

IBD Meta-analysis

Taylor Reiter Luiz Irber ... Phillip Brooks Alicia Gingrich
C. Titus Brown

September 13, 2020

Introduction

Metagenomics captures the functional potential of microbial communities through DNA sequencing of genes and organisms. Metagenomics has been used to profile many human microbial communities, including those that change in or contribute to disease. In particular, human gut microbiomes have been extensively characterized for their potential role in diseases such as obesity (Greenblum, Turnbaugh, and Borenstein 2012), type II diabetes (Qin et al. 2012), colorectal cancer (Wirbel et al. 2019), and inflammatory bowel disease (Lloyd-Price et al. 2019; Morgan et al. 2012; Hall et al. 2017; Franzosa et al. 2019). Inflammatory bowel disease (IBD) refers to a spectrum of diseases characterized by chronic inflammation of the intestines and is likely caused by host-mediated inflammatory responses at least in part elicited by microorganisms (Kostic, Xavier, and Gevers 2014). However, no causative or consistent microbial signature has been associated with IBD to date.

Statements about biology, determined once computation is all done

Although there is no consistent taxonomic or functional trend in the gut microbiome associated with IBD diagnosis, metagenomic studies conducted unto this point have left substantial portions of data unanalyzed. Reference-based pipelines commonly used to analyze metagenomic data from IBD cohorts such as HUMANN2 characterize on average 31%-60% of reads from the human gut microbiome metagenome, as many reads do not closely match sequences in reference databases (Franzosa et al. 2014; Lloyd-Price et al. 2019). To combat this issue, reference-free approaches like *de novo* assembly and binning are used to generate metagenome-assembled genome bins (MAGs) that represent species-level composites of closely related organisms in a sample. However, *de novo* approaches fail when there is low-coverage of or high strain variation in gut microbes, or with sequencing error (Olson et al. 2017). Even when performed on a massive scale, an average of 12.5% of reads fail to map to all *de novo* assembled organisms from human microbiomes (Pasolli et al. 2019).

Here we perform a meta-analysis of six studies of IBD gut metagenome cohorts comprising 260 CD, 132 UC and 213 healthy controls (see **Table 1**) (Lloyd-Price et al. 2019; Lewis et al. 2015; Hall et al. 2017; Franzosa et al. 2019; Gevers et al. 2014; Qin et al. 2010). First, we re-analyzed each study using a consistent k-mer-based, reference-free approach. We demonstrate that diagnosis accounts for a small but significant amount of variation between samples. Next, we used random forests to predict IBD diagnosis and to determine the k-mers that are predictive of UC and CD. Then, we use compact de Bruijn graph queries to reassociate k-mers with sequence context and perform taxonomic and functional characterization of these sequence neighborhoods. We find that strain variation is important (ADD MORE HERE AFTER CORNCOB). Our analysis pipeline is lightweight and is extensible to other association studies in large metagenome sequencing cohorts.

Results

Table 1: Six IBD cohorts used in this meta-analysis.

Cohort	Cohort names	Country	Total	CD	UC	nonIBD	Reference
iHMP	IBDMDB	USA	106	50	30	26	(Lloyd-Price et al. 2019)
PRJEB2054	MetaHIT	Denmark, Spain	124	4	21	99	(Qin et al. 2010)
SRP057027	NA	Canada, USA	112	87	0	25	(Lewis et al. 2015)
PRJNA385949	PRISM, STiNKi	USA	17	9	5	3	(Hall et al. 2017)
PRJNA400072	PRISM, LLDeep, and NLIBD	USA, Netherlands	218	87	76	55	(Franzosa et al. 2019)
PRJNA237362	RISK	North America	28	23	0	5	(Gevers et al. 2014)
Total			605	260	132	213	

Annotation-free approach for meta-analysis of IBD metagenomes.

Given that both reference-based and *de novo* methods suffer from substantial and biased loss of information in the analysis of metagenomes (Thomas and Segata 2019; Breitwieser, Lu, and Salzberg 2019), we sought a reference- and assembly-free pipeline to fully characterize gut metagenomes of IBD patients (**Figure 1**). K-mers, words of length k in nucleotide sequences, have previously been used for annotation-free characterization of sequencing data (reviewed by Rowe (2019)). K-mers are suitable for metagenome analysis because they do not need to be present in reference databases to be included in analysis, and because they capture information from reads even when there is low coverage or high strain variation that preclude assembly. In particular, scaled MinHash sketching produces compressed representations of k-mers in a metagenome while retaining the sequence diversity in a sample (Pierce et al. 2019). Importantly, this approach creates a consistent set of hashes across samples by retaining the same hashes when the same k-mers are observed. This enables comparisons between metagenomes. Given these attributes, we use scaled MinHash sketches to perform metagenome-wide k-mer association with IBD-subtype. We refer to scaled MinHash sketches as *signatures*, and to each subsampled k-mer in a signature as a *hash*.

We also implemented a consistent preprocessing pipeline to remove erroneous sequences that could falsely deflate similarity between samples. We removed adapters, human DNA, and erroneous k-mers, and filtered signatures to retain hashes that were present in multiple signatures. These preprocessing steps removed hashes that were likely to be errors while keeping hashes that were real but low abundance. 7,376,151 hashes remained after preprocessing and filtering.

K-mers capture variation due to disease subtype

We first sought to understand whether variation due to IBD diagnosis is detectable in gut metagenomes.

We calculated pairwise distance matrices using jaccard distance and cosine distance between filtered signatures, where jaccard distance captured sample richness and cosine distance captured sample diversity. We performed principle coordinate analysis and PERMANOVA with these distance matrices (**Figure 2**), using the variables study accession, diagnosis, library size, and number of hashes in a filtered signature (**Table 2**). Number of hashes in a filtered signature

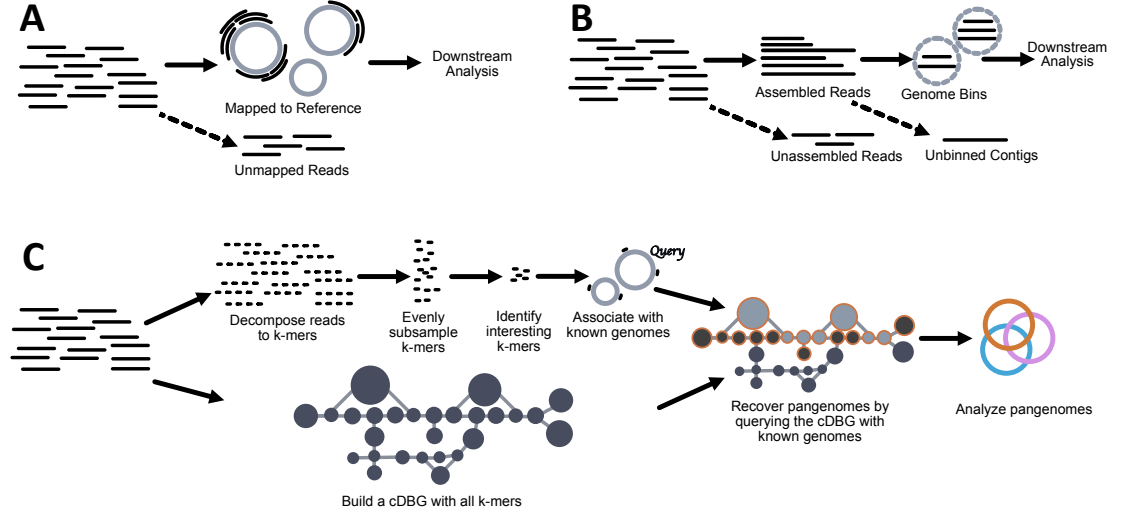


Figure 1: Comparison of common metagenome analysis techniques with the method used in this paper. Metagenomes consist of short (~ 50 - 300 bp) reads derived from sequencing DNA from environmental samples. **A** Reference-based metagenomic analysis. Reads are compared to genomes, genes, or proteins in reference databases to determine the presence and abundance of organisms and proteins in a sample. Unmapped reads are typically discarded from downstream analysis. **B** *De novo* metagenome analysis. Overlapping reads are assembled into longer contiguous sequences (~ 500 bp- 150 kbp, (Vollmers, Wiegand, and Kaster 2017)) and binned into metagenome-assembled genome bins. Bins are analyzed for taxonomy, abundance, and gene content. Reads that fail to assemble and contigs that fail to bin are usually discarded from downstream analysis. **C** Annotation-free approach for meta-analysis of metagenomes. We decompose reads into k-mers and subsample these k-mers, selecting k-mers that evenly represent the sequence diversity within a sample. We then identify interesting k-mers using random forests, and associate these k-mers with genomes in reference databases. Meanwhile, we construct a compact de Bruijn graph (cDBG) that contains all k-mers from a metagenome. We query this graph with known genomes that contain our interesting k-mers to recover sequence diversity nearby our query sequences in the cDBG. In the colored cDBG, light grey nodes indicate nodes that contain at least one identical k-mer to the query, while nodes outlined in orange indicate the nearby sequences recovered via cDBG queries. The combination of all orange nodes produces a sample-specific pangenome that represents the strain variation of closely-related organisms within a single metagenome. We repeat this process for all metagenomes and generate a single metapangenome depicted in orange, blue, and pink.

accounts for the highest variation, possibly reflecting reduced diversity in stool metagenomes of CD and UC patients (reviewed in (Schirmer et al. 2019)). Study accounts for the second highest variation, emphasizing that technical artifacts can introduce biases with strong signals. Diagnosis accounts for a similar amount of variation as study, indicating that there is a small but detectable signal of IBD subtype in stool metagenomes.

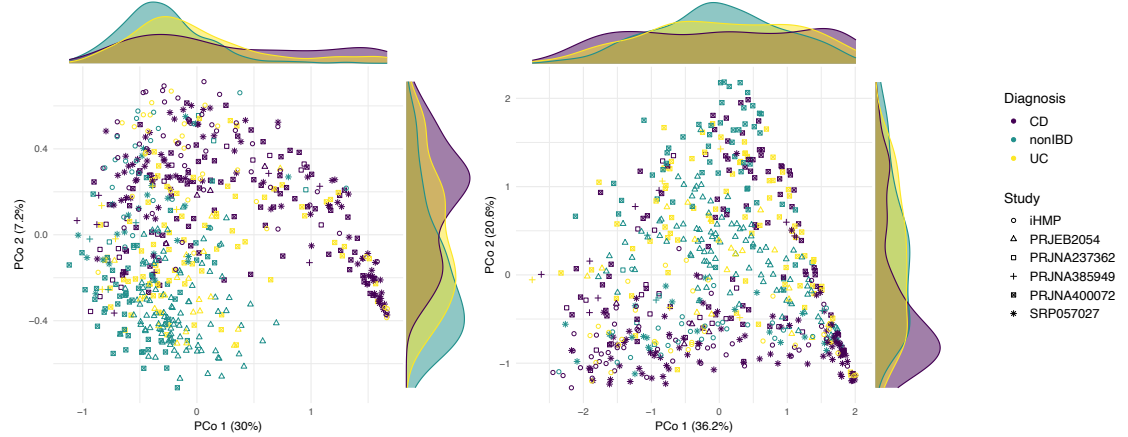


Figure 2: Principle coordinate analysis of metagenomes from IBD cohorts performed on filtered signatures. **A** Jaccard distance. **B** Angular distance.

Table 2: Results from PERMANOVA performed on Jaccard and Angular distance matrices. Number of hashes refers to the number of hashes in the filtered signature, while library size refers to the number of raw reads per sample. * denotes $p < .001$.

Variable	Jaccard distance	Angular distance
Number of hashes	9.9%*	6.2%*
Study accession	6.6%*	13.5%*
Diagnosis	6.2%*	3.3%*
Library size	0.009%*	0.01%*

Hashes are weakly predictive of IBD subtype

To evaluate whether the variation captured by diagnosis is predictive of IBD disease subtype, we built random forests classifiers to predict CD, UC, or non-IBD. We used random forests because of the interpretability of feature importance via variable importance measurements. We used a leave-one-study-out cross-validation approach where we built and optimized a classifier using five cohorts and validated on the sixth.

Given the high-dimensional structure of this dataset (e.g. many more hashes than samples), we first used the vita method to select predictive hashes in the training set (Janitza, Celik, and Boulesteix 2018; Degenhardt, Seifert, and Szymczak 2017). Vita variable selection is based on permutation of variable importance, where p-values for variable importance are calculated against a null distribution that is built from variables that are estimated as non-important (Janitza, Celik, and Boulesteix 2018). This approach retains important variables that are correlated (Janitza, Celik, and Boulesteix 2018; Seifert, Gundlach, and Szymczak 2019), which is desirable in omics-settings where correlated features are often involved in a coordinated biological response, e.g. part of the same operon, pathways, or genome (Stuart et al. 2003; Sabatti et al. 2002). Variable selection reduced the number of hashes used in each model to 29,264-41,701 (Table

91 **3**). Using this reduced set of hashes, we then optimized each random forests classifier on the
 92 training set, producing six optimized models. We validated each model on the left-out study.
 93 The accuracy on the validation studies ranged from 49.1%-75.9% (**Figure 3**), outperforming a
 94 previously published model built on metagenomic data alone (Franzosa et al. 2019).

Table 3: Number of predictive hashes after variable selection for each of 6 classifiers. Classifiers are labelled by the validation study that was held out from training.

Validation study	Selected hashes	Percent of total hashes
PRJNA385949	41701	0.57%
PRJNA237362	40726	0.55%
iHMP	39628	0.54%
PRJEB2054	35343	0.48%
PRJNA400072	32578	0.44%
SRP057027	29264	0.40%

95 We next sought to understand whether there was a consistent biological signal captured among
 96 classifiers by evaluating the fraction of shared hashes between models. We intersected each set
 97 of hashes used to build each optimized classifier (**Figure 3**). Nine hundred thirty two hashes
 98 were shared between all classifiers, while 3,859 hashes were shared between at least five studies.
 99 The presence of shared hashes between classifiers indicates that there is a weak but consistent
 100 biological signal for IBD subtype between cohorts.

101 Shared hashes accounted for 2.8% of all hashes used to build the optimized classifiers. If shared
 102 hashes are predictive of IBD subtype, we would expect that these hashes would account for an
 103 outsized proportion of variable importance in the optimized classifiers. After normalizing variable
 104 importance across classifiers, 40.2% of the total variable importance was held by shared between
 105 hashes, with 21.5% attributable to the 932 hashes shared between all six classifiers. This indicates
 106 that shared hashes contribute a large fraction of predictive power for classification of IBD subtype.

107 Many hashes were identifiable when compared against all microbial genomes in GenBank, as
 108 well as metagenome-assembled genomes from three recent *de novo* assembly efforts from human
 109 microbiome metagenomes (Pasolli et al. 2019; Nayfach et al. 2019; Almeida et al. 2019). 77.7%
 110 of hashes from all classifiers were identifiable and anchored to 1,161 genomes (**Figure 4 A**). In
 111 contrast, 69.4% of shared hashes anchored to 41 genomes (**Figure 4 B**). These shared 41 genomes
 112 held an additional 10.3% of variable importance over the shared hashes because some genomes
 113 contain additional hashes not shared across all models. Using sourmash lca classify to assign
 114 GTDB taxonomy, we find 38 species represented among the 41 genomes (**Figure 4 B**).

115 Marker genes dominate signatures of IB

116 While k-mer based signatures allow us to use all sequencing data in a metagenome and quickly
 117 compare against all known genomes, hashes lack sequence context and do not represent function.
 118 Given this, we next sought to uncover the functional potential of the shared hashes. To annotate
 119 each hash, we reasoned that a hash with predictive importance would be nearby the genomic feature
 120 driving the predictive signal in both genomic DNA and the DNA assembly graph. Therefore, we
 121 performed assembly graph queries on each metagenome using each gene in the 51 shared genomes.
 122 We then identified all gene query neighborhoods in which each hash occurred, and transferred
 123 those annotations to the hash.

124 We noticed that many hashes were annotated as marker genes: 440 hashes accounting for 7.5%
 125 of variable importance across all models annotated to bacterial single copy marker genes (Parks

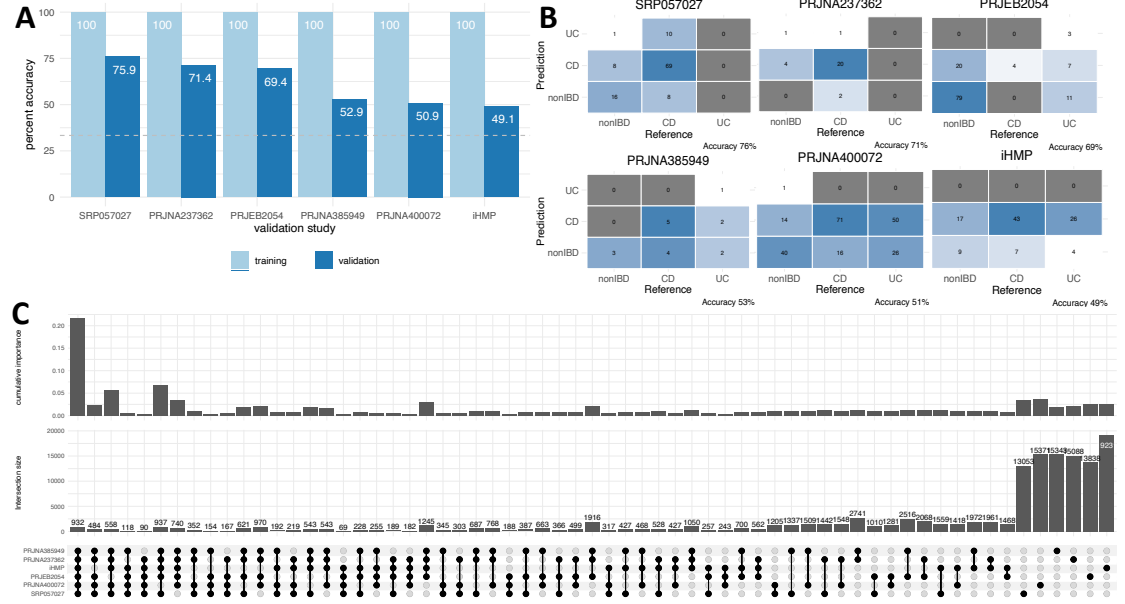


Figure 3: Random forest classifiers weakly predict IBD subtype. **A** Accuracy of leave-one-study-out random forest classifiers on training and validation sets. The validation study is on the x axis. **B** Confusion matrices depicting performance of each leave-one-study-out random forest classifier on the validation set. **C** Upset plot depicting intersections of sets of hashes as well as the cumulative normalized variable importance of those hashes in the optimized random forest classifiers. Each classifier is labelled by the left-out validation study.

et al. 2015; Na et al. 2018), as well as 16S and 23S ribosomal RNA. Given the substantial fraction of variable importance attributable to these genetic elements, we were curious how well models built from marker genes alone would perform in IBD subtype classification. However, we wanted to only look at marker gene abundances from the 41 shared genomes. We first performed cDBG queries using the 41 genomes to retrieve all reads in the assembly graph neighborhoods of those genomes. We then built random forest classifiers using the same approach as with hashes, but using abundance of 14 ribosomal marker genes and 16S rRNA (Woodcroft 2018). Classification accuracy across all models was similar to the k-mer based model (Table 4), however performs marginally better at CD classification and marginally worse at UC classification (Figure SUPPLEMENTAL CONFUSION MATRICES) Both model types ranked a 16S rRNA sequence from the genus *Acetatifactor* as having the highest variable importance across studies, demonstrating that while based on different data features, both model types extract similar information.

Table 4: Accuracy of random forest classifiers built with different underlying representations of IBD metagenomes when applied to each validation set.

Validation Study	Hash model	Ribosomal model	Hash model of ribosomal reads
SRP057027	75.9	86.4	71.7
PRJNA237362	71.4	75	64.3
PRJEB2054	69.4	19.1	15.5
PRJNA385949	52.9	52.9	41.2
PRJNA400072	50.9	48.1	47.4
iHMP	49.1	44.2	46.5

The marker gene model performed similarly as the k-mer model for all studies except PRJEB2054, however this study was sequenced with 36 base pair reads. It performed poorly due to decreased prediction of marker genes from reads. While the hash model performs well even with very short reads, this has limited technological application given ever-increasing sequencing read lengths.

While hash and ribosomal models performed similarly, we were curious whether they captured the same information or whether they captured overlapping but distinct characteristics about IBD

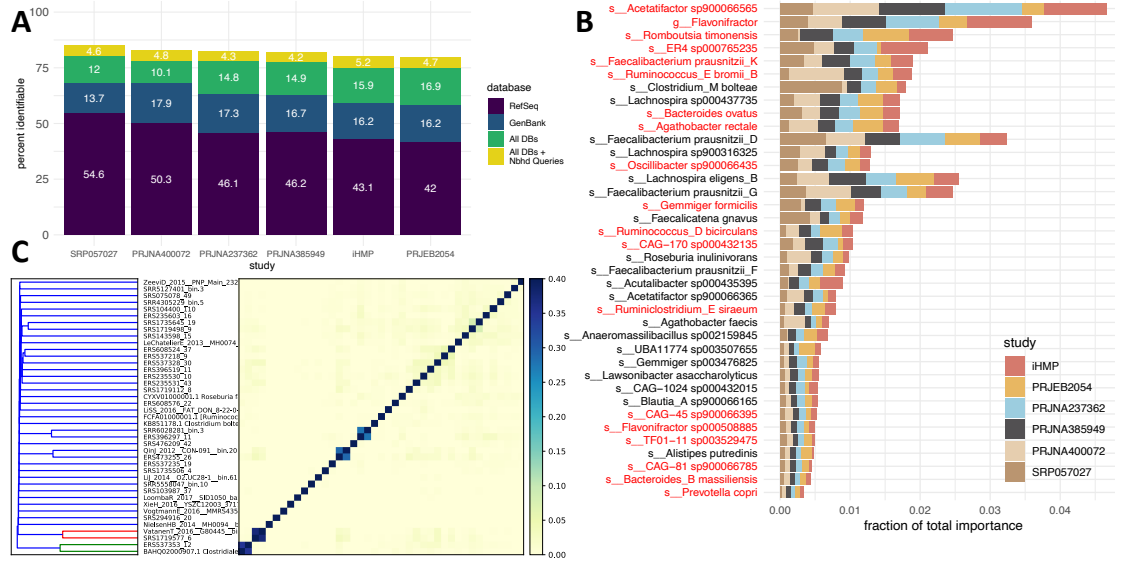


Figure 4: Some predictive hashes from random forest classifiers anchor to known genomes. **A** 75.1-80.3% of all hashes used to train classifiers anchor to known genomes in RefSeq, GenBank, or human microbiome metagenome-assembled genome databases. A further 4.2-5.6% of hashes anchor to metapangenomes of a subset of these genomes. **B** The 3,859 hashes shared between at least five classifiers anchor to 41 genomes. Genomes account for different amounts of variable importance in each model. Genomes are labelled by 38 GTDB taxonomy assignments. **C** Jaccard similarity between 41 genomes. The highest similarity between genomes is 0.37 and is shared by genomes of the same species, while most genomes have no similarity. This indicates that each genome represents distinct nucleotide sequence.

genes with decreased abundance in CD appear to be driven by general decrease in diversity of the gut microbiome, not systematic loss of specific functional potential.

Two strains of *Faecalibacterium prausnitzii* have the largest number of hashes with decreased abundance compared to nonIBD. *F. prausnitzii* is an obligate anaerobe and a key butyrate producer in the gut, thereby playing a crucial role in reducing health inflammation (Lopez-Siles et al. 2017). *F. prausnitzii* is extremely sensitive to oxygen, though may be able to withstand exposure for up to 24 hours depending on the availability of metabolites for extracellular electron transfer (Lopez-Siles et al. 2017). Hashes annotated to *Acetatifactor* (GTDB species *s__Acetatifactor sp900066365) held the largest variable importance of hashes with decreased abundance compared to nonIBD.

*Acetatifactor** is a bile-acid producing bacteria associated with a health gut, but limited evidence has associated it with decreased abundance in IBD (Pathak et al. 2018). In UC, *Gemmiger formicilis* held both the largest variable importance and number of hashes with decreased abundance compared to non-IBD. *G. formicilis* is a strictly anaerobic bacteria that produces both formic acid butyric acid (GOSSLING and Moore 1975).

While most shared genomes decrease in abundance, a substantial portion of hashes from four genomes in CD and two in UC are more abundant in disease. These four belong to *Faecalicatena gnavus* (referred to as *[Ruminococcus] gnavus* in NCBI taxonomy and IBD literature) and *Clostridium bolteae* in CD, and two genomes in the *Flavonifractor* in CD and UC. KEGG enrichment analysis on more and less abundant hashes demonstrated enrichment of pathways dominated by marker genes (e.g. Ribosomes, tRNA biosynthesis) in the less abundant hashes from CD in all four genomes, indicating a general decrease in abundance for these organisms. However, metabolic pathways such as starch and sucrose metabolism, propanoate metabolism, and peptidoglycan synthesis (MAKE FIGURE) are enriched in the more abundant hashes in CD. This shift is indicative of strain-specific enrichment in CD.

ADD:

- kmer accumulation curves
- ...

Other diff abund bio results

c bolt Given these associations, we performed differential abundance analysis on the *C. bolteae* pangenome between CD and nonIBD. We compared our results against study of virulence-causing gene in *C. bolteae* (Lozupone et al. 2012), and find that 24 of 41 previously identified orthologs are significantly induced in CD. Seven of these orthologs are associated with response to oxidative stress. (OXIDATIVE STRESS IBD BIO TIE IN).

We then performed gene enrichment analysis on the differentially abundant genes with KEGG ortholog annotations in *C. bolteae*. While many KEGG pathways are significant, flagellar assembly had the second lowest p value (17 genes). Bacterial flagellin is a dominant antigen in Crohn's disease but not ulcerative colitis (Lodes et al. 2004; Duck et al. 2007). ##### f gnavus We performed differential abundance analysis between CD and nonIBD as well as UC and non IBD to understand whether the metapangenome varied between disease states. While 5,984 genes were differentially abundant in CD, only 197 were less abundant in UC. This suggests that *F. gnavus* is different from nonIBD in CD alone.

We next investigated whether the gene cluster thought to be involved in biosynthesis of the inflammatory polysaccharide was significantly induced in CD. We identified 19 of 23 ORFs in the *F. gnavus* pangenome that matched the putative genes in the cluster, all of which were more abundant in CD. Further, two subsets, one containing 5 ORFs and one contain 7, were co-located on two contiguous sequences, indicating these genes do form a cluster. We then investigated whether this gene cluster was present in non-IBD samples, and found an average of more than

100 reads that mapped per gene in the cluster in 10 of 213 nonIBD metagenomes. This indicates that while more abundant in CD, it is also identifiable within healthy human gut microbiomes. We also genes involved in oxidative stress resistance that are more abundant in CD. This includes one super oxide dismutase and five NADH oxidases. While this evidence supports the idea that *F. gnavus* is harmful in CD, we see some genes that are more abundant in CD that are beneficial for gut health. For example, we find 10 a-L-fucosidases. Tryptophan metabolism. ?

Operons in differentially abundant genes (tmp title)

Given that all genes detected from the *F. gnavus* inflammatory polysaccharide biosynthetic gene cluster were significantly induced in CD, and that subsets of these sequences were colocated on single contiguous sequences, we reasoned that other biologically meaningful genes were likely to occur in clusters. Using results from differential abundance analysis, we searched for gene clusters of five or more genes. We selected five as a signal:noise compromise, as five was the smallest consecutive stretch detected in the *F. gnavus* cluster.

We find no evidence of gene clusters that are more abundant in UC. Conversely, we find many gene clusters in XX pangenomes that are more abundant in CD. XXX

Predictive hashes not in the metapangenomes XXX

- 9.1% of hashes
- sgc query by hash
- Assemble, deepvirfinder, mifaser, compare to viral db, etc.

Discussion

We present XXX.

In this investigation, we find that gut microbiomes from both UC and CD suffer from stochastic loss of diversity.

While *C. bolteae* and *R. gnavus* emerge as bad actors in the pathophysiology of CD, no similar signal is detected for UC. This suggests that while both diseases are associated with lower diversity, CD is uniquely exacerbated by microbes that become more abundant during disease.

Methods

All code associated with our analyses is available at www.github.com/dib-lab/2020-ibd/

IBD metagenome data acquisition and processing

We searched the NCBI Sequence Read Archive and BioProject databases for shotgun metagenome studies that sequenced fecal samples from humans with Crohn's disease, ulcerative colitis, and healthy controls. We included studies sequenced on Illumina platforms with paired-end chemistries and with sample libraries that contained greater than one million reads. For time series intervention cohorts, we selected the first time point to ensure all metagenomes came from treatment-naïve subjects.

245 We downloaded metagenomic fastq files from the European Nucleotide Archive using the
 246 “fastq_ftp” link and concatenated fastq files annotated as the same library into single files.
 247 We also downloaded iHMP samples from idbmdb.org. We used Trimmomatic (version
 248 0.39) to adapter trim reads using all default Trimmomatic paired-end adapter sequences
 249 (ILLUMINACLIP:{inputs/adapters.fa}:2:0:15) and lightly quality-trimmed the reads
 250 (MINLEN:31 LEADING:2 TRAILING:2 SLIDINGWINDOW:4:2) (Bolger, Lohse, and Usadel 2014).
 251 We then removed human DNA using BBDMap and a masked version of hg19 (Bushnell 2014).
 252 Next, we trimmed low-abundance k-mers from sequences with high coverage using khmer’s
 253 `trim-low-abund.py` (Crusoe et al. 2015).

254 Using these trimmed reads, we generated scaled MinHash signatures for each library using
 255 sourmash (k-size 31, scaled 2000, abundance tracking on) (Brown and Irber 2016). At a scaled
 256 value of 2000, an average of one k-mer will be detected in each 2000 base pair window, and 99.8%
 257 of 10,000 base pair windows will have at least one k-mer representative. We selected a k-mer size
 258 of 31 because of its species-level specificity (Koslicki and Falush 2016). A signature is composed
 259 of hashes, where each hash represents a k-mer contained in the original sequence. We retained all
 260 hashes that were present in multiple samples, and refer to these as filtered signatures.

261 Principle Coordinates Analysis

262 We used jaccard distance and cosine distance implemented in `sourmash compare` to pairwise
 263 compare filtered signatures. We then used the `dist()` function in base R to compute distance
 264 matrices. We used the `cmdscale()` function to perform principle coordinate analysis (Gower
 265 1966). We used `ggplot2` and `ggMarginal` to visualize the principle coordinate analysis (Wickham et
 266 al. 2019). To test for sources of variation in these distance matrices, we performed PERMANOVA
 267 using the `adonis` function in the R `vegan` package (Oksanen et al. 2010). The PERMANOVA
 268 was modeled as `~ diagnosis + study accession + library size + number of hashes`.

269 Random forest classifiers

270 We built random forests classifier to predict CD, UC, and non-IBD status using filtered signatures
 271 (hash models), marker genes in the shared 41 genomes (marker gene models), signatures from
 272 reads that were detected as marker genes (hash models of marker genes), and marker genes in the
 273 full metagenome (full marker gene models).

274 For models from signatures, we transformed sourmash signatures into a hash abundance table
 275 where each metagenome was a sample, each hash was a feature, and abundances were recorded for
 276 each hash for each sample. We normalized abundances by dividing by the total number of hashes
 277 in each filtered signature. We then used a leave-one-study-out validation approach where we
 278 trained six models, each of which was trained on five studies and validated on the sixth. To build
 279 each model, we first performed vita variable selection on the training set as implemented in the
 280 `Pomona` and `ranger` packages (Degenhardt, Seifert, and Szymczak 2017; Wright and Ziegler 2015).
 281 Vita variable selection reduces the number of variables (e.g. hashes) to a smaller set of predictive
 282 variables through selection of variables with high cross-validated permutation variable importance
 283 (Janitza, Celik, and Boulesteix 2018). Using this smaller set of hashes, we then built an optimized
 284 random forest model using `tuneRanger` (Probst, Wright, and Boulesteix 2019). We evaluated each
 285 validation set using the optimal model, and extracted variable importance measures for each hash
 286 for subsequent analysis. To make variable importance measures comparable across models, we
 287 normalized importance to 1 by dividing variable importance by the total number of hashes in a
 288 model and the total number of models.

289 For the marker gene models, we generated marker gene abundances for 14 ribosomal marker
 290 genes and 16S rRNA using `singleM` (Woodcroft 2018). We then followed the same model building
 291 procedure as the hash models.

292 Anchoring predictive hashes to genomes

293 We used sourmash `gather` with parameters `k 31` and `--scaled 2000` to anchor predictive hashes
294 to known genomes (Brown and Irber 2016). Sourmash `gather` searches a database of known
295 k-mers for matches with a query (Pierce et al. 2019). We used the sourmash GenBank database
296 (2018.03.29, <https://osf.io/snphy/>), and built three additional databases from medium- and high-
297 quality metagenome-assembled genomes from three human microbiome metagenome reanalysis
298 efforts (<https://osf.io/hza89/>) (Pasolli et al. 2019; Nayfach et al. 2019; Almeida et al. 2019).
299 In total, approximately 420,000 microbial genomes and metagenome-assembled genomes were
300 represented by these four databases. We used the sourmash `lca` commands against the GTDB
301 taxonomy database to taxonomically classify the genomes that contained predictive hashes. To
302 calculate the cumulative variable importance attributable to a single genome, we used an iterative
303 winner-takes-all approach. The genome with the largest fraction of predictive k-mers won the
304 variable importance for all hashes contained within its genome. These hashes were then removed,
305 and we repeated the process for the genome with the next largest fraction of predictive k-mers.

306 To identify hashes that were predictive in at least five of six models, we took the union of predictive
307 hashes from all combinations of five models, as well as from the union of all six models. We refer
308 to these hashes as shared predictive hashes. We anchored variable importance of these shared
309 predictive hashes to known genomes using sourmash `gather` as above.

310 Compact de Bruijn graph queries for predictive genes and genomes

311 To annotate hashes with functional potential, we first extracted open reading frames (ORFs)
312 from the shared 41 genomes using prokka, and annotated ORFs with EggNog (Seemann 2014;
313 Huerta-Cepas et al. 2019). When then used spacegraphcats `multifasta_query` to create a
314 hash:gene map. Spacegraphcats retrieves k-mers in the compact de Bruijn graph neighborhood
315 of a query gene, and hashing these k-mers via sourmash generates a hash:gene map (Brown et
316 al. 2020; Brown and Irber 2016). Because genomes with shared 31-mers may annotate the same
317 hash, we allowed hashes to be annotated multiple times. This was particularly appropriate for
318 hashes from highly conserved regions, e.g. 16S ribosomal RNA.

319 We used spacegraphcats `search` to retrieve k-mers in the compact de Bruijn graph neighborhood
320 of the shared genomes (Brown et al. 2020). We then used spacegraphcats `extract_reads` to
321 retrieve the reads and `extract_contigs` to retrieve unitigs in the compact de Bruijn graph that
322 contained those k-mers, respectively. These reads were used to generate marker gene abundances
323 for the 41 shared genomes for the marker gene random forest models.

324 Differential hash abundance analysis

325 To determine whether shared hashes were differentially abundant from nonIBD in UC or CD, we
326 used corncob (Martin et al. 2020). We used all hash abundances from sourmash signatures to
327 determine hash library size, and then compared hash abundances between disease groups using
328 the likelihood ratio test with the formula `study_accession + diagnosis` and the null formula
329 `study_accession` (Martin et al. 2020). We considered genes with p values < .05 after bonferonni
330 correction as statistically significant. We performed enrichment analysis using the R package
331 clusterProfiler (Yu et al. 2012).

332 References

333 Almeida, Alexandre, Alex L Mitchell, Miguel Boland, Samuel C Forster, Gregory B Gloor,
334 Aleksandra Tarkowska, Trevor D Lawley, and Robert D Finn. 2019. "A New Genomic Blueprint
335 of the Human Gut Microbiota." *Nature* 568 (7753): 499.

336 Bolger, Anthony M, Marc Lohse, and Bjoern Usadel. 2014. "Trimmomatic: A Flexible Trimmer
337 for Illumina Sequence Data." *Bioinformatics* 30 (15): 2114–20.

338 Breitwieser, Florian P, Jennifer Lu, and Steven L Salzberg. 2019. "A Review of Methods and
339 Databases for Metagenomic Classification and Assembly." *Briefings in Bioinformatics* 20 (4):
340 1125–36.

341 Brown, C Titus, and Luiz Irber. 2016. "Sourmash: A Library for Minhash Sketching of Dna." *J.*
342 *Open Source Software* 1 (5): 27.

343 Brown, C Titus, Dominik Moritz, Michael P O'Brien, Felix Reidl, Taylor Reiter, and Blair
344 D Sullivan. 2020. "Exploring Neighborhoods in Large Metagenome Assembly Graphs Using
345 Spacegraphcats Reveals Hidden Sequence Diversity." *Genome Biology* 21 (1): 1–16.

346 Bushnell, Brian. 2014. "BBMap: A Fast, Accurate, Splice-Aware Aligner." Lawrence Berkeley
347 National Lab.(LBNL), Berkeley, CA (United States).

348 Crusoe, Michael R, Hussien F Alameldin, Sherine Awad, Elmar Boucher, Adam Caldwell, Reed
349 Cartwright, Amanda Charbonneau, et al. 2015. "The Khmer Software Package: Enabling Efficient
350 Nucleotide Sequence Analysis." *F1000Research* 4.

351 Degenhardt, Frauke, Stephan Seifert, and Silke Szymczak. 2017. "Evaluation of Variable Selection
352 Methods for Random Forests and Omics Data Sets." *Briefings in Bioinformatics* 20 (2): 492–503.

353 Duck, Wayne L, Mark R Walter, Jan Novak, Denise Kelly, Maurizio Tomasi, Yingzi Cong,
354 and Charles O Elson. 2007. "Isolation of Flagellated Bacteria Implicated in Crohn's Disease."
355 *Inflammatory Bowel Diseases* 13 (10): 1191–1201.

356 Franzosa, Eric A, Xochitl C Morgan, Nicola Segata, Levi Waldron, Joshua Reyes, Ashlee M
357 Earl, Georgia Giannoukos, et al. 2014. "Relating the Metatranscriptome and Metagenome of the
358 Human Gut." *Proceedings of the National Academy of Sciences* 111 (22): E2329–E2338.

359 Franzosa, Eric A, Alexandra Sirota-Madi, Julian Avila-Pacheco, Nadine Fornelos, Henry J Haiser,
360 Stefan Reinker, Tommi Vatanen, et al. 2019. "Gut Microbiome Structure and Metabolic Activity
361 in Inflammatory Bowel Disease." *Nature Microbiology* 4 (2): 293.

362 Gevers, Dirk, Subra Kugathasan, Lee A Denson, Yoshiki Vázquez-Baeza, Will Van Treuren, Boyu
363 Ren, Emma Schwager, et al. 2014. "The Treatment-Naive Microbiome in New-Onset Crohn's
364 Disease." *Cell Host & Microbe* 15 (3): 382–92.

365 GOSSLING, JENNIFER, and WEC Moore. 1975. "Gemmiger Formicilis, N. Gen., N. Sp.,
366 an Anaerobic Budding Bacterium from Intestines." *International Journal of Systematic and*
367 *Evolutionary Microbiology* 25 (2): 202–7.

368 Gower, John C. 1966. "Some Distance Properties of Latent Root and Vector Methods Used in
369 Multivariate Analysis." *Biometrika* 53 (3-4): 325–38.

370 Greenblum, Sharon, Peter J Turnbaugh, and Elhanan Borenstein. 2012. "Metagenomic Systems
371 Biology of the Human Gut Microbiome Reveals Topological Shifts Associated with Obesity and
372 Inflammatory Bowel Disease." *Proceedings of the National Academy of Sciences* 109 (2): 594–99.

373 Hall, Andrew Brantley, Moran Yassour, Jenny Sauk, Ashley Garner, Xiaofang Jiang, Timothy
374 Arthur, Georgia K Lagoudas, et al. 2017. "A Novel Ruminococcus Gnavus Clade Enriched in
375 Inflammatory Bowel Disease Patients." *Genome Medicine* 9 (1): 103.

376 Huerta-Cepas, Jaime, Damian Szklarczyk, Davide Heller, Ana Hernández-Plaza, Sofia K Forslund,
377 Helen Cook, Daniel R Mende, et al. 2019. "EggNOG 5.0: A Hierarchical, Functionally and
378 Phylogenetically Annotated Orthology Resource Based on 5090 Organisms and 2502 Viruses."
379 *Nucleic Acids Research* 47 (D1): D309–D314.

380 Janitza, Silke, Ender Celik, and Anne-Laure Boulesteix. 2018. "A Computationally Fast Variable
381 Importance Test for Random Forests for High-Dimensional Data." *Advances in Data Analysis*

382 *and Classification* 12 (4): 885–915.

383 Koslicki, David, and Daniel Falush. 2016. “MetaPalette: A K-Mer Painting Approach for
384 Metagenomic Taxonomic Profiling and Quantification of Novel Strain Variation.” *MSystems* 1
385 (3): e00020–16.

386 Kostic, Aleksandar D, Ramnik J Xavier, and Dirk Gevers. 2014. “The Microbiome in Inflammatory
387 Bowel Disease: Current Status and the Future Ahead.” *Gastroenterology* 146 (6): 1489–99.

388 Lewis, James D, Eric Z Chen, Robert N Baldassano, Anthony R Otley, Anne M Griffiths, Dale
389 Lee, Kyle Bittinger, et al. 2015. “Inflammation, Antibiotics, and Diet as Environmental Stressors
390 of the Gut Microbiome in Pediatric Crohn’s Disease.” *Cell Host & Microbe* 18 (4): 489–500.

391 Lloyd-Price, Jason, Cesar Arze, Ashwin N Ananthakrishnan, Melanie Schirmer, Julian Avila-
392 Pacheco, Tiffany W Poon, Elizabeth Andrews, et al. 2019. “Multi-Omics of the Gut Microbial
393 Ecosystem in Inflammatory Bowel Diseases.” *Nature* 569 (7758): 655.

394 Lodes, Michael J, Yingzi Cong, Charles O Elson, Raodoh Mohamath, Carol J Landers, Stephan R
395 Targan, Madeline Fort, Robert M Hersherberg, and others. 2004. “Bacterial Flagellin Is a Dominant
396 Antigen in Crohn Disease.” *The Journal of Clinical Investigation* 113 (9): 1296–1306.

397 Lopez-Siles, Mireia, Sylvia H Duncan, L Jesús Garcia-Gil, and Margarita Martinez-Medina. 2017.
398 “Faecalibacterium Prausnitzii: From Microbiology to Diagnostics and Prognostics.” *The ISME*
399 *Journal* 11 (4): 841–52.

400 Lozupone, Catherine, Karoline Faust, Jeroen Raes, Jeremiah J Faith, Daniel N Frank, Jesse
401 Zaneveld, Jeffrey I Gordon, and Rob Knight. 2012. “Identifying Genomic and Metabolic Features
402 That Can Underlie Early Successional and Opportunistic Lifestyles of Human Gut Symbionts.”
403 *Genome Research* 22 (10): 1974–84.

404 Martin, Bryan D, Daniela Witten, Amy D Willis, and others. 2020. “Modeling Microbial
405 Abundances and Dysbiosis with Beta-Binomial Regression.” *Annals of Applied Statistics* 14 (1):
406 94–115.

407 Morgan, Xochitl C, Timothy L Tickle, Harry Sokol, Dirk Gevers, Kathryn L Devaney, Doyle V
408 Ward, Joshua A Reyes, et al. 2012. “Dysfunction of the Intestinal Microbiome in Inflammatory
409 Bowel Disease and Treatment.” *Genome Biology* 13 (9): R79.

410 Na, Seong-In, Yeong Ouk Kim, Seok-Hwan Yoon, Sung-min Ha, Inwoo Baek, and Jongsik
411 Chun. 2018. “UBCG: Up-to-Date Bacterial Core Gene Set and Pipeline for Phylogenomic Tree
412 Reconstruction.” *Journal of Microbiology* 56 (4): 280–85.

413 Nayfach, Stephen, Zhou Jason Shi, Rekha Seshadri, Katherine S Pollard, and Nikos C Kyrpides.
414 2019. “New Insights from Uncultivated Genomes of the Global Human Gut Microbiome.” *Nature*
415 568 (7753): 505.

416 Oksanen, Jari, F Guillaume Blanchet, Roeland Kindt, Pierre Legendre, RB O’hara, Gavin L
417 Simpson, Peter Solymos, M Henry H Stevens, and Helene Wagner. 2010. “Vegan: Community
418 Ecology Package. R Package Version 1.17-4.” *Http://Cran. R-Project. Org>. Acesso Em* 23:
419 2010.

420 Olson, Nathan D, Todd J Treangen, Christopher M Hill, Victoria Cepeda-Espinoza, Jay Ghurye,
421 Sergey Koren, and Mihai Pop. 2017. “Metagenomic Assembly Through the Lens of Validation: Re-
422 cent Advances in Assessing and Improving the Quality of Genomes Assembled from Metagenomes.”
423 *Briefings in Bioinformatics*.

424 Parks, Donovan H, Michael Imelfort, Connor T Skennerton, Philip Hugenholtz, and Gene W
425 Tyson. 2015. “CheckM: Assessing the Quality of Microbial Genomes Recovered from Isolates,
426 Single Cells, and Metagenomes.” *Genome Research* 25 (7): 1043–55.

427 Pasolli, Edoardo, Francesco Asnicar, Serena Manara, Moreno Zolfo, Nicolai Karcher, Federica
 428 Armanini, Francesco Beghini, et al. 2019. "Extensive Unexplored Human Microbiome Diversity
 429 Revealed by over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle."
 430 *Cell* 176 (3): 649–62.

431 Pathak, Preeti, Cen Xie, Robert G Nichols, Jessica M Ferrell, Shannon Boehme, Kristopher W
 432 Krausz, Andrew D Patterson, Frank J Gonzalez, and John YL Chiang. 2018. "Intestine Farnesoid
 433 X Receptor Agonist and the Gut Microbiota Activate G-Protein Bile Acid Receptor-1 Signaling
 434 to Improve Metabolism." *Hepatology* 68 (4): 1574–88.

435 Pierce, N Tessa, Luiz Irber, Taylor Reiter, Phillip Brooks, and C Titus Brown. 2019. "Large-Scale
 436 Sequence Comparisons with Sourmash." *F1000Research* 8.

437 Probst, Philipp, Marvin N Wright, and Anne-Laure Boulesteix. 2019. "Hyperparameters and
 438 Tuning Strategies for Random Forest." *Wiley Interdisciplinary Reviews: Data Mining and
 439 Knowledge Discovery* 9 (3): e1301.

440 Qin, Junjie, Ruiqiang Li, Jeroen Raes, Manimozhiyan Arumugam, Kristoffer Solvsten Burgdorf,
 441 Chaysavanh Manichanh, Trine Nielsen, et al. 2010. "A Human Gut Microbial Gene Catalogue
 442 Established by Metagenomic Sequencing." *Nature* 464 (7285): 59.

443 Qin, Junjie, Yingrui Li, Zhiming Cai, Shenghui Li, Jianfeng Zhu, Fan Zhang, Suisha Liang, et al.
 444 2012. "A Metagenome-Wide Association Study of Gut Microbiota in Type 2 Diabetes." *Nature*
 445 490 (7418): 55.

446 Rowe, Will PM. 2019. "When the Levee Breaks: A Practical Guide to Sketching Algorithms for
 447 Processing the Flood of Genomic Data." *Genome Biology* 20 (1): 199.

448 Sabatti, Chiara, Lars Rohlin, Min-Kyu Oh, and James C Liao. 2002. "Co-Expression Pattern
 449 from Dna Microarray Experiments as a Tool for Operon Prediction." *Nucleic Acids Research* 30
 450 (13): 2886–93.

451 Schirmer, Melanie, Ashley Garner, Hera Vlamakis, and Ramnik J Xavier. 2019. "Microbial Genes
 452 and Pathways in Inflammatory Bowel Disease." *Nature Reviews Microbiology* 17 (8): 497–511.

453 Seemann, Torsten. 2014. "Prokka: Rapid Prokaryotic Genome Annotation." *Bioinformatics* 30
 454 (14): 2068–9.

455 Seifert, Stephan, Sven Gundlach, and Silke Szymczak. 2019. "Surrogate Minimal Depth as an
 456 Importance Measure for Variables in Random Forests." *Bioinformatics* 35 (19): 3663–71.

457 Stuart, Joshua M, Eran Segal, Daphne Koller, and Stuart K Kim. 2003. "A Gene-Coexpression
 458 Network for Global Discovery of Conserved Genetic Modules." *Science* 302 (5643): 249–55.

459 Thomas, Andrew Maltez, and Nicola Segata. 2019. "Multiple Levels of the Unknown in Microbiome
 460 Research." *BMC Biology* 17 (1): 48.

461 Vollmers, John, Sandra Wiegand, and Anne-Kristin Kaster. 2017. "Comparing and Evaluating
 462 Metagenome Assembly Tools from a Microbiologist's Perspective-Not Only Size Matters!" *PloS
 463 One* 12 (1): e0169662.

464 Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy McGowan, Romain
 465 François, Garrett Grolemond, et al. 2019. "Welcome to the Tidyverse." *Journal of Open Source
 466 Software* 4 (43): 1686.

467 Wirbel, Jakob, Paul Theodor Pyl, Ece Kartal, Konrad Zych, Alireza Kashani, Alessio Milanese,
 468 Jonas S Fleck, et al. 2019. "Meta-Analysis of Fecal Metagenomes Reveals Global Microbial
 469 Signatures That Are Specific for Colorectal Cancer." *Nature Medicine* 25 (4): 679.

470 Woodcroft, B. 2018. "Singlem."

- 471 Wright, Marvin N, and Andreas Ziegler. 2015. “Ranger: A Fast Implementation of Random
472 Forests for High Dimensional Data in C++ and R.” *arXiv Preprint arXiv:1508.04409*.
- 473 Yu, Guangchuang, Li-Gen Wang, Yanyan Han, and Qing-Yu He. 2012. “ClusterProfiler: An
474 R Package for Comparing Biological Themes Among Gene Clusters.” *Omics: A Journal of*
475 *Integrative Biology* 16 (5): 284–87.