191  36   Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory bowel
1192       diseases. *Nature* **569**, 655-662, doi:10.1038/s41586-019-1237-9 (2019).
1193  37   Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and
1194       metatranscriptomes. *Nat Methods* **15**, 962-968, doi:10.1038/s41592-018-0176-y (2018).
1195  38   Caspi, R. *et al.* The MetaCyc database of metabolic pathways and enzymes. *Nucleic Acids*
1196       *Res* **46**, D633-d639, doi:10.1093/nar/gkx935 (2018).
1197  39   Truong, D. T. *et al.* MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat*
1198       *Methods* **12**, 902-903, doi:10.1038/nmeth.3589 (2015).
1199  40   Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful
1200       approach to multiple testing. *Journal of the Royal statistical society: series B*
1201       *(Methodological)* **57**, 289-300 (1995).
1202  41   Schirmer, M. *et al.* Dynamics of metatranscription in the inflammatory bowel disease gut
1203       microbiome. *Nat Microbiol* **3**, 337-346, doi:10.1038/s41564-017-0089-z (2018).
1204  42   Hawinkel, S., Rayner, J. C. W., Bijnens, L. & Thas, O. Sequence count data are poorly fit
1205       by the negative binomial distribution. *PLoS One* **15**, e0224909,
1206       doi:10.1371/journal.pone.0224909 (2020).
1207  43   Venables, W. N. & Ripley, B. D. *Modern applied statistics with S-PLUS*. (Springer Science
1208       & Business Media, 2013).
1209  44   Bates, D., Mächler, M., Bolker, B. M. & Walker, S. C. Fitting linear mixed-effects models
1210       using lme4. *Journal of Statistical Software* **67** (2015).
1211  45   Zeileis, A., Kleiber, C. & Jackman, S. Regression models for count data in R. *Journal of*
1212       *statistical software* **27**, 1-25 (2008).
1213  46   Aitchison, J. The statistical analysis of compositional data. *Journal of the Royal Statistical*
1214       *Society: Series B (Methodological)* **44**, 139-160 (1982).
1215  47   Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome Datasets
1216       Are Compositional: And This Is Not Optional. *Front Microbiol* **8**, 2224,
1217       doi:10.3389/fmicb.2017.02224 (2017).
1218  48   VanderWeele, T. J. Mediation Analysis: A Practitioner's Guide. *Annu Rev Public Health*
1219       **37**, 17-32, doi:10.1146/annurev-publhealth-032315-021402 (2016).
1220  49   Ma, S. *et al.* Population Structure Discovery in Meta-Analyzed Microbial Communities and
1221       Inflammatory Bowel Disease. *bioRxiv* (2020).
1222  50   Gibson, T. E. & Gerber, G. K. Robust and scalable models of microbiome dynamics. *arXiv*
1223       *preprint arXiv:1805.04591* (2018).
1224  51   Ren, B. S., E; Tickle, T; Huttenhower, C sparseDOSSA: Sparse Data Observations for
1225       Simulating Synthetic Abundance. R package version 1.12.0.  (2020).
1226  52   Gevers, D. *et al.* The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host*
1227       *Microbe* **15**, 382-392, doi:10.1016/j.chom.2014.02.005 (2014).
1228  53   Imhann, F. *et al.* Interplay of host genetics and gut microbiota underlying the onset and
1229       clinical presentation of inflammatory bowel disease. *Gut* **67**, 108-119, doi:10.1136/gutjnl-
1230       2016-312135 (2018).
1231  54   Cui, S., Ji, T., Li, J., Cheng, J. & Qiu, J. What if we ignore the random effects when
1232       analyzing RNA-seq data in a multifactor experiment. *Stat Appl Genet Mol Biol* **15**, 87-105,
1233       doi:10.1515/sagmb-2015-0011 (2016).
1234  55   Stasinopoulos, D. M. & Rigby, R. A. Generalized additive models for location scale and
1235       shape (GAMLSS) in R.  (2007).
1236  56   Zhang, Y. Likelihood-based and bayesian methods for tweedie compound poisson linear
1237       mixed models. *Statistics and Computing* **23**, 743-757 (2013).

1238    57    Sing, T., Sander, O., Beerenwinkel, N. & Lengauer, T. ROCR: visualizing classifier
1239        performance in R. *Bioinformatics* **21**, 3940-3941, doi:10.1093/bioinformatics/bti623
1240        (2005).
1241    58    Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell differential
1242        expression analysis. *Nat Methods* **15**, 255-261, doi:10.1038/nmeth.4612 (2018).
1243    59    Caspi, R. *et al.* The MetaCyc database of metabolic pathways and enzymes and the
1244        BioCyc collection of Pathway/Genome Databases. *Nucleic Acids Res* **42**, D459-471,
1245        doi:10.1093/nar/gkt1103 (2014).
1246    60    McIver, L. J. *et al.* bioBakery: a meta'omic analysis environment. *Bioinformatics* **34**, 1235-
1247        1237, doi:10.1093/bioinformatics/btx754 (2018).
1248    61    Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in
1249        multidimensional genomic data. *Bioinformatics* **32**, 2847-2849,
1250        doi:10.1093/bioinformatics/btw313 (2016).
1251    62    Brunson, J. C. ggalluvial: Layered Grammar for Alluvial Plots. *Journal of Open Source*
1252        *Software* **5**, 2017 (2020).
1253    63    Wickham, H. *ggplot2: elegant graphics for data analysis*. (springer, 2016).
1254    64    Lex, A., Gehlenborg, N., Strobelt, H., Vuillemot, R. & Pfister, H. UpSet: visualization of
1255        intersecting sets. *IEEE transactions on visualization and computer graphics* **20**, 1983-
1256        1992 (2014).
1257    65    Wilke, C. O., Wickham, H. & Wilke, M. C. O. Package 'cowplot'. *Streamlined Plot Theme*
1258        *and Plot Annotations for 'ggplot2* (2019).

1259

**Figure 1**

Differential abundance detection performance

Single binary metadata without repeated measures

Legend: Balanced FDR, low-powered; Inflated FDR, low-powered; Balanced FDR, well-powered; Inflated FDR, well-powered

Methods (x-axis): edgeR, limma VOOM, DESeq2, MaAsLin 1, ANCOM, Wilcoxon.TSS, MaAsLin 2, metagenomeSeq2, Compound Poisson, Negative Binomial, metagenomeSeq, ZIB, ZINB

Figure 2

## Multivariable association detection performance
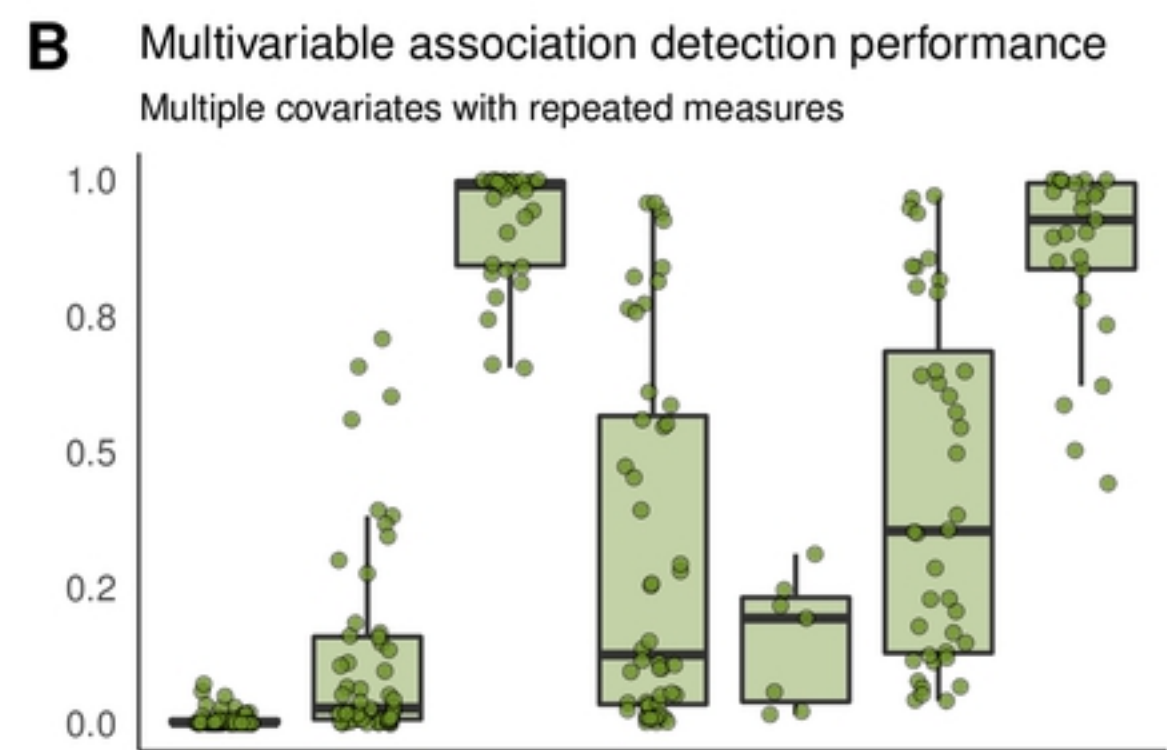### Multiple covariates without repeated measures

Figure 3

**A** Differential abundance detection performance
Single binary metadata with repeated measures

**B** Multivariable association detection performance
Multiple covariates with repeated measures

Figure 4

Figure 5