

 Open access • Posted Content • DOI:10.1101/2021.05.03.442488

## **Benchmark of data processing methods and machine learning models for gut microbiome-based diagnosis of inflammatory bowel disease** — [Source link](#)

Kubinski R, Djamén-Kepaou J, Zhanabaev T, Hernandez-Garcia A ...+7 more authors

**Institutions:** Université de Montréal, Max Planck Society, Norwich Research Park, Centre Hospitalier Universitaire Sainte-Justine

**Published on:** 04 May 2021 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Generalizability theory

Related papers:

- [Machine Learning-based Prediction Models for Diagnosis and Prognosis in Inflammatory Bowel Diseases: A Systematic Review.](#)
- [A framework for effective application of machine learning to microbiome-based classification problems](#)
- [Use of Machine Learning Approaches in Clinical Epidemiological Research of Diabetes](#)
- [Machine Learning Approaches for Inferring Liver Diseases and Detecting Blood Donors from Medical Diagnosis](#)
- [The TCGA Meta-Dataset Clinical Benchmark.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/benchmark-of-data-processing-methods-and-machine-learning-3llpy0yk9f>

1 **Benchmark of data processing methods and machine learning models for gut microbiome-**  
2 **based diagnosis of inflammatory bowel disease**

3

4 Ryszard Kubinski<sup>1\*</sup>, Jean-Yves Djamien-Kepaou<sup>1</sup>, Timur Zhanabaev<sup>1</sup>, Alex Hernandez-Garcia<sup>2</sup>,  
5 Stefan Bauer<sup>3</sup>, Falk Hildebrand<sup>4,5</sup>, Tamas Korcsmaros<sup>4,5</sup>, Sani Karam<sup>1</sup>, Prévost Jantchou<sup>6</sup>,  
6 Kamran Kafi<sup>1</sup>, Ryan D. Martin<sup>1\*</sup>

7

8 <sup>1</sup>Phyla Technologies Inc, Montreal, Canada

9 <sup>2</sup>Mila, Quebec Artificial Intelligence Institute, University of Montreal, Montreal, Canada

10 <sup>3</sup>Max Planck Institute for Intelligent Systems, Tübingen, Germany

11 <sup>4</sup>Gut Microbes & Health, Quadram Institute Bioscience, Norwich Research Park, Norwich, Norfolk,  
12 UK.

13 <sup>5</sup>Earlham Institute, Norwich Research Park, Norwich, Norfolk, UK.

14 <sup>6</sup>Centre Hospitalier Universitaire Sainte-Justine, Montréal, Canada.

15 \*Co-corresponding authors - send correspondence to [richard@phyla.ai](mailto:richard@phyla.ai) or [ryan.martin@phyla.ai](mailto:ryan.martin@phyla.ai)

16

17

18 **Keywords**

19 Inflammatory bowel disease, machine learning, gut microbiome, batch effect reduction, data  
20 normalization

21

22

23

24

25

## 26 **Abstract**

## 27 **Background**

28 Inflammatory bowel disease (IBD) patients wait months and undergo numerous invasive  
29 procedures between the initial appearance of symptoms and receiving a diagnosis. In order to  
30 reduce time until diagnosis and improve patient wellbeing, machine learning algorithms capable  
31 of diagnosing IBD from the gut microbiome's composition are currently being explored. To date,  
32 these models have had limited clinical application due to decreased performance when applied  
33 to a new cohort of patient samples. Various methods have been developed to analyze microbiome  
34 data which may improve the generalizability of machine learning IBD diagnostic tests. With an  
35 abundance of methods, there is a need to benchmark the performance and generalizability of  
36 various machine learning pipelines (from data processing to training a machine learning model)  
37 for microbiome-based IBD diagnostic tools.

## 38 **Results**

39 We collected fifteen 16S rRNA microbiome datasets (7707 samples) from North America to  
40 benchmark combinations of gut microbiome features, data normalization methods, batch effect  
41 reduction methods, and machine learning models. Pipeline generalizability to new cohorts of  
42 patients was evaluated with four binary classification metrics following leave-one dataset-out  
43 cross validation, where all samples from one study were left out of the training set and tested  
44 upon. We demonstrate that taxonomic features obtained from QIIME2 lead to better classification  
45 of samples from IBD patients than inferred functional features obtained from PICRUSt2. In  
46 addition, machine learning models that identify non-linear decision boundaries between labels are  
47 more generalizable than those that are linearly constrained. Prior to training a non-linear machine  
48 learning model on taxonomic features, it is important to apply a compositional normalization  
49 method and remove batch effects with the naive zero-centering method. Lastly, we illustrate the

50 importance of generating a curated training dataset to ensure similar performance across patient  
51 demographics.

## 52 **Conclusions**

53       These findings will help improve the generalizability of machine learning models as we move  
54 towards non-invasive diagnostic and disease management tools for patients with IBD.

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

## 72 **Introduction**

73       The human gut microbiome is a collection of microbes, viruses, and fungi residing  
74 throughout the digestive tract. The gut microbiota plays an important role in human health,  
75 influencing food digestion, the immune system, mental health, and numerous other functions  
76 (reviewed in [1]). In line with the functional role in human health, alterations in the gut microbiome  
77 have been linked to illnesses such as multiple sclerosis, type II diabetes, and inflammatory bowel  
78 disease (IBD) [2, 3]. IBD comprises two main subtypes: Crohn's disease (CD) and ulcerative  
79 colitis (UC), characterized by periodic inflammation throughout the gastrointestinal tract or  
80 localized to the colon, respectively [4]. The prevalence of IBD is increasing globally over the last  
81 several decades, from 79.5 to 84.3 per 100 000 people between 1990 and 2017, with Canada  
82 having among the highest IBD rates at 700 per 100 000 people in 2018 [5, 6]. Although the  
83 disease etiology is currently undetermined, the increasing rates of IBD have been linked to  
84 lifestyle factors, such as a Western diet [7].

85       Currently, IBD diagnosis and monitoring is primarily performed via blood tests, fecal  
86 calprotectin, and endoscopies. These methods can be costly, invasive, and display variable  
87 accuracy, all of which leads to delayed diagnosis and infrequent disease monitoring [8].  
88 Therefore, there is an unmet need for the development of further non-invasive, low-cost, and rapid  
89 methods for screening, diagnosis, and disease management for the growing number of IBD  
90 patients [9, 10]. One potential diagnostic test within these constraints involves using the gut  
91 microbiome composition to identify patients with IBD.

92       Over the past decade, several studies have compared the gut microbiome profiles of healthy  
93 individuals and those with CD or UC [2, 11–21]. Common characteristics of the gut microbiome  
94 identified in patients with IBD are the reduction in bacterial diversity and development of a  
95 dysbiotic state, referring to alterations in the structure and function of the gut microbiome  
96 compared to healthy individuals [12, 14, 19]. Principal coordinate analysis with UniFrac [16] or

97 Bray-Curtis [17] distance of the gut microbiome's composition has identified differential clustering  
98 of healthy and IBD samples. Although the dysbiotic state is commonly identified in IBD patients,  
99 it remains unknown whether the microbiome initiates IBD or is only a reflection of the patient's  
100 current health status. Larger meta-analyses have aimed to identify differentially abundant taxa  
101 between IBD patients and healthy controls in order to generate potential diagnostic biomarkers,  
102 although with limited success to date [18].

103 Due to difficulties identifying biomarkers with standard statistical methods for disease  
104 diagnosis, the field has moved to applying predictive machine learning (ML) models for  
105 classification of patient phenotypes. Several studies have demonstrated accurate classification of  
106 patients with IBD from their gut microbiome profile with ML models [2, 12, 13, 15, 18, 22–24].  
107 Common ML models employed for IBD classification include random forest (collection of decision  
108 trees for classification) [2, 15], logistic regression (binary linear classifier) [13], and neural  
109 networks (layers of differently weighted nodes contributing to a classification) [23, 24].

110 Features commonly used for IBD classification with ML models can be categorized into  
111 three groups: clinical, bacterial, and functional. Clinical features encapsulate those regarding the  
112 patient (i.e. age, sex, body mass index (BMI)) and results from other clinical tests (i.e. calprotectin,  
113 colonoscopy), which are independent of a patient's microbiome profile [25]. Taxonomy and  
114 functional features are usually determined via sequencing-based microbiome profiling, such as  
115 amplicon sequencing of the 16S rRNA gene or whole genome shotgun (WGS) sequencing of all  
116 DNA in a sample [26]. Bioinformatic tools, such as QIIME2 [27] or LotuS2 [28], provide pipelines  
117 for clustering 16S rRNA-amplicon sequences into operational taxonomic units (OTUs) which can  
118 then be compared to public databases to find taxonomy assignments [29]. WGS reads are  
119 frequently used to infer potential functions represented in the genomes of microbial community  
120 members (reviewed in [30]). Similarly, we can use known genomes in public databases to derive  
121 functional predictions in a community based solely on amplicon sequencing based taxonomy  
122 profiles, implemented in tools such as PICRUST2 [31]. Although WGS provides greater taxonomic

123 resolution and estimates of microbiome functions, 16S rRNA amplicon sequencing is currently  
124 more applicable to a diagnostic test due to its speed, affordability, and standardization of analysis  
125 tools.

126 A critical, and often under-explored, consideration for generating ML models for disease  
127 classification is their generalizability to previously unseen cohorts of patients. A ML model that  
128 underperforms when presented with data from a new patient cohort is not reliable enough to be  
129 applied in a clinical setting [32]. Despite this, models currently used in the context of microbiome  
130 data are often only trained and cross-validated with different splits of data from the same cohort.  
131 In studies where cross-validation with an unseen sample cohort is performed, the model's  
132 performance is often lower, indicative of the model overfitting to the training set [22, 23]. A  
133 proposed explanation for the reduced performance is the potential for introduction of non-  
134 biological variability to the data by wet-lab protocols and sequencing instruments during the  
135 processing of these samples, typically observed in meta-analysis of microbiome data [12].

136 In order to improve model performance on unseen data, it is necessary to apply  
137 normalization and batch effect reduction techniques prior to model training. Normalization is a  
138 critical step to remove biases to feature abundance estimates, such as the data's compositional  
139 nature, heteroskedasticity, or skewness. For example, microbiome data's compositional nature  
140 prevents the direct application of standard statistical methods as they may lead to erroneous  
141 results, and requires prior application of compositional normalization methods [33, 34]. In addition,  
142 methods have been developed to remove the technical "batch effects" commonly identified in  
143 collections of samples from different studies, such as naive zero-centering methods and the  
144 recently developed empirical Bayes' method, Meta-analysis Methods with a Uniform Pipeline for  
145 Heterogeneity (MMUPHin) [35–38]. To date, the effect of various combinations of normalization  
146 and batch effect reduction techniques on ML model generalizability remains to be benchmarked.

147 In this article, we propose a standardized approach for evaluating the performance and  
148 generalizability of data processing pipelines and ML models with microbiome data to classify

149 patients with IBD. Previous microbiome ML benchmarking studies focused on performance of  
150 various combinations of model type, normalization, and microbiome compositional features using  
151 variations of fivefold cross validation [24, 39]. Fivefold cross validation fails to assess the  
152 generalizability to new, unseen sample batches as each split potentially contains samples from  
153 all batches present in the dataset. Therefore, we implemented a leave-one-dataset-out (LODO)  
154 [40] cross-validation method to directly assess cross-batch generalizability. In this approach, the  
155 model is iteratively trained on samples of all but one dataset and then tested on the left-out  
156 dataset. Different combinations of data types, normalization methods, batch effect reduction  
157 methods, and ML models were assessed in order to establish a comprehensive performance  
158 benchmark of microbiome-based disease classification in the context of IBD.

159

160

161

162

163

164

165

166

167

168

169



## 170 **Results**

### 171 **Overview of samples and methods**

172 In order to assess the cross-batch performance of each pipeline, we implemented a LODO  
173 cross validation approach. We collected 16S rRNA gene next generation sequencing data from  
174 15 studies in North America for a total of 7707 samples, comprising 55% healthy and 45% IBD  
175 samples, of which 56% are CD and 44% are UC (**Table 1**). We completed 15 cross-validation  
176 iterations with a single dataset removed from the training set for generation of the classification  
177 model which was then used to assess model performance (**Figure 1**).

178 We evaluated the ability to classify samples from patients with IBD or non-IBD controls  
179 using different combinations of three taxonomic feature sets or six functional feature sets, eight  
180 normalization methods, four batch effect reduction methods, and nine machine learning models  
181 (**Figure 1**). The binary classification performance of each combination of feature set,  
182 normalization, batch effect reduction, and machine learning model was assessed with four  
183 classification metrics: F1 score, Matthews Correlation Coefficient (MCC), binary accuracy, and  
184 Area Under the receiver operating characteristics Curve (ROC-AUC, abbr. AUC) [60, 61]. We  
185 assessed generalizability through two methods. First, we sorted the pipeline components of  
186 interest (e.g. types of machine learning models) by the mean and standard deviation of their  
187 performance assessed by each metric. Second, in order to determine if the performance was  
188 significantly different, we performed statistical comparison of the pipelines' metrics with a Mann-  
189 Whitney U test. Therefore, the most generalizable component was defined as the top sorted  
190 method which displayed significantly better performance than baseline or other methods.

191

### 192 **Top IBD classification was obtained using taxonomic features**

193 Taxonomic features (species, genus or OTU) are predominantly used as input for ML  
194 models, whereas it is less common to use inferred functional features from PICRUSt2 as input.

195 However, previous studies have identified lower inter-individual variation of the gut microbiome's  
196 inferred functional profile than taxonomy [62, 63], suggesting that functional features may lead to  
197 better classification performance and generalizability. We processed the 16S sequencing  
198 samples with QIIME2 and PICRUSt2 to obtain taxonomy and functional feature abundance  
199 estimates, respectively.

200 For each ML model, we assessed the performance with taxonomy and functional  
201 abundance features in combination with normalization and batch effect reduction methods.  
202 Independent sorting of four classification performance metrics indicated that the taxonomic  
203 features classified IBD samples more effectively than functional features (**Figure 2A,**  
204 **Supplemental Figure 1A**). Comparison of performance with taxonomy and functional features  
205 confirmed the significantly higher performance for classification of IBD samples with taxonomic  
206 features (**Figure 2B, Supplemental Figure 1B**). Therefore, ML models using taxonomic features  
207 from this dataset lead to better classification of IBD samples than functional features.

208 Taxonomic classification with QIIME2 consists of seven hierarchical ranks, with kingdom  
209 and species at the top and bottom, respectively. Each consecutively lower taxonomy rank  
210 provides greater resolution of the gut microbiome's composition while also increasing data  
211 sparsity, which can negatively affect an ML model's performance [64]. Previous literature  
212 comparing different taxonomy ranks for disease classification indicated that lower ranks, down to  
213 genus, improved performance [65]. We assessed whether the trend for improved classification  
214 continued with the species rank and OTUs, despite their increasing sparsity. While no significant  
215 performance difference was observed between species and genus ranks, both displayed  
216 significantly higher classification performance than OTU features (**Supplemental Figure 1C**).

217

## 218 **Non-linear models achieve greatest classification performance**

219 Machine learning classification models identify decision boundaries within the feature space  
220 to separate sample labels from one another. For some ML models (BNB, Linear SVC, LR), these

221 boundaries are linearly constrained, whereas others (RF, KNN, MLP, Radial SVC, XGBoost) can  
222 identify more complex, non-linear relationships between features and class. We assessed the  
223 generalizability of three linear and five non-linear ML models across the taxonomy and functional  
224 feature sets.

225         Independent sorting of ML models for each performance metric indicated that the non-linear  
226 models had greater classification performance (top five were non-linear models) than the linear  
227 models (**Figure 3A, Supplemental Figure 2A**). Comparison of aggregate scores further  
228 confirmed non-linear models had significantly higher F1 score, MCC, AUC, and accuracy than  
229 linear models (**Figure 3B, Supplemental Figure 2B**). In order to directly assess whether the non-  
230 linearity of a model improves classification in the context of microbiome data, we compared linear  
231 and non-linear variations of a support vector machine and logistic regression. Comparison of the  
232 two variations enables direct analysis of the impact of decision boundary constraints on  
233 performance, independent of differences in model architecture. The non-linear (radial) version of  
234 logistic regression and support vector machines (Radial) had significantly greater performance  
235 than the linear version (Linear) across all four metrics (**Figure 3C, Supplemental Figure 2C**).  
236 Lastly, we assessed which non-linear model led to the highest classification performance. Across  
237 all four metrics, the random forest and XGBoost models were significantly better than MLP, KNN,  
238 and radial SVC models (**Supplemental Figure 2D**). In conclusion, non-linear models provided  
239 more accurate IBD classification, likely due to the complex relationships between features and  
240 disease labels.

241         Other ML model architectures, such as convolutional neural networks (CNNs), are  
242 commonly used for classification problems with certain structure in the input data, such as image  
243 classification. In the context of microbiome data, the CNN MDeep adds structure to OTU features  
244 through hierarchical agglomerative clustering of the phylogeny-induced correlation between  
245 OTUs [58]. As MDeep is currently only developed for OTU features, we assessed whether this  
246 CNN architecture led to greater classification performance with OTU abundance than our MLP

247 architecture. Comparison of each performance metric across all normalization and batch effect  
248 reduction methods indicated MDeep performance was not significantly different than our MLP  
249 model (**Supplemental Figure 2E**).

250 Due to the significantly better performance of non-linear classification models and  
251 taxonomic features, our subsequent analysis of normalization and batch effect reduction methods  
252 utilized only taxonomic feature sets and non-linear models.

253

### 254 **Evaluation of normalization methods**

255 We assessed normalization methods which account for different biases commonly  
256 observed in next-generation sequencing data: compositionality, heteroskedasticity, and  
257 skewness. We selected two normalizations designed for compositional data: the isometric log  
258 ratio (ILR) and centered log ratio (CLR) [66]. We selected two normalization methods which aim  
259 to reduce the heteroskedasticity: the arcsine square root (ARS) transformation [67] of the total  
260 sum scaling (TSS) values and the variance stabilized transformation (VST) from the R package  
261 DESeq2 [68]. Next, we assessed a log transformation of the TSS values (LOG), which reduces  
262 the positive skew commonly seen in the distribution of microbiome data. Lastly, we assessed  
263 normalization by TSS alone to remove differences in sequencing depth between samples or no  
264 normalization (NOT).

265 Independent sorting of each performance metric consistently identified the compositional  
266 normalization methods (CLR and ILR) as the most generalizable across non-linear models,  
267 followed by the variance/distribution modifiers (ARS, LOG, VST), and TSS as the consistently  
268 lowest performing normalization (**Figure 4, Supplemental Figure 3A**). Furthermore, the  
269 compositional methods led to significantly better performance than the other normalization types  
270 across all four metrics (**Supplemental Figure 3B**), whereas the variance/distribution modifiers  
271 and scaling method were only significantly better than no normalization. These results indicate

272 the importance of normalization methods which account for the compositional properties of  
273 microbiome data prior to model training.

274

## 275 **Evaluation of batch effect reduction methods**

276 A common issue with combining next-generation sequencing datasets for meta-analyses is  
277 the systematic differences between datasets due to differences in technical protocols. These  
278 differences add non-biological variation to the samples, decreasing the ability to ascertain  
279 biological signals [69]. Various approaches have been proposed to remove technical artifacts  
280 from dataset collections, of which we selected two relevant to microbiome data [37, 38]. First,  
281 zero-centering methods aim to reduce batch effects by centering the mean of each feature within  
282 a batch to zero. Second, Meta-analysis Methods with a Uniform Pipeline for Heterogeneity in  
283 microbiome studies (MMUPHin) [39] (microbiome specific empirical Bayes' methods), estimate  
284 and remove batch-specific parameters for each feature. Two variations of MMUPHin were  
285 implemented to simulate the scenario of obtaining a new dataset when implemented for a  
286 diagnostic test. The first (#1) applied MMUPHin to the training and test sets separately, whereas  
287 the second (#2) only applied MMUPHin to the training set (see Methods for detailed description).

288 The different batch effect reduction methods were sorted by the mean and standard  
289 deviation of their performance across taxonomic features, non-linear models, and all  
290 normalization methods. Our sorting method indicated that zero-centering was the most  
291 generalizable approach across the non-linear models. Whereas MMUPHin #1 and #2 were less  
292 generalizable than no batch reduction, with MMUPHin #2 the least generalizable (**Figure 5,**  
293 **Supplemental Figure 4A**). Additionally, zero-centering led to significantly higher binary accuracy  
294 and MCC compared to all other methods, whereas F1 score and AUC were higher only when  
295 compared to no batch effect reduction and MMUPHin #2. MMUPHin #1 had significantly better  
296 performance than MMUPHin #2, with no difference in performance observed compared to no  
297 batch effect reduction (**Supplemental Figure 4B**).

298

## 299 **Evaluation of model performance on sample and patient subgroups**

300 The samples used to assess the performance of different combinations of normalizations,  
301 batch effect reduction, and ML models were drawn from across sample collection methods (i.e.  
302 stool and biopsy) and patient demographics (i.e. paediatric and adult samples). While we did not  
303 set inclusion criteria for samples based on these differences, previous research has demonstrated  
304 distinct differences in microbiome composition between sample types and demographic groups  
305 [18, 70, 71]. For example, principal coordinate analysis (PCoA) with weighted UniFrac distance  
306 [72] and principal component analysis (PCA) of CLR-transformed taxonomic features indicated  
307 paired biopsy and stool samples from the same individual cluster separately [73].

308 We compared the model performance for the sample and patient demographics for which  
309 we were able to acquire sufficient metadata and have been associated with microbiome  
310 alterations: sample type (biopsy vs. stool), IBD subtype (CD vs. UC), sex (Female vs. Male), BMI  
311 (BMI < 30 vs. BMI > 30), and age (Adult vs. Pediatric). To assess the performance within each  
312 demographic, we included the predictions from taxonomic features (species, genus, OTU) with a  
313 compositional normalization method, zero-centering batch effect reduction, and a non-linear ML  
314 model. Our analysis focused on the MCC performance metric as it is more robust to imbalanced  
315 label distribution [61], which occurred when the samples were grouped by the five metadata  
316 categories mentioned. A logistic regression function was used to assess changes in performance  
317 corresponding to each demographic while controlling for the other metadata (**Table 2**).

318 The models displayed reduced performance for biopsy samples compared to stool samples,  
319 increased performance for samples from adult patients compared to paediatric patients, and  
320 decreased performance of samples from patients with BMI less than 30 compared to patients with  
321 BMI greater than 30. On the other hand, there was no difference in classification performance for  
322 females compared to males or for samples from patients with CD compared to patients with UC  
323 (**Table 2**). Similar results were reproduced with F1 scores, AUC, and accuracy (**Supplemental**

324 **Table 1**). The metadata groups with different performance between the two categories coincided  
325 with those that are not equally represented in our dataset, highlighting the importance of  
326 accounting for different demographic groups in a microbiome based diagnostic test.

327

### 328 **Evaluation of top performing pipeline combinations for IBD classification**

329 Our analysis identified the features, ML models, normalization methods, and batch effect  
330 reduction methods which led to the most generalizable performance. In order to determine the  
331 best overall combination of features, data processing, and ML model we assessed the top three  
332 performing models (**Table 3**). The top three models consisted of the most generalizable individual  
333 components: taxonomic features (genus), non-linear model (XGBoost or RF), compositional  
334 normalization (ILR or CLR), and zero-centering to remove batch effects. Therefore, the  
335 combination of the most generalizable methods led to the best classification performance.

336

### 337 **Identification of important features for classification with a XGBoost model**

338 In addition to predicting disease diagnoses, machine learning models can be used to identify  
339 biomarkers for disease by identifying features important for disease classification. We  
340 characterized the feature importance from the second-best overall data processing and ML model  
341 pipeline (**Table 3**). We did not analyze the feature importance of the best-performing model  
342 because the ILR normalized values no longer correspond to the starting features thereby  
343 preventing interpretation of feature importance. For an XGBoost model, the importance  
344 corresponds to a feature's contribution to the model's decision during training, referred to as the  
345 gain value [57]. We extracted the features' gain values from each of the 15 LODO iterations,  
346 sorted by the mean of all iterations, and plotted the top fifteen features (**Figure 6**). In addition, we  
347 determined the change in abundance for each taxonomy to assess whether our dataset aligned  
348 with previous findings on changes of the microbiome in IBD.

349           Amongst the top features are many taxa in the short chain fatty acid (SCFA) producing  
350 Clostridium XIVa/IV clusters, including bacteria from the *Eubacterium*, *Coprococcus*,  
351 *Lachnospira*, and *Ruminiclostridium* genera (**Figure 6A**). Aligning with previous studies, these  
352 bacteria were decreased, with the exception of *Coprococcus 3*, in IBD samples vs control samples  
353 in our dataset (**Figure 6B**) [2, 74]. *Fusobacterium* and *Veillonellaceae* genera, commonly  
354 increased in the gut microbiome of IBD patients, were also top contributors to the XGBoost  
355 classifier (**Figure 6A/B**) [2, 75]. In addition, the *Prevotellaceae* genus was the second most  
356 important feature, with the decreased abundance in IBD samples agreeing with previous studies  
357 showing decreased abundance in the gut microbiome of patients with CD and UC (**Figure 6A/B**)  
358 [76]. XGBoost classifiers have the best potential for use as a diagnostic test due to their  
359 performance as well as their interpretability and utility in identifying disease biomarkers.

360

361

362

363

364

365

366

367

368

369

370



## 371 **Discussion**

372 We assessed how different feature sets, ML models, normalization methods, and batch  
373 effect reduction methods affect predictive performance across patient cohorts in a LODO cross  
374 validation approach. The limited applicability of a PCR-based diagnostic test with a handful of  
375 microbiomes for IBD diagnosis [77] has led the field to explore the use of ML models for disease  
376 diagnosis. Our benchmark provides practical suggestions for ways to improve the performance  
377 of an IBD diagnostic test using the gut microbiome composition. First, genus abundance  
378 estimates from 16S rRNA sequencing need to be normalized by a compositional normalization  
379 method, with CLR normalization being the most appropriate as it allows for each features  
380 importance to the ML models decision to be assessed. Second, zero-centering batch effect  
381 reduction should be applied to each batch of samples collected, sequenced, and processed  
382 together to reduce systematic batch differences. Following normalization and batch effect  
383 reduction, an XGBoost or random forest classification model should be trained and optimal  
384 hyperparameters determined for implementation as a diagnostic test. With respect to the training  
385 dataset, it is important to account for patient demographics or technical differences between  
386 samples that have been associated with gut microbiome alterations. We suggest several options  
387 for optimal performance: (1) ensure balanced representation in the training dataset, (2) include  
388 the metadata labels as a feature for the model, or (3) deploy diagnostic ML models built  
389 specifically for one demographic group. In addition, the LODO cross-validation methodology is an  
390 important tool for the selection of new data preprocessing and model building methods.

391 Previous studies have demonstrated greater consistency of functional feature abundances  
392 than taxonomic feature abundance in both healthy individuals [78–80] and those with IBD [63,  
393 81]. In fact, some studies were unable to identify a single bacterium present in every IBD patient  
394 from their cohort [62]. The reduced variation and sparsity of functional features led us to  
395 hypothesize that functional abundance profiles would lead to better classification of IBD samples.

396 However, through our LODO cross validation, we found that classification performance with  
397 functional features was significantly worse than with taxonomic features (**Figure 2B**,  
398 **Supplemental Figure 1B**). We postulate the reason for the reduced classification performance  
399 with functional profiles is due to the limited recapitulation of functional profiles with PICRUSt2 [31,  
400 82] and the inability of 16S rRNA sequencing to identify strain-level functional differences of the  
401 present bacteria [83]. To overcome these limitations in future studies, measurement of the  
402 microbiome's gene content by WGS, transcriptomes by RNA-seq, or metabolites by  
403 metabolomics need to be explored. In fact, functional profiles from whole genome sequencing led  
404 to better predictions of patients with IBD who achieved remission with vedolizumab than taxonomy  
405 abundance [23]. While whole genome sequencing may improve disease classification, its much  
406 higher cost than 16S rRNA sequencing substantially hinders the technology's adoption as a  
407 diagnostic test.

408 A major hurdle in the implementation of sequencing based diagnostic tests in the clinic is  
409 the observed systematic differences between sample preparations. In a previous study, removal  
410 of these batch effects with an empirical Bayes' or zero centering approach led to improved  
411 classification [84]. However, our work only identified improved cross-batch classification  
412 performance with zero-centering and not the empirical Bayes' method MMUPHin (**Figure 5**).  
413 Current empirical Bayes' approaches are designed and optimized for disease mechanism and  
414 biomarker discovery where the disease covariate is known and incorporated in the method. The  
415 inclusion of a disease covariate is not applicable to a diagnostic scenario though, where the  
416 diagnosis label is to be determined. The lack of improvement in classification performance with  
417 MMUPHin #1 compared to no batch reduction is potentially due to its implementation in a scenario  
418 the method was not optimized for.

419 Similar to batches of samples collected for a diagnostic test, the batches in our dataset were  
420 not balanced, with some containing only a single diagnosis class (e.g. all samples coming from  
421 IBD patients). In cases where the batch and diagnosis label are confounded, batch correction

422 methods tend to reduce the disease associated differences in the process of removing the batch  
423 differences [37]. Therefore, the more advanced removal of batch effects by MMUPHIn likely led  
424 to an over-adjustment within the unbalanced batches and removal of the disease differences.  
425 Whereas, the less sophisticated removal of batch effects with the covariate naive zero-centering  
426 approach retained sufficient biological signal between disease labels for non-linear ML models to  
427 correctly classify samples across batches. Batch correction methods that do not require input of  
428 a covariate have been developed, such as frozen surrogate variable analysis or reference  
429 principal component integration (RPCI) [85, 86], although their applicability to microbiome data  
430 has not been assessed.

431         The sparse availability of metadata for the samples led to several limitations in our analysis.  
432 First, the identification of CD and UC patients relied on the accuracy of the diagnosis coding in  
433 the public databases. However, there were no studies explicitly validating the registration of CD  
434 and UC diagnosis codes. Second, although our study demonstrated reliable results, gaps in the  
435 publicly available data prevented us from several critical analyses. For instance, we lacked  
436 information on how the patients were diagnosed in every study, the timing of sample collection in  
437 relationship to their diagnosis and disease progression, current disease activity quantification,  
438 DNA extraction and sample storage information. Furthermore, there was limited information on  
439 environmental factors such as medication usage, alcohol usage, smoking, diet, and other factors  
440 known to alter the gut microbiome which could affect our analysis [81, 87]. Of the sample  
441 information and patient demographic data we obtained, clear differences in performance of our  
442 top pipelines were observed (**Table 2**). Therefore, future studies with improved lifestyle and  
443 clinical metadata are needed to systematically address how these factors affect performance of  
444 a gut microbiome diagnostic test.

445         Other non-invasive diagnostic tests for IBD, such as fecal calprotectin, continue to have  
446 significant differences between the reports on the sensitivity and specificity for classifying IBD  
447 patients from non-IBD [88, 89]. While high performance levels have been reported, one recent

448 study identified a 78% accuracy for identifying patients with IBD using fecal calprotectin [90],  
449 which is approximately 10% lower than our best model. Furthermore, while we focused solely on  
450 IBD classification, ML models using microbiome composition have wider applicability than  
451 singular biomarkers such as calprotectin. Models using microbiome data have already been  
452 implemented to predict if a patient with IBD will respond to a medication [23], to predict a patient's  
453 postprandial glycemic response [91], and for classification of other diseases, such as Parkinson's  
454 disease [48], to name a few.

455

## 456 **Conclusion**

457 With sufficient data and validation, analysis of the fecal gut microbiome can indeed be  
458 leveraged as a multi-purpose predictive tool. Given the significant delay [92–94] and associated  
459 costs of diagnosis [95, 96], it is critical to continue exploration of approaches that increase  
460 accessibility of diagnosis and decrease the cost of testing [97] in a community health or primary  
461 care setting. XGBoost and random forest machine learning models with microbiome data have  
462 the potential to achieve these goals. Further work to gather more well-annotated data, improve  
463 performance and assess models with validation studies is required.

464

465

466

467

468

469

470

471

472

## 473 **Methods**

### 474 **Acquisition of sample data**

475 Sample FASTQ files were acquired from the European Nucleotide Archive (ENA) browser.  
476 The sample metadata was acquired from the corresponding publication's supplementary  
477 materials or the QIITA microbiome platform. Only samples collected from individuals in North  
478 America were used from each dataset. The dataset accessions and technical information  
479 regarding the samples in each dataset are available in **Supplemental Table 2**.

480 The following fifteen studies were included in our dataset:

- 481 1. The American Gut cohort is from a large, open platform which collected samples from  
482 individuals in the US to identify associations between microbiomes, the environment, and  
483 individual's phenotype [41]. We included available samples that did not contain any self-  
484 reported diseases in the metadata.
- 485 2. The CVDF study determined the effect of cardiorespiratory fitness on microbiome  
486 composition and comprises a range of fitness levels [42].
- 487 3. The GEVERSM study assessed the microbiome composition of treatment naive, newly  
488 diagnosed, paediatric patients with IBD and adult patients diagnosed with IBD for 0 to 57  
489 years [2].
- 490 4. The GEVERSC cohort consists of additional samples from paediatric and adult patients  
491 added to the GEVERSM study [2].
- 492 5. The GLS study longitudinally sampled 19 patients with CD (Crohn's disease activity index  
493 (CDAI) between 44 and 273) and 12 healthy control individuals [20].
- 494 6. The Human Microbiome Project (HMP) study longitudinal tracked paediatric and adult  
495 patients ranging from newly diagnosed to diagnosed for 39 years. Diagnosis was  
496 confirmed by colonoscopy prior to enrollment in the study along with several other  
497 inclusion criteria listed in the corresponding publication [43].

- 498 7. The MUC study collected mucosal biopsies from 44 paediatric patients with CD and 62  
499 non-IBD paediatric control patients [21].
- 500 8. PRJNA418765 was a longitudinal study of patients with CD that were refractory to anti-  
501 TNF initiating ustekinumab assessed at week 0, 4, 6 and 22. To be included, patients  
502 required at least three months Crohn's disease history and a CDAI between 220 and 450  
503 [44].
- 504 9. PRJNA436359 was a longitudinal study of new onset and treatment naive paediatric  
505 patients with UC receiving a variety of medications at week 0, 4, 12, and 52. Inclusion  
506 criteria consisted of presence of disease beyond the rectum, Paediatric Ulcerative Colitis  
507 Activity Index (PUCAI) of 10 or more, and no previous therapy [45].
- 508 10. QIITA10184 was a study comparing five different fecal collection methods and their effect  
509 on the healthy participant's microbiome composition identified with 16S rRNA gene  
510 sequencing [46].
- 511 11. QIITA10342 study assessed the microbiome composition and function of healthy  
512 individuals in two American Indian communities in the United States [47].
- 513 12. QIITA10567 samples consist of the control individuals in a study linking alterations in  
514 microbiome composition to Parkinson's disease [48].
- 515 13. The QIITA1448 study compared microbiome composition of individuals in traditional  
516 agricultural societies in Peru to those in industrialized cities in the United States [49].
- 517 14. The QIITA2202 study collected longitudinal stool samples from two healthy individuals  
518 alongside detailed lifestyle characteristics to correlate with microbiome composition [48,  
519 50].
- 520 15. The QIITA550 study collected longitudinal stool samples from two individuals to assess  
521 temporal changes in microbiome composition [51].

522

## 523 **Taxonomy classification with QIIME2**

524 Taxonomy abundance tables were generated from the FASTQ files using QIIME2 (v2020.2)  
525 [27]. Reads were trimmed to remove low quality reads (trimming parameters listed in  
526 Supplemental Table 1), chimeras removed, and sequences denoised using Dada2 [52] or Deblur  
527 (for GLS and AG only). The processed sequences were clustered into OTUs and the centroid  
528 sequences classified with a Naive Bayes classifier [53] at 99% identity using the Silva 132 99%  
529 reference database [29, 54, 55]. For classification, the corresponding 16S rRNA gene  
530 hypervariable region's sequences were extracted from the Silva 132 99% reference database  
531 with the QIIME2 plugin feature-classifier's extract-reads function using the primers from the  
532 respective study. The extracted reads and the corresponding taxonomy were used to train the  
533 Naive Bayes classifier with the QIIME2 plugin feature-classifier's fit-classifier-naive-bayes  
534 function. Taxonomic feature tables were collapsed to species (level 7) and genus (level 6)  
535 classification for further analysis.

536

## 537 **Inferring Functional Abundance with PICRUST2**

538 Functional abundance tables were generated using PICRUST2 (v2.3.0) from the OTU  
539 abundance table and representative OTU sequences from QIIME2. We generated abundance  
540 tables from the six different databases incorporated into PICRUST2: Clusters of Orthologous  
541 Groups of proteins (COG), Kyoto Encyclopedia of Genes and Genomes (KEGG) orthologs (KO),  
542 Enzyme Commission (EC), Pfam protein domain (PFAM), TIGR protein family (TIGRFAM) and  
543 MetaCyc pathways. Each database is independently curated and provides information on different  
544 aspects of the functional properties present in the microbiome.

545

## 546 **Leave-One-Dataset-Out (LODO) cross validation**

547 The generalizability of each model, normalization, and batch effect reduction method, was  
548 determined through a cross validation strategy which assessed predictive performance on

549 previously unseen batches of samples (**Supplemental Figure 5**). As there were 15 datasets, we  
550 iterated through the full dataset 15 times, generating the training set by removing all samples from  
551 a single dataset to a separate test set. The training set was used to prune features that were not  
552 present in at least 10% of samples from one dataset. Following pruning, the remaining features  
553 were selected from the test set and the samples were normalized and batch reduced with the  
554 respective methods. Lastly, the training set was balanced to have the same number of healthy  
555 and IBD samples by subsampling the label with the greater number of samples, while maintaining  
556 the proportion of samples from each collection site, disease label (UC/CD/Control), and sample  
557 type (stool/biopsy).

558 For our modified implementation of MMUPHin, the data processing was adjusted to ensure  
559 the training and test sets were batch reduced separately. For MMUPHin #1, the test dataset's  
560 samples were removed and the data from the remaining studies batch reduced with MMUPHin  
561 prior to training the model. Independently, the full dataset (with training studies and test dataset)  
562 was batch reduced and the test dataset's samples then used to assess the model's classification  
563 performance (**Supplemental Figure 5**, MMUPHin #1). For MMUPHin #2, the training studies  
564 were batched reduced with MMUPHin prior to model training and the model's classification  
565 performance then assessed on non-batch reduced samples from the test dataset (**Supplemental**  
566 **Figure 5**, MMUPHin #2). Lastly, feature abundance for some samples following MMUPHin batch  
567 effect reduction on the training set when QIITA2202 was left out and the test set when HMP was  
568 left out were all zero. Rows with all zero are not appropriate input for the compositional  
569 normalization methods, therefore we replaced the feature values for these samples with equal  
570 relative abundance prior to normalization.

571



## 572 **Feature Selection**

573           Following taxonomy classification and inference of functional abundance, features present  
574 in less than 10% of the samples within each dataset were pruned from the dataset. The feature  
575 pruning was performed on the training set only, with the features then selected from the test set.

576

## 577 **Normalization methods**

578           When possible, normalization methods were implemented using python (v3.6.12) and R  
579 (v3.6.3) packages with the methods already incorporated. For CLR and ILR normalization, zero  
580 values were first replaced with the multiplicative replacement function prior to normalization with  
581 the clr and ilr functions, respectively, from the python package SciKit-Bio (v0.5.2). CLR performs  
582 a log transformation of abundance values, which are normalized by the geometric mean of all  
583 features. ILR uses a change of coordinate space projection calculation to transform proportional  
584 data (or relative abundances) to a new space with an orthonormal basis.

585           For TSS normalization, the counts for each feature were divided by the sum of all feature  
586 counts in the sample with a custom python function. The method constrains the sample row sum  
587 to one, aiming to similarly scale all samples while maintaining biological information of microbial  
588 abundances. For ARS normalization, the TSS normalized values were transformed with sqrt  
589 function followed by the arcsin function from the python package numpy (v1.19.2). The LOG  
590 normalization was also applied to the TSS normalized values using the log function from numpy  
591 following replacement of all 0s with 1.

592           For VST normalization, we used the varianceStabilizingTransformation function in the R  
593 package DESeq2 (v1.26.0). VST aims to factor out the dependence of the variance in the mean  
594 abundance of a feature. The method numerically integrates the dispersion relation of the feature  
595 mean fitted with a spline, evaluating the transformation for each abundance in the feature. VST  
596 normalization was performed similarly to the previously described modified MMUPHin

597 implementation, with the training set normalized separately from the test set as the normalization  
598 is dependent on all samples present in the dataset.

599

## 600 **Batch effect reduction methods**

601 We explored two methods for batch effect reduction: naive zero-centering and an empirical  
602 Bayes method. The naive zero-centering batch effect reduction entails centering the mean of  
603 each feature within each batch to zero [37]. We also assessed MMUPHin, a recently developed  
604 empirical Bayes method designed specifically for zero-inflated microbial abundance data.  
605 MMUPHin estimates parameters for the additive and multiplicative batch effects, using normal  
606 and inverse gamma distributions, respectively. The estimated parameters are then used to  
607 remove the batch effects from the dataset [38, 56]. For MMUPHin, the sample type (stool/biopsy)  
608 was used as a covariate for MMUPHin #1 and the sample type and disease label (UC/CD/Control)  
609 were covariates for MMUPHin #2. We considered a batch as the whole dataset or split a dataset  
610 into multiple batches when the metadata indicated different sample preprocessing methods or  
611 samples were processed in different locations.

612

## 613 **Standard machine learning models**

614 We assessed the classification performance of standard machine learning and deep  
615 learning models. The standard models were implemented using the python package SciKit-Learn  
616 (v0.23.2). Hyperparameters were not optimized and decided prior to experimentation.

617

### 618 *Bernoulli Naive Bayes Classifier*

619 The Bernoulli Naive Bayes Classifier (BNB) model converts the feature space to binary  
620 values and then estimates parameters of a Bernoulli distribution for classification purposes. We  
621 implemented the BNB model using the default settings in SciKit-Learn.

622

623 *Random Forest*

624 Random Forest (RF) models use an ensemble of decision trees that discriminate the feature  
625 space by a sequence of threshold conditional statements. The power of the model comes from  
626 its non-linear classification capabilities and the number of trees used to label classification. We  
627 implemented the Random Forest classifier with the following modifications to the default SciKit-  
628 learn settings: `n_estimators = 500`, `max_features = sqrt`, and `class_weight = balanced`.

629

630 *K-Nearest Neighbour Classifier*

631 The K-Nearest Neighbour Classifier (KNN) classifies each sample by majority vote of the K  
632 nearest neighbours in its surrounding. We implemented the K Neighbors classifier with the  
633 following modifications to the default SciKit-learn settings: `n_neighbors = 6`, `weights = distance`,  
634 and `metric = manhattan`.

635

636 *Support Vector Machine Classifier*

637 The Support Vector Machine Classifier (SVC) identifies multivariate decision boundaries  
638 that separate class labels. We implemented two SVC variations, the first with a linear kernel,  
639 constraining the decision boundary to a linear hyperplane, using the `SGDClassifier` class from  
640 SciKit-learn with the following modifications to default settings: `loss = modified_huber`, `tol = 10e-`  
641 `5`, and `max_iter = 10000`. The second variation used the radial basis function kernel with the SVC  
642 class from SciKit-Learn, which removes the linear constraint of the decision boundary, with the  
643 following modifications to the default settings: `tol = 10e-6`, `class_weight = balanced`, and `max_iter`  
644 `= 100000`.

645

646 *Logistic Regression*

647 Logistic Regression classification estimates the probability of a certain class in a binary  
648 classification problem using a statistical fit to the logistic function. We implemented the

649 LogisticRegression class from SciKit-Learn with the following modifications to the default settings:  
650 solver = sag, class\_weight = balanced, and max\_iter = 10000. For the non-linear variation, the  
651 feature space was first transformed with the radial basis function kernel implemented with the  
652 rbf\_kernel function from SciKit-Learn prior to fitting a logistic regression model.

653

654 Gradient Boosted Trees (XGBoost)

655 Gradient boosted trees consist of a collection of sequential decision trees, where each tree  
656 learns and reduces the error of the previous tree [57]. The gradient boosted trees model was  
657 implemented with the XGBoost package's (v1.2.0) XGBoostClassifier class with the following  
658 modifications to default settings: n\_estimators = 500.

659

660 **Deep learning models**

661 The deep learning models were built with the python package Tensorflow (v2.2.0). The  
662 models were trained for up to 100 epochs with a batch size of 16 and samples shuffled. The best  
663 weights were selected using early stopping (EarlyStopping callback) by monitoring the validation  
664 loss (5% split of the training set) with a min\_delta =  $1 \times 10^{-3}$  and patience = 10.

665

666 *Multilayer Perceptron (MLP)*

667 A MLP is a neural network architecture composed of one or more layers of fully connected  
668 neurons that take as input the weights of the previous layer and output the result of an activation  
669 function to the subsequent layer. For binary classification, the final layer contains a single node  
670 that predicts the class probability. We implemented an MLP architecture with three hidden layers  
671 of 256 neurons using a rectified linear unit (ReLU) activation function followed by a Dropout layer  
672 with a dropout rate of 50%. The final layer predicted the class label with a sigmoid activation  
673 function. The model was trained using a binary cross entropy loss function and the Adam  
674 optimizer with a learning rate of 0.001.

675

## 676 *Convolutional Neural Network*

677 We implemented MDeep, a CNN architecture recently designed for microbiome data [58].  
678 CNNs require an inherent structure to present in the data, which is added to the OTU dataset by  
679 hierarchical agglomerative clustering of the phylogeny-induced correlation between OTUs. We  
680 built a phylogenetic tree with the `align_to_tree_mafft_fasttree` function in the QIIME2 phylogeny  
681 python plugin using the OTU representative sequences obtained from clustering 16S rRNA  
682 sequences with QIIME2. The phylogenetic tree was imported into R using the phyloseq package  
683 and the cophenetic distance between OTUs determined with the R package ape. The cophenetic  
684 distance was then used to calculate the phylogeny-induced correlation as described in the original  
685 study and OTUs clustered using the HAC function from the MDeep GitHub repository  
686 (<https://github.com/lichen-lab/MDeep>).

687

## 688 **Performance metrics**

689 To measure the performance of the various normalization, batch effect reduction, and model  
690 combinations we used four commonly used metrics for binary classification: F1 score, Area Under  
691 the receiver operating characteristic Curve (AUC), binary accuracy and Matthews Correlation  
692 Coefficient (MCC). Since the number of samples in each dataset ranged from 23 to 1279, we first  
693 balanced the number of samples from each dataset by up sampling each (with replacement) to  
694 100 000 samples while maintaining the confusion matrix proportions for each individual dataset.  
695 Balancing the number of samples ensured that altered performance with a single, large dataset  
696 did not control the overall score and changes in performance for small studies was still observed.  
697 The up sampled dataset was then used to calculate the respective metrics using the functions  
698 implemented in SciKit-Learn.

699

## 700 **Sample subgroup performance analysis**

701 We assessed the performance of our algorithm for five different metadata variables, each  
702 with two categorical labels. The samples were grouped by the five variables, with the two  
703 categories for each variable coded as 0 or 1. The performance metric was calculated within each  
704 grouping for classification of control samples and either UC or CD (depending on the specific  
705 grouping). For the logistic regression analysis, the metric was input as the dependent variable  
706 and the five metadata groups as the independent variables. The MCC score was scaled with the  
707 MinMaxScaler from SciKit-Learn to scale the range from 0 to 1 as required for the logistic function.

708

## 709 **Feature Importance from XGBoost Classifier**

710 In order to determine the importance of each taxonomy, we collected the features' gain  
711 value from our second-best pipeline composed of CLR normalized, zero-centered, genus  
712 abundance features with an XGBoost Classifier. The gain values were collected from the trained  
713 XGBoost classifier in each LODO iteration separately.

714

## 715 **Taxonomy Differential Abundance**

716 Differential taxonomy abundance was performed with Analysis of Compositions of  
717 Microbiomes with Bias Correction (ANCOM-BC) (v1.0.5) [59]. The fold change between control  
718 samples and IBD samples (UC and CD) was determined with a Bonferroni multiple comparison  
719 correction applied to the p-values.

720

721

722

723

## 724 **Abbreviations**

725	ANCOM-BC	Analysis of Compositions of Microbiomes with Bias Correction
726	ARS	Arcsine square root transformation
727	AUC	Area Under the receiver operating characteristics Curve
728	BMI	Body mass index
729	BNB	Bernoulli Naive Bayes Classifier
730	CD	Crohn's disease
731	CDAI	Crohn's disease activity index
732	CLR	Centered-log ratio
733	CNN	Convolutional neural network
734	COG	Clusters of Orthologous Groups of proteins
735	EC	Enzyme Commission
736	HMP	Human Microbiome Project
737	IBD	Inflammatory Bowel Disease
738	ILR	Isometric-log ratio
739	KEGG	Kyoto Encyclopedia of Genes and Genomes
740	KNN	K-Nearest Neighbour Classifier
741	KO	KEGG orthologs
742	LODO	Leave-one-dataset-out
743	LOG	Log transformation
744	LR	Logistic Regression
745	MCC	Matthews Correlation Coefficient
746	ML	Machine learning
747	MLP	Multilayer Perceptron
748	MMUPHin	Meta-analysis Methods with a Uniform Pipeline for Heterogeneity

749	NOT	No normalization
750	OTU	Operational taxonomic unit
751	PCA	Principal component analysis
752	PCoA	Principal coordinate analysis
753	PFAM	Pfam protein domain
754	PUCAI	Paediatric Ulcerative Colitis Activity Index
755	RF	Random Forest
756	RPCI	Reference principal component integration
757	SCFA	Short chain fatty acid
758	SVC	Support Vector Machine Classifier
759	TSS	Total sum scaling
760	TIGRFAM	TIGR protein family
761	UC	Ulcerative Colitis
762	VST	Variance stabilized transformation
763	XGBoost	eXtreme Gradient Boosting
764	WGS	Whole genome shotgun
765		
766		
767		
768		
769		
770		
771		



772 **Declarations**

773 **Ethics approval and consent to participate**

774 Not applicable.

775

776 **Consent for publication**

777 Not applicable.

778

779 **Availability of data and material**

780 Publicly available datasets were analyzed in this study. The raw sequencing data for the following  
781 16S rRNA datasets were downloaded from European Nucleotide Archive at the following  
782 accession numbers: American Gut (PRJEB11419), CVDF (PRJNA308319), GEVERSC  
783 (PRJEB13680), GEVERSM (PRJEB13679), GLS (PRJEB23009), MUC (PRJNA317429),  
784 PRJNA418765, PRJNA436359, QIITA10184 (PRJEB13895), QIITA10342 (PRJEB13619),  
785 QIITA10567 (PRJEB14674), QIITA1448 (PRJEB13051), QIITA2202 (PRJEB6518), QIITA550  
786 (PRJEB19825). The raw sequencing data for the HMP 16S rRNA dataset was downloaded from  
787 [ibdmdb.org](http://ibdmdb.org).

788

789 **Competing interests**

790 RK is a founder of Phyla Technologies Inc and is currently the Chief Scientific Officer. RM, JD,  
791 and TZ were employed by Phyla Technologies Inc at the time of the manuscript.

792

793 **Funding**

794 The work in this manuscript was funded by Investissement Québec Programme innovation – volet  
795 1 and Quebec Ministry of Economy and Innovation’s Entrepreneurship Assistance Program  
796 (PAEN) - component 3a. The work of FH and TK were supported by the Earlham Institute

797 (Norwich, UK) in partnership with the Quadram Institute Bioscience (Norwich, UK) and  
798 strategically supported by a UKRI BBSRC UK grant (BB/CSP17270/1). FH and TK were also  
799 supported by a BBSRC ISP grant for Gut Microbes and Health BB/R012490/1 and its constituent  
800 projects, BBS/E/F/000PR10353 and BBS/E/F/000PR10355. FH received funding from the  
801 European Research Council (ERC) under the European Union's Horizon 2020 research and  
802 innovation programme (grant agreement No. 948219)

803

#### 804 **Authors' contributions**

805 RK, JD, TZ, and RM designed the data processing pipeline, performed the experiments and  
806 analyzed the pipelines' performance. RK and RM wrote the manuscript. AHG and SB contributed  
807 to the experimental design. AHG, SB, FH, TK, SK, PJ, and KK contributed to interpretation of the  
808 results and editing and revising the manuscript. All authors reviewed, revised, and approved the  
809 final manuscript.

810

#### 811 **Acknowledgements**

812 We would like to thank Luca Cuccia, Laura Minkova, Houman Farzin, Michael Golfi, Paul Godin,  
813 and Yasmine Mouley (Phyla Technologies Inc.) for their feedback and support as the manuscript  
814 was completed. We would also like to thank Sébastien Giguère (Valence Discovery) for his  
815 guidance during our methodology development.

816

#### 817 **Author Information**

##### 818 **Phyla Technologies Inc, Montréal, Canada**

819 Ryszard Kubinski, Jean-Yves Képaou Djamen, Timur Zhanabaev, Sani Karam, Kamran Kafi,  
820 Ryan D. Martin

##### 821 **Mila (Québec Artificial Intelligence Institute), University of Montreal, Montreal, Canada**

822 Alex Hernandez-Garcia

823 **Max Planck Institute for Intelligent Systems, Tübingen, Germany**

824 Stefan Bauer

825 **Gut Microbes & Health, Quadram Institute Bioscience, Norwich Research Park, Norwich,**

826 **Norfolk, UK.**

827 Falk Hildebrand, Tamas Korcsmaros

828 **Earlham Institute, Norwich Research Park, Norwich, Norfolk, UK.**

829 Falk Hildebrand, Tamas Korcsmaros

830 **Centre Hospitalier Universitaire Sainte-Justine, Montréal, Canada.**

831 Prévost Jantchou

832

833

834

835

836

837

838

839

840

841

842

843

844

845

## 846 References

- 847 1. Mohajeri MH, Brummer RJM, Rastall RA, Weersma RK, Harmsen HJM, Faas M, et al. The role  
848 of the microbiome for human health: from basic science to clinical applications. *Eur J Nutr.*  
849 2018;57 Suppl 1:1–14.
- 850 2. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The  
851 treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe.* 2014;15:382–92.
- 852 3. Opazo MC, Ortega-Rocha EM, Coronado-Arrázola I, Bonifaz LC, Boudin H, Neunlist M, et al.  
853 Intestinal Microbiota Influences Non-intestinal Related Autoimmune Diseases. *Front Microbiol.*  
854 2018;9:432.
- 855 4. Caruso R, Lo BC, Núñez G. Host-microbiota interactions in inflammatory bowel disease. *Nat*  
856 *Rev Immunol.* 2020;20:411–26.
- 857 5. Benchimol EI, Bernstein CN, Bitton A, Murthy SK, Nguyen GC, Lee K, et al. The Impact of  
858 Inflammatory Bowel Disease in Canada 2018: A Scientific Report from the Canadian Gastro-  
859 Intestinal Epidemiology Consortium to Crohn's and Colitis Canada. *J Can Assoc Gastroenterol.*  
860 2019;2 Suppl 1:S1–5.
- 861 6. GBD 2017 Inflammatory Bowel Disease Collaborators. The global, regional, and national  
862 burden of inflammatory bowel disease in 195 countries and territories, 1990-2017: a systematic  
863 analysis for the Global Burden of Disease Study 2017. *Lancet Gastroenterol Hepatol.* 2020;5:17–  
864 30.
- 865 7. Rizzello F, Spisni E, Giovanardi E, Imbesi V, Salice M, Alvisi P, et al. Implications of the  
866 Westernized Diet in the Onset and Progression of IBD. *Nutrients.* 2019;11.  
867 doi:10.3390/nu11051033.
- 868 8. Ricciuto A, Mack DR, Huynh HQ, Jacobson K, Otley AR, deBruyn J, et al. Diagnostic Delay Is  
869 Associated With Complicated Disease and Growth Impairment in Paediatric Crohn's Disease. *J*  
870 *Crohns Colitis.* 2021;15. doi:10.1093/ecco-jcc/jjaa197.
- 871 9. Noiseux I, Veilleux S, Bitton A, Kohen R, Vachon L, Guay BW, et al. Inflammatory bowel  
872 disease patient perceptions of diagnostic and monitoring tests and procedures. *BMC*  
873 *Gastroenterol.* 2019;19:1–11.
- 874 10. Armstrong D, Barkun AN, Chen Y, Daniels S, Hollingworth R, Hunt RH, et al. Access to  
875 specialist gastroenterology care in Canada: the Practice Audit in Gastroenterology (PAGE) Wait  
876 Times Program. *Can J Gastroenterol.* 2008;22:155–60.
- 877 11. Pittayanon R, Lau JT, Leontiadis GI, Tse F, Yuan Y, Surette M, et al. Differences in Gut  
878 Microbiota in Patients With vs Without Inflammatory Bowel Diseases: A Systematic Review.  
879 *Gastroenterology.* 2020;158:930–46.e1.
- 880 12. Duvallat C, Gibbons SM, Gurry T, Irizarry RA, Alm EJ. Meta-analysis of gut microbiome  
881 studies identifies disease-specific and shared responses. *Nat Commun.* 2017;8:1784.
- 882 13. de Meij TGJ, de Groot EFJ, Peeters CFW, de Boer NKH, Kneepkens CMF, Eck A, et al.  
883 Variability of core microbiota in newly diagnosed treatment-naïve paediatric inflammatory bowel  
884 disease patients. *PLoS One.* 2018;13:e0197649.

- 885 14. Pascal V, Pozuelo M, Borrueal N, Casellas F, Campos D, Santiago A, et al. A microbial  
886 signature for Crohn's disease. *Gut*. 2017;66:813–22.
- 887 15. Tedjo DI, Smolinska A, Savelkoul PH, Masclee AA, van Schooten FJ, Pierik MJ, et al. The  
888 fecal microbiota as a biomarker for disease activity in Crohn's disease. *Sci Rep*. 2016;6:35216.
- 889 16. Halfvarson J, Brislawn CJ, Lamendella R, Vázquez-Baeza Y, Walters WA, Bramer LM, et al.  
890 Dynamics of the human gut microbiome in inflammatory bowel disease. *Nat Microbiol*.  
891 2017;2:17004.
- 892 17. Clooney AG, Eckenberger J, Laserna-Mendieta E, Sexton KA, Bernstein MT, Vagianos K, et  
893 al. Ranking microbiome variance in inflammatory bowel disease: a large longitudinal  
894 intercontinental study. *Gut*. 2020. doi:10.1136/gutjnl-2020-321106.
- 895 18. Walters WA, Xu Z, Knight R. Meta-analyses of human gut microbes associated with obesity  
896 and IBD. *FEBS Lett*. 2014;588:4223–33.
- 897 19. McHardy IH, Goudarzi M, Tong M, Ruegger PM, Schwager E, Weger JR, et al. Integrative  
898 analysis of the microbiome and metabolome of the human intestinal mucosal surface reveals  
899 exquisite inter-relationships. *Microbiome*. 2013;1:17.
- 900 20. Vázquez-Baeza Y, Gonzalez A, Xu ZZ, Washburne A, Herfarth HH, Sartor RB, et al. Guiding  
901 longitudinal sampling in IBD cohorts. *Gut*. 2018;67:1743–5.
- 902 21. Liu T-C, Gurram B, Baldridge MT, Head R, Lam V, Luo C, et al. Paneth cell defects in Crohn's  
903 disease patients promote dysbiosis. *JCI Insight*. 2016;1:e86907.
- 904 22. Douglas GM, Hansen R, Jones CMA, Dunn KA, Comeau AM, Bielawski JP, et al. Multi-omics  
905 differentially classify disease state and treatment outcome in pediatric Crohn's disease.  
906 *Microbiome*. 2018;6:13.
- 907 23. Ananthakrishnan AN, Luo C, Yajnik V, Khalili H, Garber JJ, Stevens BW, et al. Gut Microbiome  
908 Function Predicts Response to Anti-integrin Biologic Therapy in Inflammatory Bowel Diseases.  
909 *Cell Host Microbe*. 2017;21:603–10.e3.
- 910 24. Topçuoğlu BD, Lesniak NA, Ruffin MT 4th, Wiens J, Schloss PD. A Framework for Effective  
911 Application of Machine Learning to Microbiome-Based Classification Problems. *MBio*. 2020;11.  
912 doi:10.1128/mBio.00434-20.
- 913 25. Waljee AK, Lipson R, Wiitala WL, Zhang Y, Liu B, Zhu J, et al. Predicting Hospitalization and  
914 Outpatient Corticosteroid Use in Inflammatory Bowel Disease Patients Using Machine Learning.  
915 *Inflamm Bowel Dis*. 2017;24:45–53.
- 916 26. Berg G, Rybakova D, Fischer D, Cernava T, Vergès M-CC, Charles T, et al. Microbiome  
917 definition re-visited: old concepts and new challenges. *Microbiome*. 2020;8:103.
- 918 27. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, et al. Reproducible,  
919 interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol*.  
920 2019;37:852–7.
- 921 28. Hildebrand F, Tadeo R, Voigt AY, Bork P, Raes J. LotuS: an efficient and user-friendly OTU  
922 processing pipeline. *Microbiome*. 2014;2:30.

- 923 29. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal  
924 RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.*  
925 2013;41 Database issue:D590–6.
- 926 30. Frioux C, Singh D, Korcsmaros T, Hildebrand F. From bag-of-genes to bag-of-genomes:  
927 metabolic modelling of communities in the era of metagenome-assembled genomes. *Comput*  
928 *Struct Biotechnol J.* 2020;18:1722–34.
- 929 31. Douglas GM, Maffei VJ, Zaneveld JR, Yurgel SN, Brown JR, Taylor CM, et al. PICRUSt2 for  
930 prediction of metagenome functions. *Nature Biotechnology.* 2020;38:685–8. doi:10.1038/s41587-  
931 020-0548-6.
- 932 32. Ho DSW, Schierding W, Wake M, Saffery R, O’Sullivan J. Machine Learning SNP Based  
933 Prediction for Precision Medicine. *Front Genet.* 2019;10:267.
- 934 33. Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, et al. Normalization and  
935 microbial differential abundance strategies depend upon data characteristics. *Microbiome.*  
936 2017;5:27.
- 937 34. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome Datasets Are  
938 Compositional: And This Is Not Optional. *Frontiers in Microbiology.* 2017;8.  
939 doi:10.3389/fmicb.2017.02224.
- 940 35. Gibbons SM, Duvallet C, Alm EJ. Correcting for batch effects in case-control microbiome  
941 studies. *PLoS Comput Biol.* 2018;14:e1006102.
- 942 36. Wang Y, LêCao K-A. Managing batch effects in microbiome data. *Briefings in Bioinformatics.*  
943 2019. doi:10.1093/bib/bbz105.
- 944 37. Nygaard V, Rødland EA, Hovig E. Methods that remove batch effects while retaining group  
945 differences may lead to exaggerated confidence in downstream analyses. *Biostatistics.*  
946 2015;17:29–39.
- 947 38. Ma S, Shungin D, Mallick H, Schirmer M, Nguyen LH, Kolde R, et al. Population Structure  
948 Discovery in Meta-Analyzed Microbial Communities and Inflammatory Bowel Disease.  
949 doi:10.1101/2020.08.31.261214.
- 950 39. Song K, Wright FA, Zhou Y-H. Systematic Comparisons for Composition Profiles, Taxonomic  
951 Levels, and Machine Learning Methods for Microbiome-Based Disease Prediction. *Front Mol*  
952 *Biosci.* 2020;7:610845.
- 953 40. Thomas AM, Manghi P, Asnicar F, Pasolli E, Armanini F, Zolfo M, et al. Metagenomic analysis  
954 of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures and a link  
955 with choline degradation. *Nat Med.* 2019;25:667–78.
- 956 41. McDonald D, Hyde E, Debelius JW, Morton JT, Gonzalez A, Ackermann G, et al. American  
957 Gut: an Open Platform for Citizen Science Microbiome Research. *mSystems.* 2018;3.  
958 doi:10.1128/mSystems.00031-18.
- 959 42. Estaki M, Pither J, Baumeister P, Little JP, Gill SK, Ghosh S, et al. Cardiorespiratory fitness  
960 as a predictor of intestinal microbial diversity and distinct metagenomic functions. *Microbiome.*  
961 2016;4:42.

- 962 43. Lloyd-Price J, Arze C, Ananthkrishnan AN, Schirmer M, Avila-Pacheco J, Poon TW, et al.  
963 Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature*.  
964 2019;569:655–62.
- 965 44. Doherty MK, Ding T, Koumpouras C, Telesco SE, Monast C, Das A, et al. Fecal Microbiota  
966 Signatures Are Associated with Response to Ustekinumab Therapy among Crohn’s Disease  
967 Patients. *MBio*. 2018;9. doi:10.1128/mBio.02120-17.
- 968 45. Schirmer M, Denson L, Vlamakis H, Franzosa EA, Thomas S, Gotman NM, et al.  
969 Compositional and Temporal Changes in the Gut Microbiome of Pediatric Ulcerative Colitis  
970 Patients Are Linked to Disease Course. *Cell Host Microbe*. 2018;24:600–10.e4.
- 971 46. Vogtmann E, Chen J, Amir A, Shi J, Abnet CC, Nelson H, et al. Comparison of Collection  
972 Methods for Fecal Samples in Microbiome Studies. *Am J Epidemiol*. 2017;185:115–23.
- 973 47. Sankaranarayanan K, Ozga AT, Warinner C, Tito RY, Obregon-Tito AJ, Xu J, et al. Gut  
974 Microbiome Diversity among Cheyenne and Arapaho Individuals from Western Oklahoma. *Curr  
975 Biol*. 2015;25:3161–9.
- 976 48. Hill-Burns EM, Debelius JW, Morton JT, Wissemann WT, Lewis MR, Wallen ZD, et al.  
977 Parkinson’s disease and Parkinson’s disease medications have distinct signatures of the gut  
978 microbiome. *Mov Disord*. 2017;32:739–49.
- 979 49. Obregon-Tito AJ, Tito RY, Metcalf J, Sankaranarayanan K, Clemente JC, Ursell LK, et al.  
980 Subsistence strategies in traditional societies distinguish gut microbiomes. *Nat Commun*.  
981 2015;6:6505.
- 982 50. David LA, Materna AC, Friedman J, Campos-Baptista MI, Blackburn MC, Perrotta A, et al.  
983 Host lifestyle affects human microbiota on daily timescales. *Genome Biol*. 2014;15:R89.
- 984 51. Caporaso JG, Lauber CL, Costello EK, Berg-Lyons D, Gonzalez A, Stombaugh J, et al.  
985 Moving pictures of the human microbiome. *Genome Biol*. 2011;12:R50.
- 986 52. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. DADA2: High  
987 resolution sample inference from amplicon data. doi:10.1101/024034.
- 988 53. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment of  
989 rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol*. 2007;73:5261–7.
- 990 54. Yilmaz P, Parfrey LW, Yarza P, Gerken J, Priesse E, Quast C, et al. The SILVA and “All-  
991 species Living Tree Project (LTP)” taxonomic frameworks. *Nucleic Acids Research*.  
992 2014;42:D643–8. doi:10.1093/nar/gkt1209.
- 993 55. Glöckner FO, Yilmaz P, Quast C, Gerken J, Beccati A, Ciuprina A, et al. 25 years of serving  
994 the community with ribosomal RNA gene reference databases and tools. *J Biotechnol*.  
995 2017;261:169–76.
- 996 56. Ma S. MMUPHin: Meta-analysis Methods with Uniform Pipeline for Heterogeneity in  
997 Microbiome Studies. R package version 0.99.3. 2019.
- 998 57. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the 22nd  
999 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York,  
1000 NY, USA: ACM; 2016. doi:10.1145/2939672.2939785.

- 1001 58. Wang Y, Bhattacharya T, Jiang Y, Qin X, Wang Y, Liu Y, et al. A novel deep learning method  
1002 for predictive modeling of microbiome data. *Brief Bioinform.* 2020. doi:10.1093/bib/bbaa073.
- 1003 59. Lin H, Peddada SD. Analysis of compositions of microbiomes with bias correction. *Nat*  
1004 *Commun.* 2020;11:3514.
- 1005 60. Cross-validation Metrics for Evaluating Classification Performance on Imbalanced Data - IEEE  
1006 Conference Publication.  
1007 [https://ieeexplore.ieee.org/abstract/document/8949568?casa\\_token=GzthpwK9bOkAAAAA:vV-](https://ieeexplore.ieee.org/abstract/document/8949568?casa_token=GzthpwK9bOkAAAAA:vV-LF2CYOeiUi4xgtw_R1B0aAPaQWkUkgBpEYqac4bsB6OceUWdp2kgTuBRLMLAUdS6idoYz0H)  
1008 [LF2CYOeiUi4xgtw\\_R1B0aAPaQWkUkgBpEYqac4bsB6OceUWdp2kgTuBRLMLAUdS6idoYz0H](https://ieeexplore.ieee.org/abstract/document/8949568?casa_token=GzthpwK9bOkAAAAA:vV-LF2CYOeiUi4xgtw_R1B0aAPaQWkUkgBpEYqac4bsB6OceUWdp2kgTuBRLMLAUdS6idoYz0H)  
1009 s. Accessed 27 Nov 2020.
- 1010 61. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1  
1011 score and accuracy in binary classification evaluation. *BMC Genomics.* 2020;21:6.
- 1012 62. Moustafa A, Li W, Anderson EL, Wong EHM, Dulai PS, Sandborn WJ, et al. Genetic risk,  
1013 dysbiosis, and treatment stratification using host genome and gut microbiome in inflammatory  
1014 bowel disease. *Clin Transl Gastroenterol.* 2018;9:e132.
- 1015 63. Davenport M, Poles J, Leung JM, Wolff MJ, Abidi WM, Ullman T, et al. Metabolic alterations  
1016 to the mucosal microbiota in inflammatory bowel disease. *Inflamm Bowel Dis.* 2014;20:723–31.
- 1017 64. Karlsson I, Bostrom H. Handling Sparsity with Random Forests When Predicting Adverse  
1018 Drug Events from Electronic Health Records. 2014 IEEE International Conference on Healthcare  
1019 Informatics. 2014. doi:10.1109/ichi.2014.10.
- 1020 65. Bang S, Yoo D, Kim S-J, Jhang S, Cho S, Kim H. Establishment and evaluation of prediction  
1021 model for multiple disease classification based on gut microbial data. *Sci Rep.* 2019;9:1–9.
- 1022 66. Pawlowsky-Glahn V, Egozcue JJ. Compositional data and their analysis: an introduction.  
1023 Geological Society, London, Special Publications. 2006;264:1–10.
- 1024 67. Bailey K, Sokal RR, James Rohlf F. Biometry: The Principles and Practice of Statistics in  
1025 Biological Research (2nd ed.). *Journal of the American Statistical Association.* 1982;77:946.  
1026 doi:10.2307/2287349.
- 1027 68. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-  
1028 seq data with DESeq2. *Genome Biol.* 2014;15:550.
- 1029 69. Taub MA, Bravo HC, Irizarry RA. Overcoming bias and systematic errors in next generation  
1030 sequencing data. *Genome Med.* 2010;2:1–5.
- 1031 70. Kim YS, Unno T, Kim BY, Park MS. Sex Differences in Gut Microbiota. *World J Mens Health.*  
1032 2020;38:48–60.
- 1033 71. Radjabzadeh D, Boer CG, Beth SA, van der Wal P, Kieft-De Jong JC, Jansen MAE, et al.  
1034 Diversity, compositional and functional differences between gut microbiota of children and adults.  
1035 *Sci Rep.* 2020;10:1040.
- 1036 72. Durbán A, Abellán JJ, Jiménez-Hernández N, Ponce M, Ponce J, Sala T, et al. Assessing gut  
1037 microbial diversity from feces and rectal mucosa. *Microb Ecol.* 2011;61:123–33.
- 1038 73. Mas-Lloret J, Obón-Santacana M, Ibáñez-Sanz G, Guinó E, Pato ML, Rodríguez-Moranta F,



- 1039 et al. Gut microbiome diversity detected by high-coverage 16S and shotgun sequencing of paired  
1040 stool and colon sample. *Sci Data*. 2020;7:92.
- 1041 74. Nagao-Kitamoto H, Kamada N. Host-microbial Cross-talk in Inflammatory Bowel Disease.  
1042 *Immune Netw*. 2017;17:1–12.
- 1043 75. Glassner KL, Abraham BP, Quigley EMM. The microbiome and inflammatory bowel disease.  
1044 *J Allergy Clin Immunol*. 2020;145:16–27.
- 1045 76. Chen L, Wang W, Zhou R, Ng SC, Li J, Huang M, et al. Characteristics of fecal and mucosa-  
1046 associated microbiota in Chinese patients with inflammatory bowel disease. *Medicine* .  
1047 2014;93:e51.
- 1048 77. Wyatt A, Kellermayer R. PCR Based Fecal Pathogen Panel Testing Should Be Interpreted  
1049 with Caution at Diagnosis of Pediatric Inflammatory Bowel Diseases. *Ann Clin Lab Sci*.  
1050 2018;48:674–6.
- 1051 78. Human Microbiome Project Consortium. Structure, function and diversity of the healthy human  
1052 microbiome. *Nature*. 2012;486:207–14.
- 1053 79. Zhernakova A, Kurilshikov A, Bonder MJ, Tigchelaar EF, Schirmer M, Vatanen T, et al.  
1054 Population-based metagenomics analysis reveals markers for gut microbiome composition and  
1055 diversity. *Science*. 2016;352:565–9.
- 1056 80. Turnbaugh PJ, Hamady M, Yatsunencko T, Cantarel BL, Duncan A, Ley RE, et al. A core gut  
1057 microbiome in obese and lean twins. *Nature*. 2009;457:480–4.
- 1058 81. Zhou Y, Xu ZZ, He Y, Yang Y, Liu L, Lin Q, et al. Gut Microbiota Offers Universal Biomarkers  
1059 across Ethnicity in Inflammatory Bowel Disease Diagnosis and Infliximab Response Prediction.  
1060 *mSystems*. 2018;3. doi:10.1128/mSystems.00188-17.
- 1061 82. Sun S, Jones RB, Fodor AA. Inference-based accuracy of metagenome prediction tools varies  
1062 across sample types and functional categories. *Microbiome*. 2020;8:46.
- 1063 83. De Filippis F, Pasolli E, Ercolini D. Newly Explored Faecalibacterium Diversity Is Connected  
1064 to Age, Lifestyle, Geography, and Disease. *Curr Biol*. 2020. doi:10.1016/j.cub.2020.09.063.
- 1065 84. Luo J, Schumacher M, Scherer A, Sanoudou D, Megherbi D, Davison T, et al. A comparison  
1066 of batch effect removal methods for enhancement of prediction performance using MAQC-II  
1067 microarray gene expression data. *Pharmacogenomics J*. 2010;10:278–91.
- 1068 85. Parker HS, Corrada Bravo H, Leek JT. Removing batch effects for prediction problems with  
1069 frozen surrogate variable analysis. *PeerJ*. 2014;2:e561.
- 1070 86. Liu Y, Wang T, Zhou B, Zheng D. Robust integration of multiple single-cell RNA sequencing  
1071 datasets using a single reference space. *Nat Biotechnol*. 2021. doi:10.1038/s41587-021-00859-  
1072 x.
- 1073 87. Bryrup T, Thomsen CW, Kern T, Allin KH, Brandslund I, Jørgensen NR, et al. Metformin-  
1074 induced changes of the gut microbiota in healthy young men: results of a non-blinded, one-armed  
1075 intervention study. *Diabetologia*. 2019;62:1024–35.
- 1076 88. Lewis JD. The utility of biomarkers in the diagnosis and therapy of inflammatory bowel

- 1077 disease. *Gastroenterology*. 2011;140:1817–26.e2.
- 1078 89. Ma C, Battat R, Parker CE, Khanna R, Jairath V, Feagan BG. Update on C-reactive protein  
1079 and fecal calprotectin: are they accurate measures of disease activity in Crohn’s disease? *Expert*  
1080 *Rev Gastroenterol Hepatol*. 2019;13:319–30.
- 1081 90. E Penna FGC, Rosa RM, da Cunha PFS, de Souza SCS, de Abreu Ferrari M de L. Faecal  
1082 calprotectin is the biomarker that best distinguishes remission from different degrees of  
1083 endoscopic activity in Crohn’s disease. *BMC Gastroenterol*. 2020;20:35.
- 1084 91. Zeevi D, Korem T, Zmora N, Israeli D, Rothschild D, Weinberger A, et al. Personalized  
1085 Nutrition by Prediction of Glycemic Responses. *Cell*. 2015;163:1079–94.
- 1086 92. Nguyen VQ, Jiang D, Hoffman SN, Guntaka S, Mays JL, Wang A, et al. Impact of Diagnostic  
1087 Delay and Associated Factors on Clinical Outcomes in a U.S. Inflammatory Bowel Disease  
1088 Cohort. *Inflamm Bowel Dis*. 2017;23:1825–31.
- 1089 93. Zaharie R, Tantau A, Zaharie F, Tantau M, Gheorghe L, Gheorghe C, et al. Diagnostic Delay  
1090 in Romanian Patients with Inflammatory Bowel Disease: Risk Factors and Impact on the Disease  
1091 Course and Need for Surgery. *J Crohns Colitis*. 2016;10:306–14.
- 1092 94. Vavricka SR, Spigaglia SM, Rogler G, Pittet V, Michetti P, Felley C, et al. Systematic  
1093 evaluation of risk factors for diagnostic delay in inflammatory bowel disease. *Inflamm Bowel Dis*.  
1094 2012;18:496–505.
- 1095 95. Vadstrup K, Alulis S, Borsi A, Gustafsson N, Nielsen A, Wennerström ECM, et al. Cost Burden  
1096 of Crohn’s Disease and Ulcerative Colitis in the 10-Year Period Before Diagnosis-A Danish  
1097 Register-Based Study From 2003-2015. *Inflamm Bowel Dis*. 2020;26:1377–82.
- 1098 96. Park KT, Ehrlich OG, Allen JI, Meadows P, Szigethy EM, Henrichsen K, et al. The Cost of  
1099 Inflammatory Bowel Disease: An Initiative From the Crohn’s & Colitis Foundation. *Inflamm Bowel*  
1100 *Dis*. 2020;26:1–10.
- 1101 97. Zhang W, Wong CH, Chavannes M, Mohammadi T, Rosenfeld G. Cost-effectiveness of faecal  
1102 calprotectin used in primary care in the diagnosis of inflammatory bowel disease. *BMJ Open*.  
1103 2019;9:e027043.
- 1104
- 1105
- 1106

1107 **Tables**

1108 **Table 1. Overview of 15 datasets used to compare the effect of different features, data**

1109 **preprocessing methods, and machine learning models on IBD classification performance.**

1110 Available metadata (age, sex, BMI), disease activity, and medication use) is provided for each

1111 dataset. Blank spaces indicate that the respective metadata was not available for the dataset's

1112 samples. The following abbreviations are used: female (F), male (M), and other (O).

1113

Study	Accession	Disease Type	Number of Samples	Sample Type		Age		Sex			BMI		Disease Activity		Medications		
				Stool	Biopsy	Mean	SD	F	M	O	Mean	SD	Active	Rem-ission	Biologics	Immuno-suppressants	5-ASA
American Gut	PRJEB11419	Control	1279	1279	0	46.5	12.2	600	595	1	23.3	2.7					
CVDF	PRJNA308319	Control	39	39	0	25.4	4.2	15	24		24.0	2.9					
		CD	219	219	0	12.0	2.9	87	132								
GEVERSC	PRJEB13680	Control	28	28	0	12.3	3.5	10	18								
		UC	37	37	0	11.8	3.6	22	15								
		CD	689	166	523	19.6	14.2	312	377					15	31	51	
GEVERSM	PRJEB13679	Control	320	7	313	14.0	9.8	157	163								
		UC	268	106	162	24.9	17.5	121	147					2	5	52	
GLS	PRJEB23009	CD	340	340	0	30.2	9.0	215	102		25.7	7.2	43	297	145	74	15
		Control	335	335	0	48.6	14.4	152	166		32.8	8.4					
HMP	ibdmdb.org	Control	43	0	43	28.7	22.0	20	23								
		UC	36	0	36	27.7	17.4	20	16								
MUC	PRJNA317429	CD	35	0	35	14.5	3.5	13	22								
		Control	47	0	47	11.9	3.4	21	25								
PRJNA418765	PRJNA418765	CD	589	589	0	40.4	13.2	332	257		26.4	6.6		589	416		
PRJNA436359	PRJNA436359	UC	1178	917	261	12.6	3.3	582	596				875	303			
QIITA10184	PRJEB13895	Control	962	962	0												
QIITA10342	PRJEB13619	Control	58	58	0	43.2	15.3				31.0	7.5					
QIITA10567	PRJEB14674	Control	133	133	0	70.3	8.6				28.3	5.7					
QIITA1448	PRJEB13051	Control	23	23	0												
QIITA2202	PRJEB6518	Control	516	516	0	29.6	4.8	516									
QIITA550	PRJEB19825	Control	467	467	0	32.8	0.5	131	336								
<b>Total</b>			<b>7707</b>	<b>6221</b>	<b>1486</b>			<b>3358</b>	<b>3048</b>	<b>1</b>			<b>918</b>	<b>1189</b>	<b>578</b>	<b>110</b>	<b>118</b>

1114

1115

1116 **Table 2. Model performance for different sample types and patient demographics.**

1117 Samples with available metadata were categorized into groups based on the collection method  
 1118 or the patient's specific demographic group based on sex, age, and BMI. Predictive performance  
 1119 for all combinations of taxonomic features, compositional normalizations, zero-centering batch  
 1120 effect reduction, and non-linear models were included in the analysis. Logistic regression was  
 1121 performed to assess the performance differences within each sample and demographic group  
 1122 while adjusting for the remaining covariates. \*\*\*\* indicates p-value < 0.0001, and \* indicates p-  
 1123 value < 0.05. Coefficient refers to the corresponding independent variable's coefficient for the  
 1124 logistic regression function and SE refers to the standard error of the coefficient.

Group	Variable	Coefficient	SE
Sample Type	Biopsy (vs. Stool)	-0.44 *	0.2
Life Stage	Adult (vs. Pediatric)	1.39 ****	0.18
BMI Stratification	BMI <30 (vs. BMI > 30)	-0.85 ****	0.19
Sex	Female (vs. Male)	-0.02	0.18
IBD Type	CD (vs. UC)	0.05	0.18

1125

1126

1127 **Table 3. Top three data processing and model pipelines for classifying IBD samples.**

1128 Three combinations which appeared most frequently when all models were sorted by F1 score,  
 1129 accuracy, AUC, or MCC.

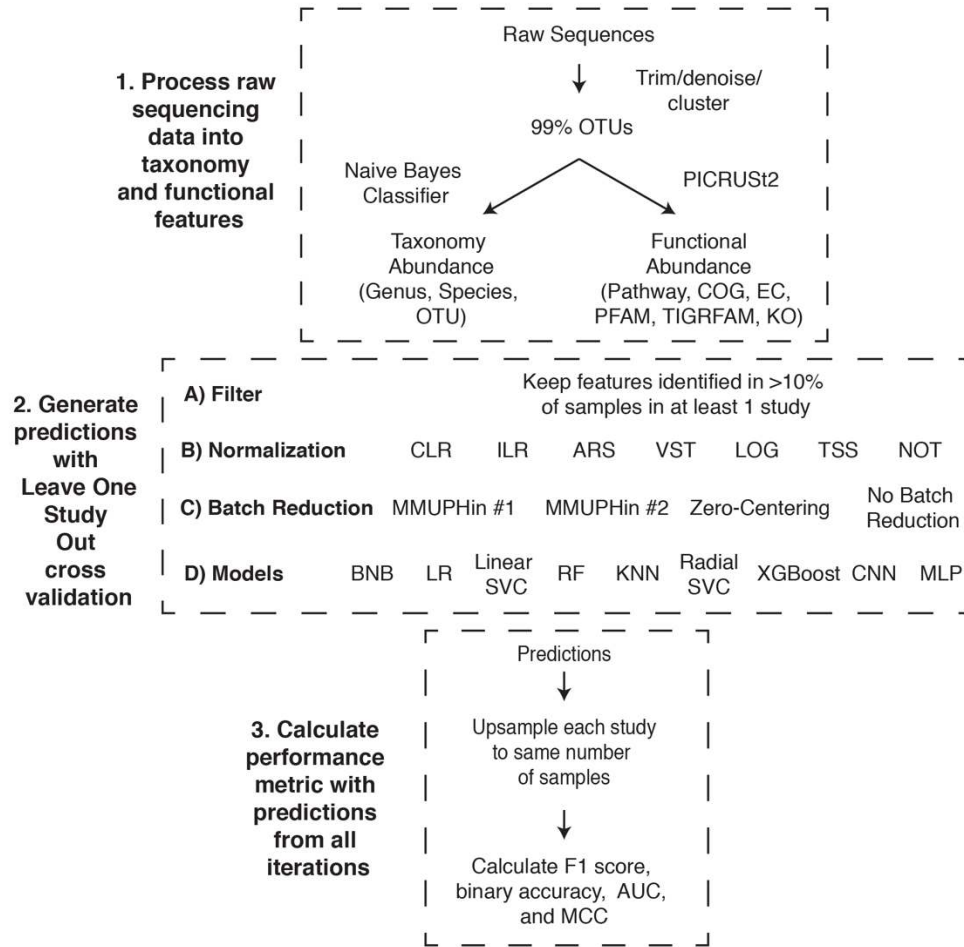
Features	Normalization	Batch Reduction	Model	F1 Score	Accuracy	AUC	MCC
Genus	ILR	Zero-Centering	XGBoost	83.67	87.92	87.9	74.3
Genus	CLR	Zero-Centering	XGBoost	82.96	87.41	87.3	73.2
Genus	ILR	Zero-Centering	Random Forest	82.66	87.4	86.9	72.9

1130

1131

1132

1133 **Figures**



1134

1135 **Figure 1. Leave-one-dataset-out cross-validation pipeline.**

1136 The experiments comprised three different stages to go from raw sequence files to the  
 1137 performance metrics. 1) Raw sequences were processed with Dada2 or Deblur and close-  
 1138 reference clustered into OTUs at 99% identity. The OTUs were classified to taxonomy at 99%  
 1139 identity with QIIME2 and used to infer functional profiles with PICRUSt2. 2) Generating predictions  
 1140 for the 15 iterations of our LODO cross validation consisted of all possible combinations of the  
 1141 listed filtering method, normalization methods, batch effect reduction methods, and models. 3)  
 1142 The predictions from each iteration were combined and the number of samples from each dataset  
 1143 up sampled to 100 000 prior to calculating the performance metrics. The descriptions of acronyms  
 1144 and abbreviations are the following: Clusters of Orthologous Groups of proteins (COG), Kyoto

1145 Encyclopedia of Genes and Genomes (KEGG) orthologs (KO), Enzyme Commission (EC), Pfam  
1146 protein domain (PFAM), TIGR protein family (TIGRFAM) and MetaCyc pathways (pathway),  
1147 centered log-ratio (CLR), isometric log-ratio (ILR), arcsine square root transformation (ARS),  
1148 variance stabilizing transformation (VST), log transformation (LOG), total sum scaling (TSS), no  
1149 normalization (NOT), Bernoulli Naive Bayes (BNB), logistic regression (LR), linear support vector  
1150 machine (Linear SVC), random forest (RF), K nearest neighbours (KNN), radial support vector  
1151 machine (Radial SVC), eXtreme Gradient Boosting (XGBoost), convolutional neural network  
1152 (CNN), multilayer perceptron (MLP).

1153

1154

1155

1156

1157

1158

1159

1160

1161

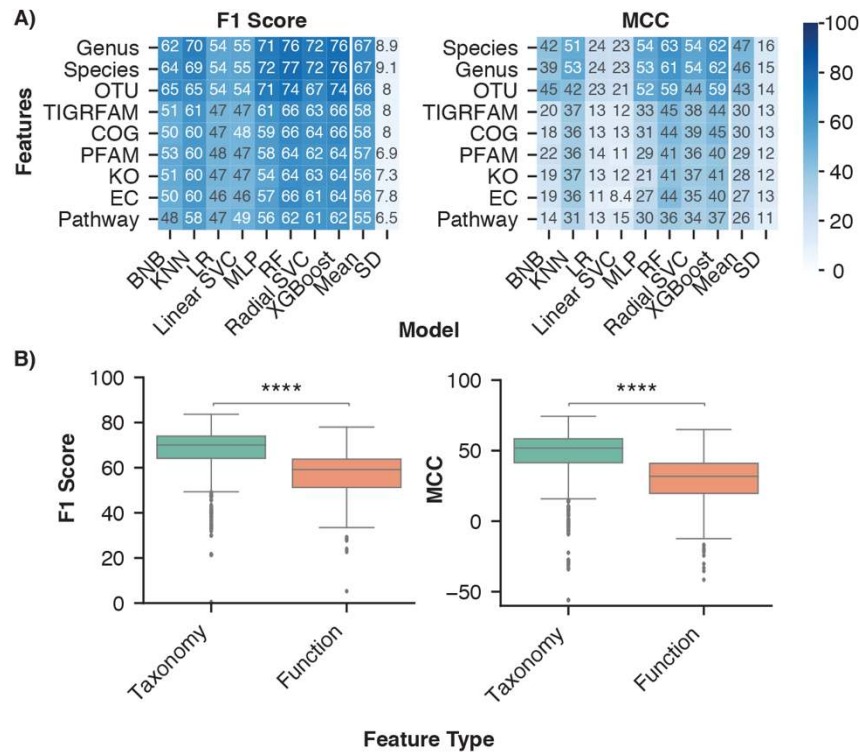
1162

1163

1164

1165

1166

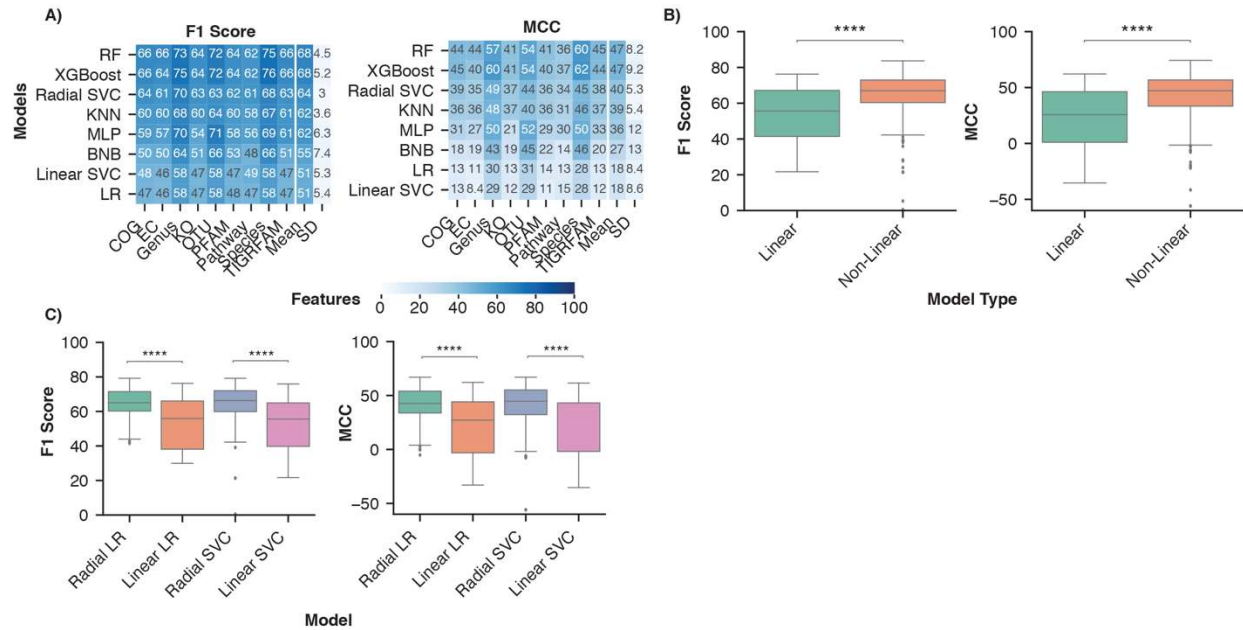


1167

1168 **Figure 2. Optimal disease classification of microbiome samples obtained with taxonomic**  
 1169 **features.**

1170 A) Average performance of the taxonomy and functional feature sets for each ML model  
 1171 architecture. Rows were sorted in descending order by the mean column followed by the standard  
 1172 deviation (SD) column. B) Distribution of performance metrics for taxonomy and functional  
 1173 features across all normalization, batch effect reduction, and model combinations. Independent  
 1174 Mann-Whitney U tests were performed to compare aggregate performance of taxonomy and  
 1175 functional features. The analysis was limited to normalization (ILR, CLR, VST, ARS, LOG, TSS,  
 1176 NOT) and batch effect reduction (no batch reduction or Zero-Centering) methods that were  
 1177 performed on all feature sets. \*\*\*\* indicates p-value < 0.0001.

1178



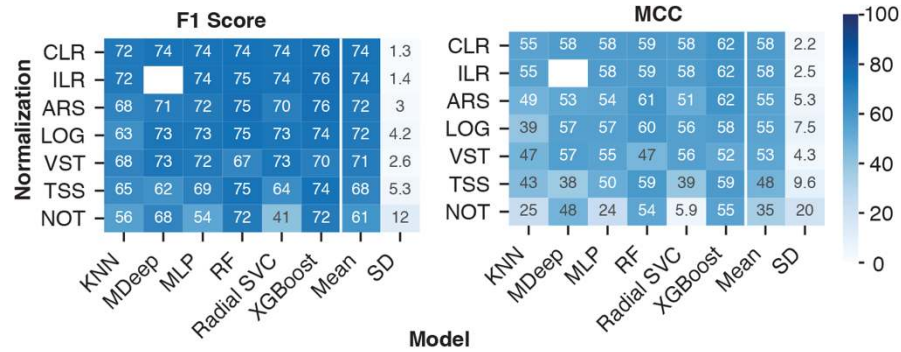
1179

1180 **Figure 3. Non-linear models are better suited to identify decision boundaries between**  
 1181 **control and IBD samples than linear models.**

1182 A) Average model performance for each feature set across normalization and batch effect  
 1183 reduction methods. Rows were sorted in descending order by mean followed by the standard  
 1184 deviation of performance across all feature sets. B) Distribution of performance of non-linear (RF,  
 1185 MLP, KNN, XGBoost, radial SVC) and linear (BNB, Linear SVC, LR) models. Independent Mann-  
 1186 Whitney U-tests were performed to compare each performance metric. The analysis was limited  
 1187 to normalization (ILR, CLR, VST, ARS, LOG, TSS, NOT) and batch effect reduction (no batch  
 1188 effect reduction or zero centering) methods performed on all feature types. C) Distribution of  
 1189 classification performance with the non-linear and linear variations of logistic regression and  
 1190 support vector machines across all feature sets. A Mann-Whitney U test with Bonferroni correction  
 1191 was performed to compare the linear and non-linear variation of each model respectively. \*\*\*\*  
 1192 indicates p-value < 0.0001.

1193





1194

1195 **Figure 4. Compositional normalization methods lead to the highest model performance for**  
 1196 **IBD classification.**

1197 Average model performance with each normalization method across all batch effect reduction  
 1198 methods. Performance following data processing with all pairwise combinations of the  
 1199 normalization methods (ILR, CLR, LOG, ARS, VST, TSS and NOT) and batch effect reduction  
 1200 methods (No batch reduction, MMUPHin #1, MMUPHin #2, and Zero-Centering) were included.  
 1201 Rows were sorted in descending order by the mean and standard deviation of each performance  
 1202 metric across the non-linear models. No analysis was performed for MDeep paired with ILR as  
 1203 the ILR normalized values no longer map directly to a feature, therefore removing the  
 1204 phylogenetic structure required for MDeep.

1205

1206

1207

1208

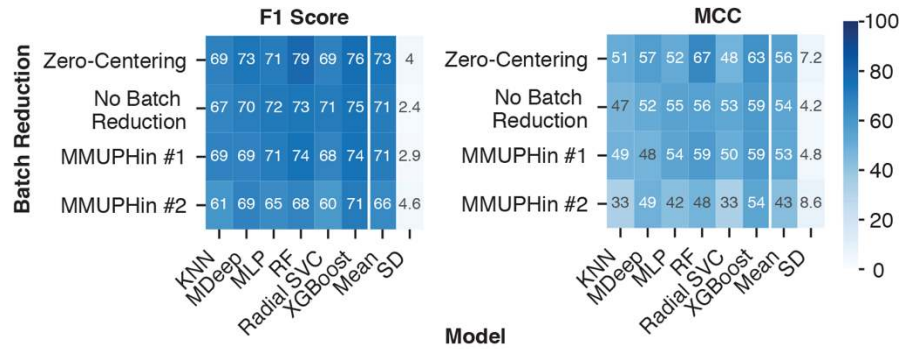
1209

1210

1211

1212

1213



1214

1215 **Figure 5. Batch effect reduction with the naive zero-centering method improved IBD**  
 1216 **classification.**

1217 Average performance of each batch effect reduction method across all combinations of  
 1218 normalization methods, taxonomic features, and non-linear ML models. Performance following  
 1219 data processing with all pairwise combinations of the normalization methods (ILR, CLR, LOG,  
 1220 ARS, VST, TSS and NOT) and batch effect reduction methods (No batch reduction, MMUPHin  
 1221 #1, MMUPHin #2, and Zero-Centering) were included. Rows were sorted in descending order by  
 1222 the mean and standard deviation of performance across all non-linear models.

1223

1224

1225

1226

1227

1228

1229

1230

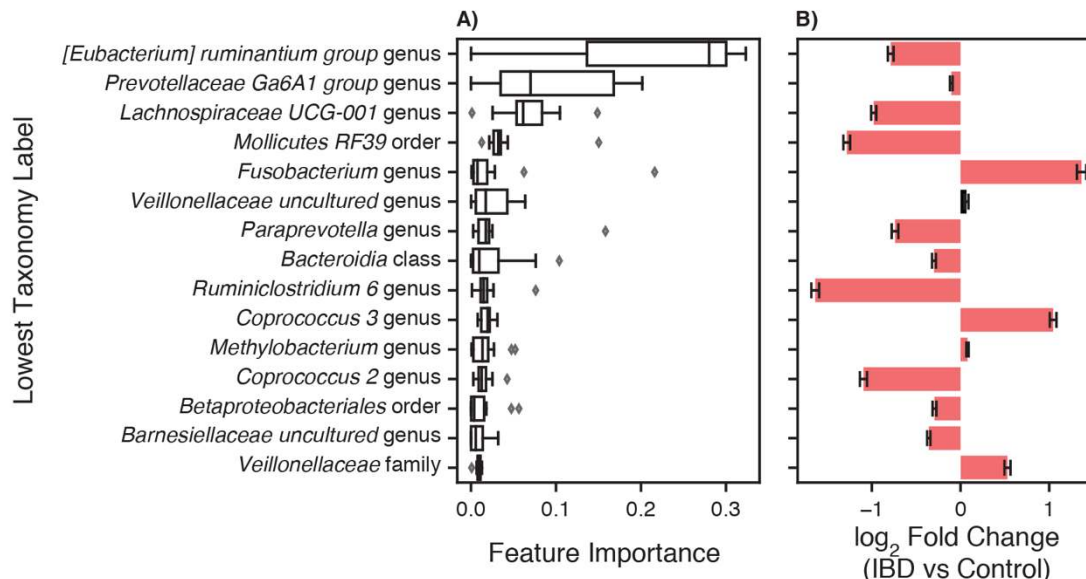
1231

1232

1233

1234

1235



1236

1237 **Figure 6. Features with greatest contribution to IBD classification by an XGBoost classifier.**

1238 A) An XGBoost classifier was trained with CLR normalized genus abundance features and zero-  
1239 centered batch effect reduction for fifteen LODO iterations. The features' gain values for each  
1240 iteration were extracted and sorted by the mean gain across all iterations. The lowest  
1241 classification rank for each feature was used as the label for the corresponding bar. B) Changes  
1242 in taxonomy abundance between control samples and those from patients with IBD. Bars  
1243 represent the fold change  $\pm$  the standard error determined with Analysis of Compositions of  
1244 Microbiomes with Bias Correction (ANCOM-BC). Red indicates a significant fold change between  
1245 IBD and control samples ( $p < 0.05$ ) and black indicates non-significant fold change.

1246

1247

1248

1249

1250

1251

1252 **Supplemental Tables**

1253 **Supplemental Table 1. Model performance on sample and patient demographics.**

1254 Samples were grouped by the five different metadata categories and the classification  
1255 performance with the indicated metric determined for control and CD or UC (depending on the  
1256 metadata group). Coefficient refers to the corresponding independent variable's coefficient for the  
1257 logistic regression function and SE refers to the standard error of the coefficient. \*\*\*\* indicates p-  
1258 value < 0.0001, \*\*\* indicates p-value < 0.001, \*\* indicates p-value < 0.01, and \* indicates p-value  
1259 < 0.05.

<b>Metric</b>	<b>Group</b>	<b>Variable</b>	<b>Coefficient</b>	<b>SE</b>
F1 Score	Sample Type	Biopsy (vs. Stool)	-0.33	0.21
F1 Score	Life Stage	Adult (vs. Pediatric)	0.47 **	0.18
F1 Score	BMI Stratification	BMI <30 (vs. BMI > 30)	0.49 **	0.19
F1 Score	Sex	Female (vs. Male)	0.31	0.19
F1 Score	IBD Type	CD (vs. UC)	0.54 **	0.18
AUC	Sample Type	Biopsy (vs. Stool)	-0.18	0.2
AUC	Life Stage	Adult (vs. Pediatric)	1.25 ****	0.17
AUC	BMI Stratification	BMI <30 (vs. BMI > 30)	-0.17	0.18
AUC	Sex	Female (vs. Male)	0.1	0.18
AUC	IBD Type	CD (vs. UC)	0.07	0.17
Accuracy	Sample Type	Biopsy (vs. Stool)	-0.7 ***	0.2
Accuracy	Life Stage	Adult (vs. Pediatric)	1.03 ****	0.18
Accuracy	BMI Stratification	BMI <30 (vs. BMI > 30)	0.41 *	0.19
Accuracy	Sex	Female (vs. Male)	0.2	0.19
Accuracy	IBD Type	CD (vs. UC)	0.17	0.18

1260

1261

1262

1263

1264

1265

1266 **Supplemental Table 2. Overview of QIIME2 processing for 15 microbiome datasets.**

1267 Samples were collected from the listed ENA accession, with only samples corresponding to  
1268 individuals in North America retained. Trim length was used as input for the trunc\_len parameter,  
1269 forward trim as the trim\_left input for single end read and trim\_left\_f for paired-end reads, and  
1270 reverse trim as the trim\_left\_r input for paired-end reads in the python API for QIIME2's Dada2  
1271 plugin.

Study ID	ENA Accession	Hypervariable Region	Trim Length	Forward Trim	Reverse Trim
American Gut	ERP012803	V4	124	0	0
CVDF	PRJNA308319	V3-V4	290	40	40
GEVERSC	PRJEB13680	V4	174	0	0
GEVERSM	PRJEB13679	V4	174	0	
GLS	PRJEB23009	V4	99	0	
HMP	ibdmdb.org	V4	249	0	0
MUC	PRJNA317429	V4	174	19	21
PRJNA418765	PRJNA418765	V4	245	0	3
PRJNA436359	PRJNA436359	V4	170	0	3
QIITA10184	PRJEB13895	V4	120	0	
QIITA10342	PRJEB13619	V4	100	0	
QIITA10567	PRJEB14674	V4	99	0	
QIITA1448	PRJEB13051	V4	99	0	
QIITA2202	PRJEB6518	V4	99	0	
QIITA550	PRJEB19825	V4	149	0	

1272

1273

1274

1275

1276

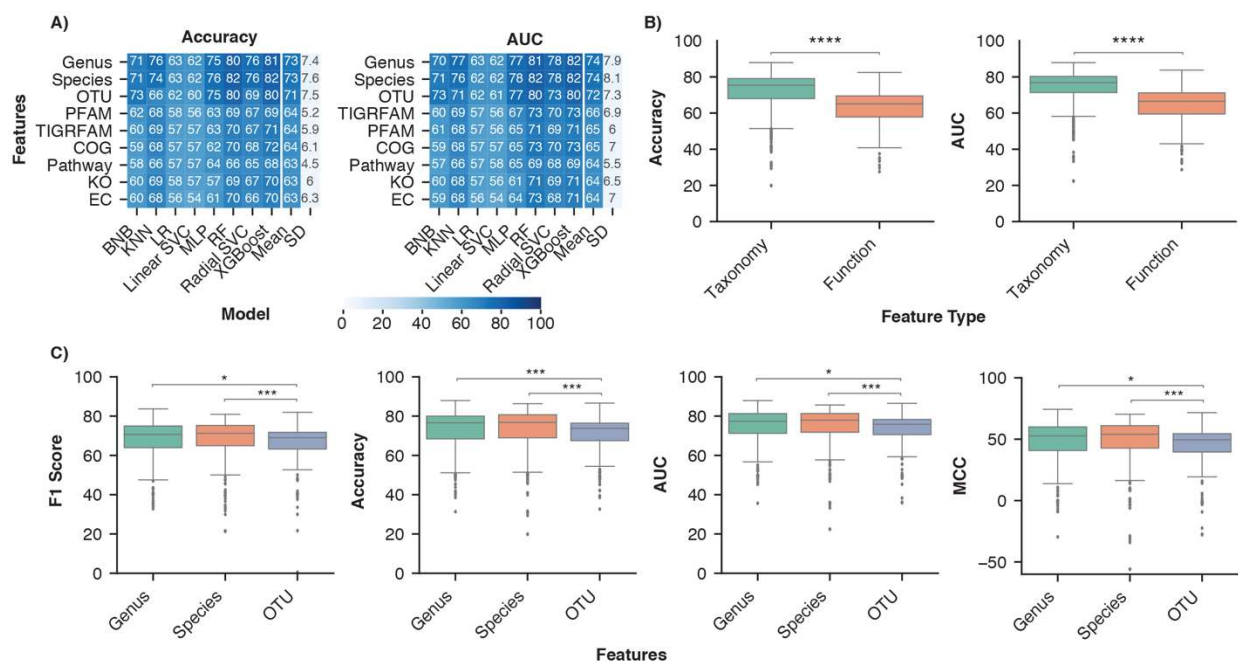
1277

1278

1279

1280

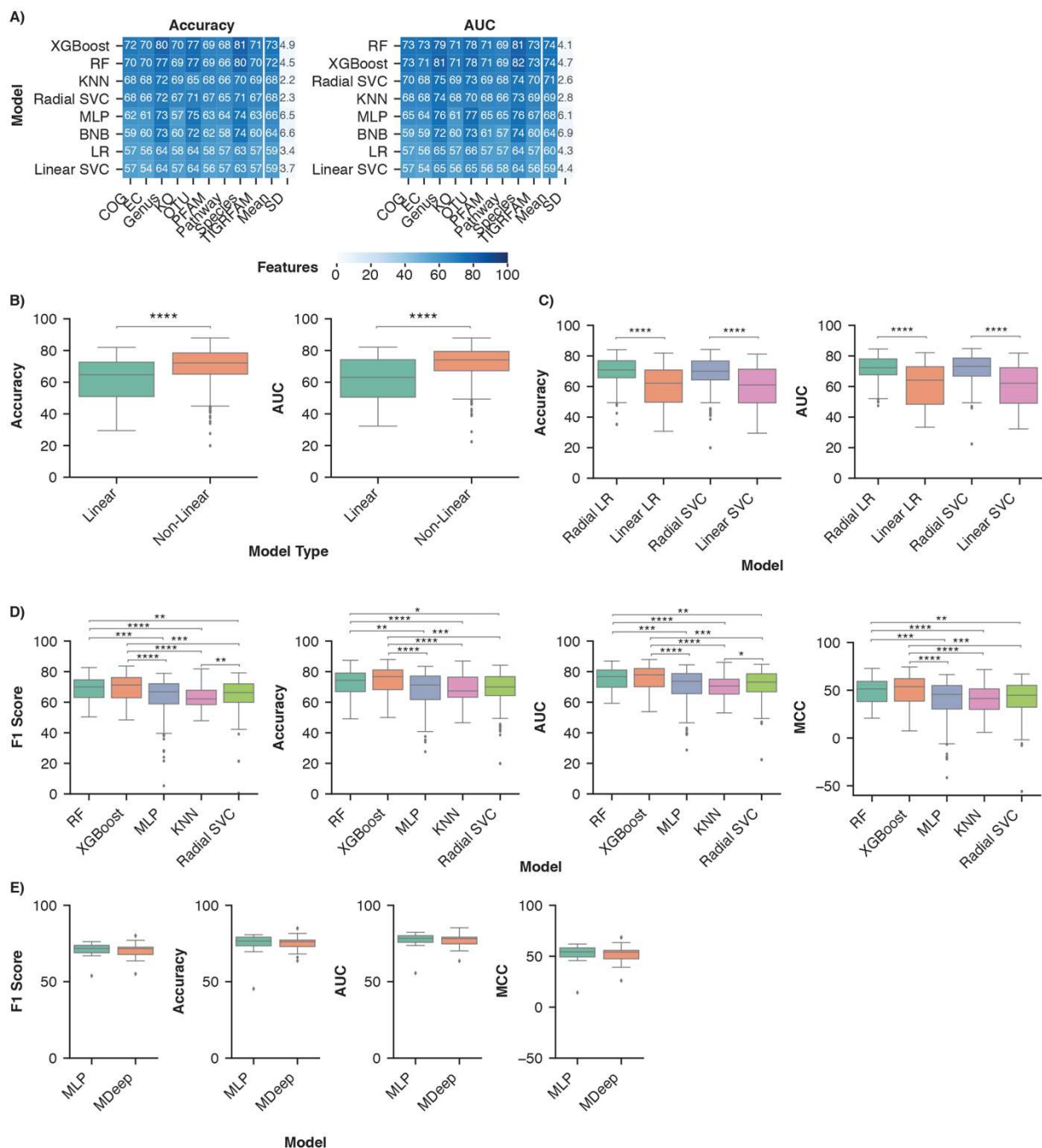
## 1281 Supplemental Figures



1282  
 1283 **Supplemental Figure 1. Greater classification of IBD samples with taxonomic features than**  
 1284 **functional features.**

1285 A) Average performance for each type of feature across different model architectures. Rows were  
 1286 sorted in descending order by the mean column followed by the standard deviation (SD) column.  
 1287 B) Distribution of performance metrics across all normalization, batch effect reduction, and model  
 1288 combinations. Independent Mann-Whitney U-tests were performed to compare aggregate  
 1289 performance measured by each metric of taxonomy and functional features. C) Comparison of  
 1290 classification performance with the three taxonomic feature sets. All pairwise comparisons were  
 1291 performed with a Mann-Whitney U-test followed by Bonferroni correction and the significant  
 1292 comparisons are indicated. The analysis was limited to normalization (ILR, CLR, VST, ARS, LOG,  
 1293 TSS, NOT) and batch effect reduction (No batch reduction or Zero-Centering) methods that were  
 1294 performed on all feature sets. \*\*\*\* indicates p-value < 0.0001, \*\*\* indicates p-value < 0.001, \*  
 1295 indicates p-value < 0.05.

1296



1297

1298 **Supplemental Figure 2. Greater IBD classification performance with non-linear than linear**

1299 **ML models.**

1300 A) Average model performance for each feature set across normalization and batch effect

1301 reduction methods. B) Distribution of performance of non-linear and linear models. Independent

1302 Mann-Whitney U-tests were performed to compare each performance metric. Analysis was limited

1303 to datasets preprocessed using normalization (ILR, CLR, VST, ARS, LOG, TSS, NOT) and batch  
1304 effect reduction (No batch reduction or Zero-Centering) methods performed on all feature types.  
1305 C) Distribution of classification performance with the non-linear and linear variations of logistic  
1306 regression and support vector machines across all feature sets. A Mann-Whitney U test with  
1307 Bonferroni correction was performed to compare the linear and non-linear variation of each model  
1308 respectively. D) Comparison of IBD classification performance between the non-linear models.  
1309 All pairwise comparisons were performed by Mann-Whitney U test with a Bonferroni correction  
1310 and the significant comparisons were labelled. E) Comparison of two neural network  
1311 architectures: the convolutional neural network MDeep or a MLP. A Mann-Whitney U test was  
1312 used to compare each performance metric and the significant comparisons were labelled. \*\*\*\*  
1313 indicates p-value < 0.0001, \*\*\* indicates p-value < 0.001, \*\* indicates p-value < 0.01, and \*  
1314 indicates p-value < 0.05.

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

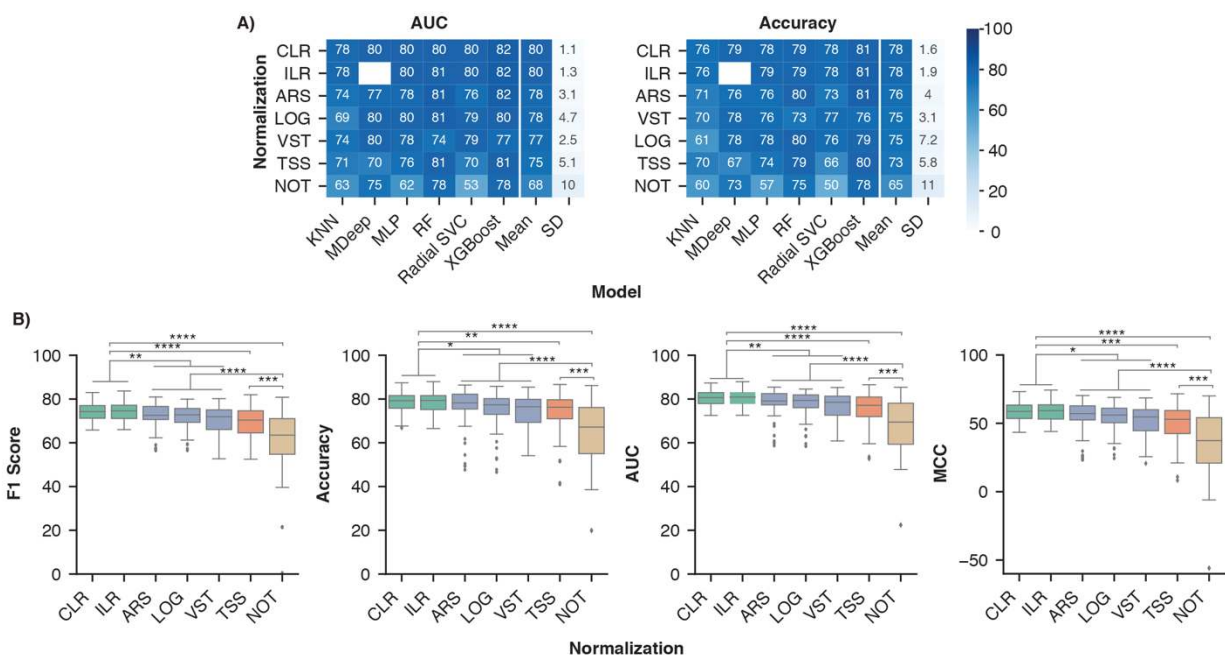
1325

1326

1327

1328





1329

1330 **Supplemental Figure 3. Compositional normalization methods lead to the highest model**  
 1331 **performance for IBD classification.**

1332 A) Average model performance with each normalization method across all batch effect reduction  
 1333 methods. B) Comparing the effect of different classes of normalization methods. The  
 1334 compositional category consists of CLR and ILR (green), variance/distribution modifiers consists  
 1335 of VST, ARS, and LOG (blue), scaling consists of TSS (orange), and no normalization consists of  
 1336 NOT (brown). All pairwise combinations were compared with a Mann-Whitney U test with a  
 1337 Bonferroni correction and the significant comparisons labelled. \*\*\*\* indicates p-value < 0.0001,  
 1338 \*\*\* indicates p-value < 0.001, \*\* indicates p-value < 0.01, and \* indicates p-value < 0.05.

1339

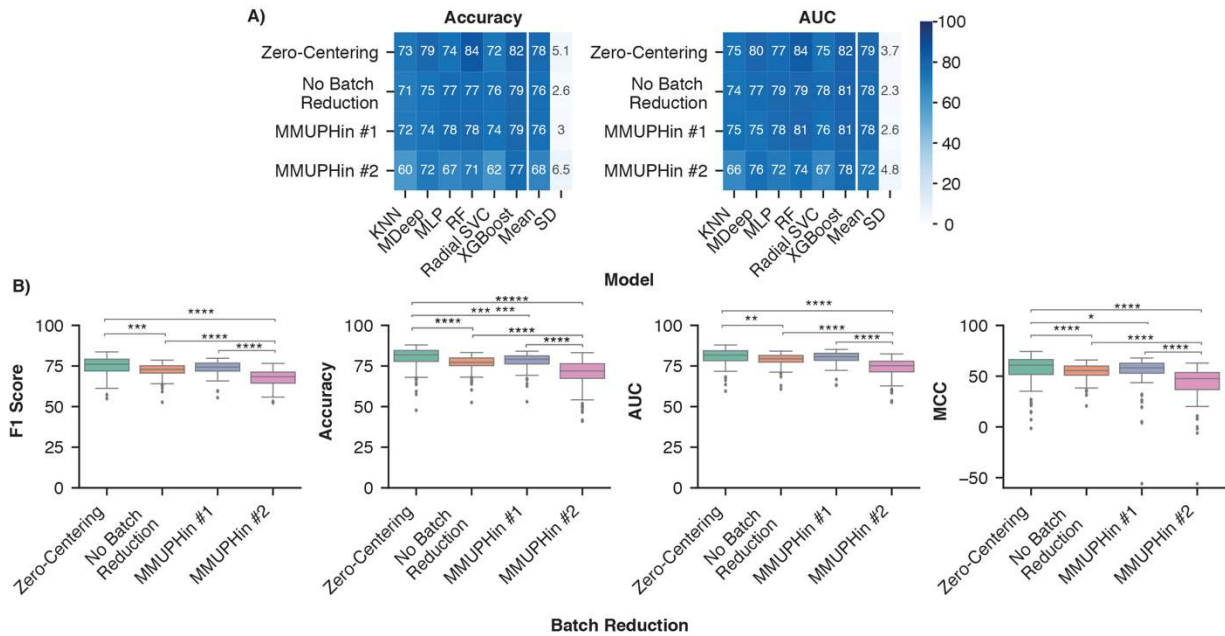
1340

1341

1342

1343

1344



1345

1346 **Supplemental Figure 4. Removing batch effects with zero-centering improved IBD**  
 1347 **classification.**

1348 A) Batch effect reduction methods sorted by average performance across all combinations  
 1349 normalization methods, taxonomic features, and non-linear ML models. B) Comparing the effect  
 1350 of different batch effect reduction methods on classification of IBD samples. All pairwise  
 1351 combinations were compared with a Mann-Whitney U test and the significant comparisons were  
 1352 labelled. \*\*\*\* indicates p-value < 0.0001, \*\*\* indicates p-value < 0.001, \*\* indicates p-value < 0.01,  
 1353 and \* indicates p-value < 0.05.

1354

1355

1356

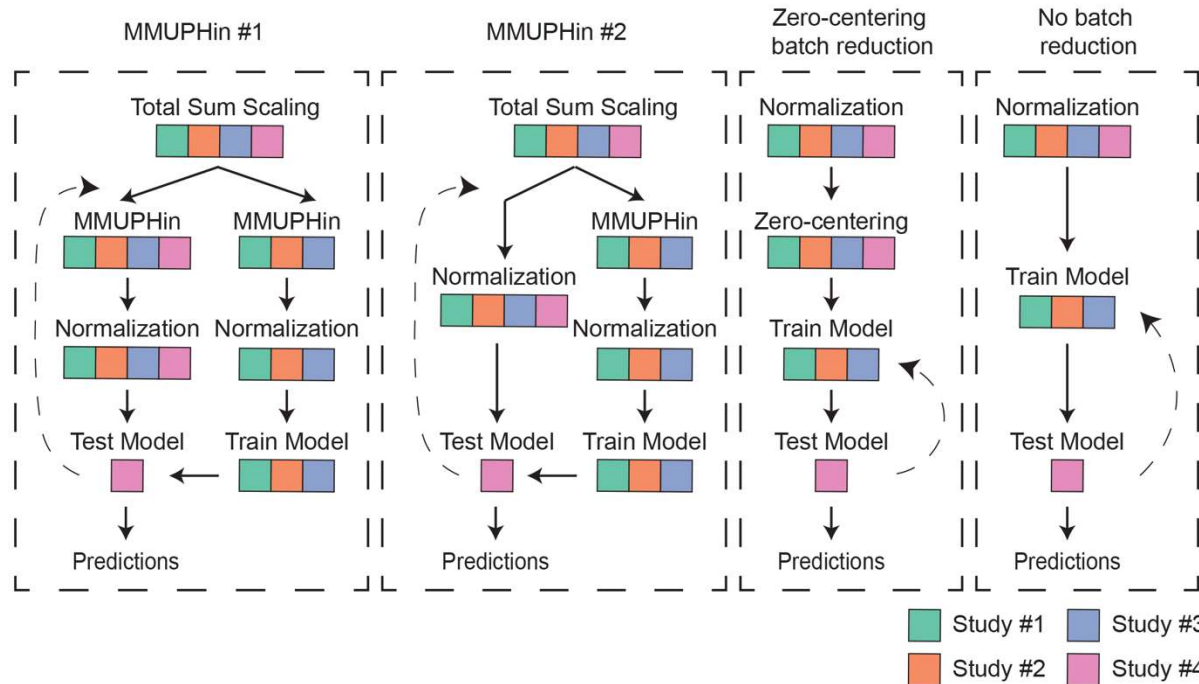
1357

1358

1359

1360

1361



1362

1363 **Supplemental Figure 5. LODO cross validation pipeline for different batch effect reduction**  
 1364 **methods.**

1365 Pipeline performance (combinations of a normalization method, batch effect reduction method,  
 1366 and ML model) was determined with LODO cross validation. Four variations were used to account  
 1367 for different requirements in the batch effect reduction step to ensure the training and test set  
 1368 were independent. The coloured boxes under the text indicate which studies the corresponding  
 1369 step was performed with. The diagram illustrates how a single dataset (illustrated by the different  
 1370 colours) is removed from the dataset the model is trained with and then tested on. The dashed  
 1371 arrow indicates the step returned to in each iteration, with a different dataset removed in each  
 1372 iteration.

1373