# Rational Design of Live Biotherapeutic Products for the Prevention of *Clostridioides difficile* Infection

Shanlin Ke[1], Javier A Villafuerte Gálvez[2], Zheng Sun[1], Yangchun Cao[2,3], Nira R Pollock[4,5], Xinhua Chen[2,#], Ciarán P Kelly[2,#], Yang-Yu Liu[1,6,#]


[1]*Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts 02115, USA.*
[2]*Division of Gastroenterology, Department of Medicine, Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, Massachusetts 02215, USA.*
[3]*College of Animal Science and Technology, Northwest A&F University, Yangling, Shaanxi 712100, People's Republic of China.*
[4]*Division of Infectious Disease, Department of Medicine, Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, Massachusetts 02215, USA.*
[5]*Department of Laboratory Medicine, Boston Children's Hospital, Boston, Massachusetts, USA.*
[6]*Center for Artificial Intelligence and Modeling, The Carl R. Woese Institute of Genomic Biology, University of Illinois at Urbana-Champaign, Champaign, IL, USA.*


[#] Correspondence: Y.-Y.L. (yyl@channing.harvard.edu); C.P.K. (ckelly2@bidmc.harvard.edu); X.C. (xchen1@bidmc.harvard.edu).

1

## Abstract

*Clostridioides difficile* infection (CDI) is one of the leading causes of healthcare- and antibiotic-associated diarrhea. While fecal microbiota transplantation (FMT) has emerged as a promising therapy for recurrent CDI, its exact mechanisms of action and long-term safety are not fully understood. Defined consortia of clonal bacterial isolates, known as live biotherapeutic products (LBPs), have been proposed as an alternative therapeutic option. However, the rational design of LBPs remains challenging. Here, we employ a computational pipeline and three independent metagenomic datasets to systematically identify microbial strains that have the potential to inhibit CDI. We first constructed the CDI-related microbial genome catalog, comprising 3,741 non-redundant metagenome-assembled genomes (nrMAGs) at the strain level. We then identified multiple potential protective nrMAGs that can be candidates for the design of microbial consortia targeting CDI, including strains from *Dorea formicigenerans*, *Oscillibacter welbionis*, and *Faecalibacterium prausnitzii*. Importantly, some of these potential protective nrMAGs were found to play an important role in the success of FMT, and the majority of the top protective nrMAGs can be validated by various previously reported findings. Our results demonstrate a computational framework for the rational selection of microbial strains targeting CDI, paving the way for the computational design of microbial consortia against other enteric infections.

## Introduction

*Clostridioides difficile* infection (CDI) is one of the leading causes of healthcare- and antibiotic-associated diarrhea, affecting roughly 500,000 patients and leading to almost 30,000 deaths annually in the United States[1,2]. Exposure to toxinogenic *C. difficile* can lead to a spectrum of clinical outcomes, including asymptomatic colonization, mild diarrhea, and more severe disease syndromes such as pseudomembranous colitis, toxic megacolon, bowel perforation, sepsis, and death[3]. Antibiotics serve as the standard treatment for primary CDI[4,5]. However, CDI recurrence occurs in approximately a quarter of cases after antibiotic treatment[6,7]. Once CDI recurs, patients may get into a vicious cycle of antibiotic therapy and relapse[8]. Moreover, the use of antibiotics has been identified as the primary risk factor for developing CDI, and reports of strains with decreased sensitivity to vancomycin are becoming more frequent.

The human gut microbiome is critical in providing colonization resistance against exogenous pathogens through complex mechanisms such as nutrient competition, competitive metabolic interactions, niche exclusion, and induction of the host immune response[9]. Intestinal microbiota restoration, such as fecal microbiota transplantation (FMT), has been shown to be effective for CDI treatment as well as the restoration of colonization resistance against *C. difficile*[10,11]. While FMT has emerged as a promising therapy for recurrent CDI (rCDI), its exact mechanisms of action are not fully understood[12]. In addition, FMT has the potential to transmit undetected or emerging pathogens, which may result in hospitalization or even death[13,14]. Recently, the FDA has approved fecal microbiota products (e.g., Rebyota[15] and Vowst[16]) for the prevention of rCDI in individuals 18 years of age and older, following antibiotic treatment for rCDI. Rebyota is a room temperature shelf stable suspension of healthy donor stool[17], although its clinical effect size for the prevention of rCDI is modest (RR, 1.17; 95% CI, 0.99–1.39)[18] and its microbial composition is not predefined[19]. Although Vowst is a formulation of live fecal microbiota consisting of a highly purified collection of about 50 species of *Firmicutes* spores with a more robust clinical effect size (1.46; 95% CI, 1.21–1.75)[18], the ecological principle underlying the selection of these microbial strains is unclear.

94

95 The variability of biological properties among bacterial strains within the same species
96 underscores the significance of conducting strain-level composition analysis to
97 understand the role of the human microbiome in human health and disease[20]. For
98 example, some strains from *Escherichia coli* (e.g., *E. coli O157:H7*) cause severe
99 abdominal pain, bloody diarrhea, and vomiting[21]. In contrast, *E. coli Nissle 1917* is a
100 non-pathogenic strain that has been utilized as a probiotic agent to treat gastrointestinal
101 infections in humans[22,23]. Whole metagenome shotgun (WMS) sequencing is a rapid,
102 cost-effective, and high-throughput technology for profiling microbial communities in
103 human microbiome studies[24]. However, precise identification of microorganisms at the
104 strain level remains challenging. Additionally, traditional strain-level profilers can only
105 identify strains within the reference genome databases[25]. These databases are subject
106 to limitations and biases and are unable to characterize microbes that do not have high-
107 quality reference genomes. To resolve these limitations, an alternative strategy for
108 WMS data analysis involves reconstructing metagenome-assembled genomes (MAGs)
109 through *de novo* assembly and binning, offering the advantage of recovering genomes
110 for uncultured microorganisms absent from current reference databases[26].

111

112 In this study, we leveraged a novel computational framework[27] we previously developed
113 to rationally design a bacterial consortium against CDI (**Fig. 1**). The metagenome
114 assembly and binning strategies were applied to reconstruct microbial population
115 genomes directly from the microbiome samples of two independent CDI-related cohorts
116 as well as the healthy controls from the Human Microbiome Project (HMP). Specifically,
117 we sought to identify known and unknown taxa at the strain level, quantify the degree of
118 donor strain engraftment, and design a candidate bacterial consortium against CDI.

119

## Results

### Study cohorts and metagenomic datasets

122 To rationally design microbial consortia against *C. difficile*, we aimed to infer species
123 that may inhibit *C. difficile* from CDI-related microbiome samples. We first collected
124 WMS sequencing data from our in-house clinical cohort (denoted as BIDMC-cohort

125 hereafter) [28,29]. Our BIDMC cohort consists of 104 well-characterized recruited

126 participants divided into four groups (**Table S1** and **Fig. 2a**): (1) Control (CON, n = 26);

127 (2) Non-CDI Diarrhea (NCD, n = 14); (3) Asymptomatic Carriage of *C. difficile* (ASC, n =

128 17); (4) CDI (n = 47). Given the fact that the participants from the CON group were not

129 healthy people, we retrieved an FMT study[30] (denoted as Verma-cohort hereafter) with

130 publicly available data that assessed the microbiome composition of donors (n=21) and

131 recipients (n=22, pre-and post-FMT) through WMS sequencing (**Methods** and **Fig. 2a**).

132 In addition, we included two sets of randomly selected metagenome samples of healthy

133 adults (n = 94) from the Human Microbiome Project (HMP)[31].

134

## A high-quality microbial genome catalog

136 Following quality control, we performed metagenomic assembly and binning on those

137 microbiome samples from three cohorts, yielding 7,769 MAGs. To evaluate the highest

138 quality representative genomes, we dereplicated the 7,769 MAGs at an average

139 nucleotide identity (ANI) threshold of 99%, resulting in a final set of 3,741 non-

140 redundant MAGs (nrMAGs) with strain-level resolution. The nrMAGs were contributed

141 by HMP, Verma-cohort, and the BIDMC-cohort in proportions of approximately 37%,

142 23%, and 40%, respectively **(Fig. 2b)**. In particular, our findings indicate that recipients

143 prior to FMT made a smaller contribution to the nrMAG collection compared to donors

144 and recipients after FMT, suggesting a reduced microbial diversity **(Fig. S1)**. These

145 nrMAGs exhibited a mean completeness of 88%, mean contamination of 0.93%, mean

146 genome size of 2.5 megabases (Mb), and mean N50 of 65.7 kilobases (kb) (**Fig. 1b-c**

147 and **Fig. 2c-f**). Out of the 3,741 strain-level nrMAGs, 1,390 (37.16%) nrMAGs met

148 medium-quality criteria (50% ≤ completeness < 90%, and ≤5% contamination), while

149 2,351 (62.84%) nrMAGs exhibited high-quality ( ≥ 90% completeness, and ≤ 5%

150 contamination)[32,33] (**Fig. 1b**).

151

152 Using the Genome Taxonomy Database[34], these nrMAGs were taxonomically assigned

153 to 17 phyla, 22 classes, 47 orders, 104 families, and 408 genera, spanning across 883

154 species. Most of them belonged to Firmicutes_A (60.97%), followed by Bacteroidetes

155 (10.18%), and Actinobacteria (10.00%). The phyla information of nrMAGs was

156 summarized in **Fig. 1b-c**. Among those 883 species, *Agathobacter rectalis*, *Blautia_A*

157 *wexlerae*, *Gemmiger formicilis*, *Fusicatenibacter saccharivorans*, and *Bifidobacterium*

158 *longum* were the top five species with the highest strain-level diversity (i.e., number of

159 nrMAGs identified within a specific species, **Fig. 2g**).

160

161 **Microbial diversity**

162 We initially investigated alpha diversity in the human microbiome at the nrMAG level.

163 Alpha diversity measures (i.e., Richness and Shannon index) were compared among

164 different groups from the same cohorts (**Fig. 2h-j**). No significant differences were found

165 between the two randomly selected sets from the HMP (**Fig. 2h** and **Fig. S2a**). In

166 accordance with the original study[30], we found that the Richness and Shannon indices

167 of the gut microbiome in the recipients of pre-FMT were significantly lower than those in

168 donors. After FMT, those recipients showed similar alpha diversity to donors (**Fig. 2i**

169 and **Fig. S2b**). In the BIDMC-cohort, we found that only the CDI group showed

170 significantly lower alpha diversity than the CON group (**Fig. 2j** and **Fig. S2c**).

171 Participants from the ASC group only showed a significantly lower number of identified

172 nrMAGs than CON group participants.

173

174 Principal coordinate analysis (PCoA) based on robust Aitchison distance, combined with

175 PERMANOVA (permutational multivariate analysis of variance, a statistical method

176 commonly used for testing the association between the microbiome and a covariate of

177 interest), revealed no significant difference in the gut microbial community structure at

178 the nrMAG-level between the two datasets from HMP (**Fig. 2k**). We found that the

179 microbiomes of the donor and the recipients from pre-and post-FMT were

180 compositionally distinct in the Verma-cohort ($P = 0.0001$, PERMANOVA **Fig. 2l**).

181 Consistent with our previous study using 16S rRNA gene sequencing data[29], the overall

182 microbial composition differed significantly among different groups in the BIDMC-cohort

183 ($P = 0.0001$, PERMANOVA **Fig. 2m**).

184

**Identify potential permissive and protective nrMAGs for the rational design of live biotherapeutics.**

To identify candidate strains for the development of microbiota-derived biotherapeutics , we applied the generalized microbe-phenotype triangulation (GMPT) method, moving beyond the standard association analysis[27]. The GMPT relies on the following core hypothesis: species that are differentially abundant in most pairwise phenotype-based comparisons and whose abundances display a strong negative (or positive) correlation with the abundance of the pathogen tend to be causal preventive (or permissive) species that directly inhibit (or promote) the growth of the pathogen[27]. Our GMPT analysis incorporated microbiome data from the BIDMC-cohort, donor data from the Verma-cohort[30], and one set from HMP. Since we have two sets of randomly selected metagenome samples of healthy adults from the HMP, we systematically included one set of HMP microbiome data at a time to cross-validate the results between the two datasets. Then, we conducted pairwise comparisons for six individual phenotype groups, which encompassed CDI, ASC, NCD, and CON from the BIDMC-cohort, donors from the Verma-cohort, and a dataset from HMP.

Applying this approach to the data with the first set of HMP data, 15 pairwise differential abundance analyses generated a total of 1,349 nrMAGs present in at least one pairwise comparison (**Table S6**). To explore the potential relationship between those candidate nrMAGs and CDI, we calculated Spearman correlation coefficients between the average relative abundances of nrMAGs and pragmatic severity scores in a continuum of non-CDI controls and *C. difficile* colonized and infected subjects (i.e., HMP healthy controls: 0; Donor from Verma et al.[30]: 1; CON: 2; NCD: 3; ASC: 4 and CDI: 5) in different phenotypes. Similarly, we identified a total of 1,390 nrMAGs present in at least one pairwise comparison with the second set of HMP data (**Table S7**). Among the protective nrMAGs between the two runs with HMP data, 80.77% (525/650) and 81.14% (525/647) of them were overlapped, respectively. We then computed the average rank between two runs based on the frequency (**Table S8)**. Among the top 40 potential protective nrMAGs, the dominant species were *Dorea formicigenerans*, *Oscillibacter welbionis*, *Faecalibacterium prausnitzii*, *GCA-900066135 sp900066135*, *Bariatricus comes*,

*Phocaeicola dorei*, *Anaerobutyricum hallii*, *Bacteroides ovatus*, *Blautia_A obeum*, *Mediterraneibacter faecis*, *Alistipes putredinis*, *Odoribacter splanchnicus*, *Streptococcus salivarius*, and *Dorea longicatena* (**Table 1**). Through a systematic review of literature, we found that most of our candidate strains have been reported to be protective from CDI or non-CDI antibiotic associated diarrhea at higher taxonomical levels (e.g., species and genus levels) across existing studies (**Table 1**). These findings support the validity of our methods.

**The protective strains play an important role in FMT.**

To further validate the potential role of the protective strains we identified from the GMPT pipeline, we systemically tracked the microbiome changes of the recipients who underwent FMT in the Verma-cohort. We aimed to investigate if those protective strains also play an important role in the success of FMT. Notably, the microbiome samples from the recipients in the Verma cohort were not included in our previous GMPT analysis.

First, we examined the gain and loss of microbial strains before and after FMT to assess the transfer and engraftment of the donor microbiome in the recipient. For donor, pre-FMT recipients, and post-FMT recipients, we identified 3,129, 2,093, and 3,054 nrMAGs, respectively. Notably, post-FMT recipients showed a loss of 33 nrMAGs (**Fig. 3a**), with the majority of the lost strains attributed to species such as *Anaeroglobus micronuciformis*, *Phascolarctobacterium faecium*, *Fusobacterium polymorphum*, and *Duodenibacillus sp003472385* from Firmicutes_C, Proteobacteria, and Fusobacteriota (**Fig. 3b**). On the contrary, all recipients exhibited a gain of 923 nrMAGs from their donors (**Fig. 3a**). The majority of these engrafted strains were taxonomically annotated to Actinobacteriota and Firmicutes_A, such as species like *Ruminococcus_D bicirculans*, *Faecalibacterium prausnitzii*, *Faecalibacterium prausnitzii_G*, *Agathobacter rectalis*, *Agathobacter faecis*, *Acetatifactor sp900066565*, *Bifidobacterium adolescentis*, and *Collinsella aerofaciens_G* (**Fig. 3c-d**).

246 For each donor-recipient pair, we then calculated the difference in their gut microbial
247 community structure before and after FMT using the robust Aitchison distance. Our
248 findings indicate that the distance between donors and recipients was significantly
249 reduced after FMT compared to the pre-FMT state (**Fig. 3e**). Additionally, calculating the
250 strain share rate for each donor-recipient pair before and after FMT revealed agreement
251 with our previous finding that recipients gained more strains, shared a greater number
252 of strains with the donor after FMT (**Fig. 3f**).

253

254 Changes in the microbiome induced by FMT not only indicate the transfer and
255 engraftment of the donor microbiome but also involve alterations in the abundance of
256 coexisting strains. To address this question, we conducted the differential abundance
257 analysis among three groups. Consistent with the robust Aitchison distance and strain
258 share rate analyses, we only identified less differential abundant strains between donor
259 and post-FMT recipients (**Fig. 3g and Table S2**). We have identified 223 and 238
260 differential abundant nrMAGs from the comparison of donor vs. pre-FMT recipients (**Fig.**
261 **3g and Table S3**) and pre-FMT recipients vs. post-FMT recipients (**Fig. 3g and Table**
262 **S4**), respectively. Among these differential abundant nrMAGs, we found 179 overlapped
263 strains (**Fig. 3h, Table S5**), including strains from *Blautia_A wexlerae*, *Veillonella*
264 *parvula_A*, *Veillonella parvula*, *Blautia_A sp900066165*, *Escherichia coli_D*, *Escherichia*
265 *flexneri*, *Anaeroglobus micronuciformis*, *Blautia_A obeum*, *Lacticaseibacillus rhamnosus*,
266 and *Veillonella dispar_A.* Specifically, we observed significant increases in some
267 candidate protective strains following FMT. These include multiple strains from *Dorea*
268 *formicigenerans*, *Mediterraneibacter faecis*, *Phocaeicola dorei*, *Blautia_A wexlerae*, and
269 *Blautia_A obeum*. This finding further validated the potential role of protective strains in
270 treating CDI.

271

## Discussion

273 The growing interest in FMT as a therapeutic approach stems from its high success rate
274 in treating recurrent CDI, leading to an exploration of its potential for addressing various
275 human diseases[35]. However, FMT remains an unstandardized procedure with unclear
276 mechanisms and long-term safety concerns[35,36]. Therefore, an advantage of microbial

277 consortia over "whole stool" FMT is the introduction of a group of specific microbiota

278 that can precisely target and effectively treat a disease while minimizing clinical risks. In

279 this study, we used a computational pipeline to directly identify candidate bacterial

280 strains from a diverse CDI-related metagenomic dataset, thereby facilitating the

281 targeted development of microbial therapies and advancing our understanding of CDI

282 pathogenesis and treatment.

283

284 By tracking the dynamic changes in gut microbiome data undergoing FMT, we identified

285 significant shifts in the microbial structure of the recipients. Although we did not utilize

286 the microbiome data from recipients before and after FMT in our GMPT pipeline, we

287 found that some of the top ranked candidate protective strains showed significant

288 increases after FMT, including multiple strains from *Dorea formicigenerans*,

289 *Mediterraneibacter faecis*, *Phocaeicola dorei*, and *Blautia_A obeum*. This finding

290 provides an additional layer of validity to our method. In addition, we performed a

291 systematic literature review on the highest ranked candidate protective strains and

292 found that the majority of them have been reported to have various protective roles at

293 species or genus levels in the CDI continuum: negative association with *C. difficile*

294 colonization, infection and severity. We found clustering of the main protective species

295 within the families Lachnospiraceae, Bacteroidaceae, and Oscillospiraceae. For

296 example, *F. prausnitzii*, a beneficial human gut microbe touted as a candidate for next-

297 generation probiotics[37], was found to have reduced abundance in CDI patients, which

298 was restored after FMT[38]. Interestingly, we have also identified a protective strain of the

299 species *Dorea longicatena,* which is a component of a defined bacterial consortium

300 (VE303) with encouraging Phase 2 clinical data, consisting of eight, nonpathogenic,

301 nontoxigenic, commensal strains of Clostridia[39].

302

303 In addition to the potential protective strains, we also identified multiple permissive

304 strains of *C. difficile*, including strains from *Enterococcus_B faecium* and *Eggerthella*

305 *lenta*. This aligns with a previous study reporting that enterococci (including *E. faecium*)

306 can enhance the fitness and pathogenicity of *C. difficile* via shaping the metabolic

307 environment in the gut and reprograming *C. difficile* metabolism[40]. Additionally, *E. lenta*

308 is an anaerobic gram-positive bacillus associated with polymicrobial intraabdominal

309 infections[41]. Therefore, the potentially permissive strains that we identified from this

310 study offer the opportunity to further understand how *C. difficile* interacts with the rich

311 community of microorganisms in the colon. Moreover, in alignment with the variation in

312 biological properties among bacterial strains within the same species[42], we have

313 observed distinct roles played by different strains of *F. prausnitzii_D* in the context of

314 CDI. This underscores the critical importance of conducting studies at the strain level.

315

316 The current study has some limitations. First, we leveraged metagenomic data from

317 three independent datasets with technological variations, including differences in

318 sequencing depth. Second, we did not pre-define a strict threshold to select potential

319 protective strains from the candidate list for further experimental validations. Lastly, the

320 inference of the efficacy of candidate protective strains against CDI is limited by the

321 current computational algorithm. To test the efficacy of our proposed microbial consortia

322 and gain a deeper understanding of exact mechanisms, the utilization of techniques of

323 metabolomics and immunological approaches, along with direct *in vitro* and *in vivo*

324 experiments, are necessary.

325

326 Taken together, our results provide compelling evidence for the rational design of

327 microbial consortia against *C. difficile*. Many of the candidates detected here replicate

328 previously reported findings, supporting the validity of our results. Importantly, our work

329 paves the way for the design of LBPs against general microbiome-related diseases.

330

331 **Methods**

332 **Study cohorts**

333 **Dataset I: BIDMC-cohort.** The background and design of this cohort have been

334 detailed in our previous studies[28,29]. This clinical cohort consists of 104 well-

335 characterized recruited participants, who were divided into four groups associated with

336 different *C. difficile* infection/colonization statuses: (1) *C. difficile* infection (CDI, n=47):

337 Eligible patients were inpatients ≥ 18 years old with new-onset diarrhea, positive clinical

338 stool NAAT (Xpert *C. difficile*/Epi) result, and a decision to treat for CDI; (2)

11

339    Asymptomatic Carriage (ASC, n = 17): Eligible patients were inpatients ≥ 18 years old,

340    admitted for at least 72 hours, who had received at least one dose of an antibiotic within

341    the past seven days, and did not have diarrhea in the 48 hours prior to stool specimen

342    submission, and positive clinical stool NAAT result; (3) Non-CDI Diarrhea (NCD, n = 14):

343    patients with diarrhea (confirmed using the same definition used for the CDI cohort) but

344    had NAAT-negative stool on clinical C. difficile testing; and (4) Control (CON, n = 26):

345    patients without diarrhea who had screened as eligible for the ASC cohort but were

346    NAAT-negative on research stool testing. DNA of fecal samples (200 mg) were

347    extracted using Mag-Bind® Universal Metagenomics Kit (Product# M5633-01, Omega

348    Biotek) and DNeasy PowerSoil Kit (Catalog# 12888-100, Qiagen) according to

349    manufacturer's instructions. The quality of the extracted DNA was measured by 1%

350    agarose gel electrophoresis and Qubit® 3.0 Fluorometer (ThermoFisher). Subsequently,

351    the extracted DNAs were used for shotgun metagenomic library construction, and

352    sequencing was performed on the Illumina HiSeq X Ten platform, generating a 150 bp

353    paired-end library for each sample.

354    **Dataset II: Verma-cohort.** In the study conducted by Verma et al[30], fecal samples were

355    collected from 22 patients with recurrent CDI before and after FMT and their

356    corresponding healthy donors (n=21, with one donor providing fecal samples for two

357    different recipients). Eight-seven WMS human gut metagenomes were downloaded

358    from this study via NCBI Sequence Read Archive (BioProject ID PRJNA705895). The

359    clinical outcome in recurrent CDI patients after FMT was determined by the

360    symptomatic resolution of CDI[30]. Clinical symptoms such as diarrhea, bloating,

361    abdominal pain, and cramping were alleviated in all patients within 3–7 days following

362    FMT[30].

363    **Dataset III:** Human Microbiome Project. Human gut metagenomes (Ninety-eight

364    individuals) were randomly selected from HMP data (https://portal.hmpdacc.org/)[43]. All

365    samples are from the HMP study[31] and are healthy adult subjects. In total, ninety-four

366    human gut metagenomes were randomly selected based on the largest group size in

367    our clinical cohort. To cross-validate the main findings, we randomly divided the HMP

368    data into two sets in the downstream analyses.

369

**Metagenome assembly and binning**

Genome reconstruction of the human microbiome using metagenomic sequencing data was executed through the functional modules of metaWRAP (v1.3.2)[44]. All metagenomic sequencing data underwent quality control and removal of human contamination using metaWRAP-Read_qc. Clean reads were then assembled with the metaWRAP-Assembly module using metaSPAdes (v3.13.0)[45]. The assembled contigs were binned into bins using three metagenomic binning tools: MetaBAT (v2.12.1)[46], MaxBin (v2.2.6)[47], and CONCOCT (v1.0.0)[48]. The default minimum length of contigs used for constructing bins with MaxBin2 and CONCOCT was 1000 bp, and metaBAT2 was defaulted to 1500 bp[44]. The bins from each binning tool were integrated and refined with Bin_refinement module of metaWRAP with options "-c 50 -x 10", corresponding to the criterion of medium-quality draft MAGs[32]. CheckM (v1.0.12)[49] was used to estimate the completeness and contamination of the bins, and the minimum completion and maximum contamination were 50% and 10%, respectively.

**De-replication of MAGs and genome annotation.** All 7,776 MAGs underwent de-replication into non-redundant MAGs (nrMAGs) using dRep (v3.0.0) (≥50% genome completeness and ≤5% contamination)[50]. Initially, MAGs from three cohorts were divided into primary clusters using Mash at a 90% Mash ANI. Then, each primary cluster was used to form secondary clusters at the threshold of 99% ANI with at least 30% overlap between genomes[51]. Taxonomic annotation of all nrMAGs was conducted using GTDB-Tk (v.1.4.1)[52] based on the Genome Taxonomy Database (http://gtdb.ecogenomic.org/)[34], providing standardized taxonomic labels for subsequent analysis in this study.

**Abundance estimation and phylogenetic analysis of nrMAGs**

The metaWRAP-Quant_bins module coupled with Salmon (v0.13.1)[53] was employed to access the abundance of each nrMAGs in each metagenomic sample. The phylogenetic tree of the nrMAGs was constructed using PhyloPhlAn (v3.0.58)[54] and visualized through iTOL (https://itol.embl.de/)[55].

**Statistical analysis**

Microbial alpha diversity measures were calculated at the nrMAGs level using R vegan package (v2.5.7), and principal coordinates analysis (PCoA) plots were generated using robust Aitchison distance[56]. Differences in microbiome compositions across different groups were tested by the permutational multivariate analysis of variance (PERMANOVA) using the "adonis" function in R vegan package. All PERMANOVA tests were performed with 9999 permutations based on the robust Aitchison distance. We defined strain-sharing rates as the total number of shared strains between two samples divided by the number of common species identified from the two samples. Differences between the groups were analyzed using a Wilcoxon–Mann–Whitney test. For differential abundance analysis and GMPT (Generalized Microbe Phenotype Triangulation) pipeline[27], we used ANCOM (analysis of composition of microbiomes)[57], with a Benjamini–Hochberg correction at a 5% level of significance. All statistical analysis was performed with R (version 3.6.3).

**Data availability**

Metagenomic data from HMP are available via https://portal.hmpdacc.org. The metagenomic data from the study of Verma et al.[30] can be downloaded via NCBI Sequence Read Archive (BioProject ID PRJNA705895). The metagenomic data from the BIDMC-cohort is available in the NCBI Bioproject under accession code PRJNA1067975. Metagenome-assembled genomes for all samples are available on Figshare (https://doi.org/10.6084/m9.figshare.25355857).

**Code availability**

The code for the construction of the MAGs catalog and statistical analysis and visualization is available in the GitHub repository (https://github.com/ShanlinKe/CDI).

**Acknowledgements**

## Author contributions

Y.-Y.L., X.C. and C.P.K. conceived and designed the project. N.R.P. and C.P.K. planned and performed the human studies and sample collections. S.K. performed all the data analysis. S.K., J.A.V.G., and Y.-Y.L. interpreted the results and prepared the manuscript. Z.S., X.C., and C.P.K. interpreted the results, reviewed and edited the manuscript. Y.C. acquired the raw sequencing data, reviewed and edited the manuscript. All authors have read and approved the manuscript.

## Competing interests

C.P.K. has acted as a paid consultant to: Artugen Therapeutics, Facile Therapeutics, Finch, Fzata, Glaxo Smith Kline, Immunics Therapeutics, Recursion Pharmaceuticals, RVAC Medicines, Sanofi Pasteur, Seres Therapeutics, and Summit Therapeutics; C.P.K. has acted as a paid consultant and member of the Scientific Advisory Board to: Acurx Pharmaceuticals, Anokion, Ferring Pharma, Inova Diagnostics, Janssen Pharmaceuticals, Merck & Company, Milky Way Life Sciences, Pfizer, Takeda, and Vedanta Biosciences; C.P.K. has acted as an unpaid consultant and had private equity in Glutenostics; has acted as a paid consultant and has private equity in Cour Pharmaceuticals, and First Light Biosciences, Inc; and has acted as a study investigator for: Milky Way Life Sciences and Merck. X.C. is a paid consultant and board member of Milky Way Life Sciences and has acted as a consultant to Artugen Therapeutics and RVAC Medicines. X.C. is a co-founder with private equity in Milky Way Life Sciences

465

466

467

468

469

470

471 **References**

472 1    Magill, S. S. *et al.* Changes in Prevalence of Health Care-Associated Infections in U.S.
473    Hospitals. *N Engl J Med* **379**, 1732-1744, doi:10.1056/NEJMoa1801550 (2018).
474 2    Sehgal, K. & Khanna, S. Gut microbiome and Clostridioides difficile infection: a closer
475    look at the microscopic interface. *Therap Adv Gastroenterol* **14**, 1756284821994736,
476    doi:10.1177/1756284821994736 (2021).
477 3    Rupnik, M., Wilcox, M. H. & Gerding, D. N. Clostridium difficile infection: new
478    developments in epidemiology and pathogenesis. *Nat Rev Microbiol* **7**, 526-536,
479    doi:10.1038/nrmicro2164 (2009).
480 4    Nelson, R. L., Suda, K. J. & Evans, C. T. Antibiotic treatment for Clostridium difficile-
481    associated diarrhoea in adults. *Cochrane Database Syst Rev* **3**, CD004610,
482    doi:10.1002/14651858.CD004610.pub5 (2017).
483 5    Ooijevaar, R. E. *et al.* Update of treatment algorithms for Clostridium difficile infection.
484    *Clin Microbiol Infect* **24**, 452-462, doi:10.1016/j.cmi.2017.12.022 (2018).
485 6    Rineh, A., Kelso, M. J., Vatansever, F., Tegos, G. P. & Hamblin, M. R. Clostridium difficile
486    infection: molecular pathogenesis and novel therapeutics. *Expert Rev Anti Infect Ther* **12**,
487    131-150, doi:10.1586/14787210.2014.866515 (2014).
488 7    Bagdasarian, N., Rao, K. & Malani, P. N. Diagnosis and treatment of Clostridium difficile
489    in adults: a systematic review. *JAMA* **313**, 398-408, doi:10.1001/jama.2014.17103 (2015).
490 8    Yokoyama, Y. *et al.* Risk factors of first recurrence of Clostridioides difficile infection.
491    *Anaerobe* **75**, 102556, doi:10.1016/j.anaerobe.2022.102556 (2022).
492 9    Ducarmon, Q. R. *et al.* Gut Microbiota and Colonization Resistance against Bacterial
493    Enteric Infection. *Microbiol Mol Biol Rev* **83**, doi:10.1128/MMBR.00007-19 (2019).
494 10    Laffin, M., Millan, B. & Madsen, K. L. Fecal microbial transplantation as a therapeutic
495    option in patients colonized with antibiotic resistant organisms. *Gut Microbes* **8**, 221-
496    224, doi:10.1080/19490976.2016.1278105 (2017).
497 11    Woodworth, M. H., Hayden, M. K., Young, V. B. & Kwon, J. H. The Role of Fecal
498    Microbiota Transplantation in Reducing Intestinal Colonization With Antibiotic-Resistant

16

499     Organisms: The Current Landscape and Future Directions. *Open Forum Infect Dis* **6**,
500     doi:10.1093/ofid/ofz288 (2019).
501  12 Segal, J. P. *et al.* Mechanisms underpinning the efficacy of faecal microbiota
502     transplantation in treating gastrointestinal disease. *Therap Adv Gastroenterol* **13**,
503     1756284820946904, doi:10.1177/1756284820946904 (2020).
504  13 DeFilipp, Z. *et al.* Drug-Resistant E. coli Bacteremia Transmitted by Fecal Microbiota
505     Transplant. *N Engl J Med* **381**, 2043-2050, doi:10.1056/NEJMoa1910437 (2019).
506  14 Blaser, M. J. Fecal Microbiota Transplantation for Dysbiosis - Predictable Risks. *N Engl J*
507     *Med* **381**, 2064-2066, doi:10.1056/NEJMe1913807 (2019).
508  15 Administration, T. U. S. F. a. D. FDA Approves First Fecal Microbiota Product.  (2022).
509  16 Administration, T. U. S. F. a. D. FDA Approves First Orally Administered Fecal Microbiota
510     Product for the Prevention of Recurrence of Clostridioides difficile Infection.  (2023).
511  17 Dubberke, E. R. *et al.* Results From a Randomized, Placebo-Controlled Clinical Trial of a
512     RBX2660-A Microbiota-Based Drug for the Prevention of Recurrent Clostridium difficile
513     Infection. *Clin Infect Dis* **67**, 1198-1204, doi:10.1093/cid/ciy259 (2018).
514  18 Peery, A. F. *et al.* AGA Clinical Practice Guideline on Fecal Microbiota-Based Therapies
515     for Select Gastrointestinal Diseases. *Gastroenterology* **166**, 409-434,
516     doi:10.1053/j.gastro.2024.01.008 (2024).
517  19 Walter, J. & Shanahan, F. Fecal microbiota-based treatment for recurrent Clostridioides
518     difficile infection. *Cell* **186**, 1087, doi:10.1016/j.cell.2023.02.034 (2023).
519  20 Liao, H., Ji, Y. & Sun, Y. High-resolution strain-level microbiome composition analysis
520     from short reads. *Microbiome* **11**, 183, doi:10.1186/s40168-023-01615-w (2023).
521  21 Besser, M. D. R. E., Griffin, M. D. P. M. & Slutsker, M. D. M. P. H. L. Escherichia
522     coliO157:H7 Gastroenteritis and the Hemolytic Uremic Syndrome: An Emerging
523     Infectious Disease. *Annual Review of Medicine* **50**, 355-367,
524     doi:10.1146/annurev.med.50.1.355 (1999).
525  22 Sonnenborn, U. Escherichia coli strain Nissle 1917-from bench to bedside and back:
526     history of a special Escherichia coli strain with probiotic properties. *FEMS Microbiol Lett*
527     **363**, doi:10.1093/femsle/fnw212 (2016).
528  23 Arribas, B. *et al.* A probiotic strain of Escherichia coli, Nissle 1917, given orally exerts
529     local and systemic anti-inflammatory effects in lipopolysaccharide-induced sepsis in
530     mice. *Br J Pharmacol* **157**, 1024-1033, doi:10.1111/j.1476-5381.2009.00270.x (2009).
531  24 Tu, Q., He, Z. & Zhou, J. Strain/species identification in metagenomes using genome-
532     specific markers. *Nucleic Acids Res* **42**, e67, doi:10.1093/nar/gku138 (2014).
533  25 Truong, D. T., Tett, A., Pasolli, E., Huttenhower, C. & Segata, N. Microbial strain-level
534     population structure and genetic diversity from metagenomes. *Genome Res* **27**, 626-638,
535     doi:10.1101/gr.216242.116 (2017).
536  26 Ke, S., Weiss, S. T. & Liu, Y. Y. Dissecting the role of the human microbiome in COVID-19
537     via metagenome-assembled genomes. *Nat Commun* **13**, 5235, doi:10.1038/s41467-022-
538     32991-w (2022).
539  27 Ke, S. *et al.* A computational method to dissect colonization resistance of the gut
540     microbiota against pathogens. *Cell Reports Methods*, 100576,
541     doi:https://doi.org/10.1016/j.crmeth.2023.100576 (2023).

542   28   Pollock, N. R. *et al.* Comparison of Clostridioides difficile Stool Toxin Concentrations in
543        Adults With Symptomatic Infection and Asymptomatic Carriage Using an Ultrasensitive
544        Quantitative Immunoassay. *Clin Infect Dis* **68**, 78-86, doi:10.1093/cid/ciy415 (2019).
545   29   Ke, S. *et al.* Integrating gut microbiome and host immune markers to understand the
546        pathogenesis of Clostridioides difficile infection. *Gut Microbes* **13**, 1-18,
547        doi:10.1080/19490976.2021.1935186 (2021).
548   30   Verma, S. *et al.* Identification and engraftment of new bacterial strains by shotgun
549        metagenomic sequence analysis in patients with recurrent Clostridioides difficile
550        infection before and after fecal microbiota transplantation and in healthy human
551        subjects. *PLoS One* **16**, e0251590, doi:10.1371/journal.pone.0251590 (2021).
552   31   Human Microbiome Project, C. Structure, function and diversity of the healthy human
553        microbiome. *Nature* **486**, 207-214, doi:10.1038/nature11234 (2012).
554   32   Bowers, R. M. *et al.* Minimum information about a single amplified genome (MISAG) and
555        a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nat Biotechnol* **35**,
556        725-731, doi:10.1038/nbt.3893 (2017).
557   33   Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes
558        substantially expands the tree of life. *Nat Microbiol* **2**, 1533-1542, doi:10.1038/s41564-
559        017-0012-7 (2017).
560   34   Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
561        substantially revises the tree of life. *Nat Biotechnol* **36**, 996-1004, doi:10.1038/nbt.4229
562        (2018).
563   35   Wortelboer, K., Nieuwdorp, M. & Herrema, H. Fecal microbiota transplantation beyond
564        Clostridioides difficile infections. *EBioMedicine* **44**, 716-729,
565        doi:10.1016/j.ebiom.2019.05.066 (2019).
566   36   El-Matary, W. Fecal microbiota transplantation: long-term safety issues. *Am J
567        Gastroenterol* **108**, 1537-1538, doi:10.1038/ajg.2013.208 (2013).
568   37   Bai, Z. *et al.* Comprehensive analysis of 84 Faecalibacterium prausnitzii strains uncovers
569        their genetic diversity, functional characteristics, and potential risks. *Front Cell Infect
570        Microbiol* **12**, 919701, doi:10.3389/fcimb.2022.919701 (2022).
571   38   Bjorkqvist, O. *et al.* Faecalibacterium prausnitzii increases following fecal microbiota
572        transplantation in recurrent Clostridioides difficile infection. *PLoS One* **16**, e0249861,
573        doi:10.1371/journal.pone.0249861 (2021).
574   39   Louie, T. *et al.* VE303, a Defined Bacterial Consortium, for Prevention of Recurrent
575        Clostridioides difficile Infection: A Randomized Clinical Trial. *JAMA* **329**, 1356-1366,
576        doi:10.1001/jama.2023.4314 (2023).
577   40   Smith, A. B. *et al.* Enterococci enhance Clostridioides difficile pathogenesis. *Nature* **611**,
578        780-786, doi:10.1038/s41586-022-05438-x (2022).
579   41   Ugarte-Torres, A., Gillrie, M. R., Griener, T. P. & Church, D. L. Eggerthella lenta
580        Bloodstream Infections Are Associated With Increased Mortality Following Empiric
581        Piperacillin-Tazobactam (TZP) Monotherapy: A Population-based Cohort Study. *Clin
582        Infect Dis* **67**, 221-228, doi:10.1093/cid/ciy057 (2018).
583   42   Hilau, S., Katz, S., Wasserman, T., Hershberg, R. & Savir, Y. Density-dependent effects
584        are the main determinants of variation in growth dynamics between closely related

585     bacterial strains. *PLoS Comput Biol* **18**, e1010565, doi:10.1371/journal.pcbi.1010565
586     (2022).
587  43  Gevers, D. *et al.* The Human Microbiome Project: a community resource for the healthy
588     human microbiome. *PLoS Biol* **10**, e1001377, doi:10.1371/journal.pbio.1001377 (2012).
589  44  Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP-a flexible pipeline for genome-
590     resolved metagenomic data analysis. *Microbiome* **6**, 158, doi:10.1186/s40168-018-0541-
591     1 (2018).
592  45  Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile
593     metagenomic assembler. *Genome Res* **27**, 824-834, doi:10.1101/gr.213959.116 (2017).
594  46  Wu, Y. W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning algorithm
595     to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605-607,
596     doi:10.1093/bioinformatics/btv638 (2016).
597  47  Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately
598     reconstructing single genomes from complex microbial communities. *PeerJ* **3**, e1165,
599     doi:10.7717/peerj.1165 (2015).
600  48  Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat*
601     *Methods* **11**, 1144-1146, doi:10.1038/nmeth.3103 (2014).
602  49  Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
603     assessing the quality of microbial genomes recovered from isolates, single cells, and
604     metagenomes. *Genome Res* **25**, 1043-1055, doi:10.1101/gr.186072.114 (2015).
605  50  Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate
606     genomic comparisons that enables improved genome recovery from metagenomes
607     through de-replication. *ISME J* **11**, 2864-2868, doi:10.1038/ismej.2017.126 (2017).
608  51  Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using
609     MinHash. *Genome Biol* **17**, 132, doi:10.1186/s13059-016-0997-x (2016).
610  52  Chaumeil, P. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify
611     genomes with the Genome Taxonomy Database. *Bioinformatics*,
612     doi:10.1093/bioinformatics/btz848 (2019).
613  53  Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and
614     bias-aware quantification of transcript expression. *Nat Methods* **14**, 417-419,
615     doi:10.1038/nmeth.4197 (2017).
616  54  Asnicar, F. *et al.* Precise phylogenetic analysis of microbial isolates and genomes from
617     metagenomes using PhyloPhlAn 3.0. *Nat Commun* **11**, 2500, doi:10.1038/s41467-020-
618     16366-7 (2020).
619  55  Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic
620     tree display and annotation. *Nucleic Acids Res* **49**, W293-W296,
621     doi:10.1093/nar/gkab301 (2021).
622  56  Martino, C. *et al.* A Novel Sparse Compositional Technique Reveals Microbial
623     Perturbations. *mSystems* **4**, doi:10.1128/mSystems.00016-19 (2019).
624  57  Mandal, S. *et al.* Analysis of composition of microbiomes: a novel method for studying
625     microbial composition. *Microb Ecol Health Dis* **26**, 27663, doi:10.3402/mehd.v26.27663
626     (2015).

627    58    Crobach, M. J. T. *et al.* The Bacterial Gut Microbiota of Adult Patients Infected, Colonized
628          or Noncolonized by Clostridioides difficile. *Microorganisms* **8**,
629          doi:10.3390/microorganisms8050677 (2020).
630    59    Mullish, B. H. *et al.* Microbial bile salt hydrolases mediate the efficacy of faecal
631          microbiota transplant in the treatment of recurrent Clostridioides difficile infection. *Gut*,
632          doi:10.1136/gutjnl-2018-317842 (2019).
633    60    Berkell, M. *et al.* Microbiota-based markers predictive of development of Clostridioides
634          difficile infection. *Nat Commun* **12**, 2241, doi:10.1038/s41467-021-22302-0 (2021).
635    61    Milani, C. *et al.* Gut microbiota composition and Clostridium difficile infection in
636          hospitalized elderly individuals: a metagenomic study. *Sci Rep* **6**, 25945,
637          doi:10.1038/srep25945 (2016).
638    62    Han, S. H., Yi, J., Kim, J. H., Lee, S. & Moon, H. W. Composition of gut microbiota in
639          patients with toxigenic Clostridioides (Clostridium) difficile: Comparison between
640          subgroups according to clinical criteria and toxin gene load. *PLoS One* **14**, e0212626,
641          doi:10.1371/journal.pone.0212626 (2019).
642    63    Feuerstadt, P. *et al.* SER-109, an Oral Microbiome Therapy for Recurrent Clostridioides
643          difficile Infection. *N Engl J Med* **386**, 220-229, doi:10.1056/NEJMoa2106516 (2022).
644    64    Tvede, M. & Rask-Madsen, J. Bacteriotherapy for chronic relapsing Clostridium difficile
645          diarrhoea in six patients. *Lancet* **1**, 1156-1160, doi:10.1016/s0140-6736(89)92749-9
646          (1989).
647    65    Yoon, S. *et al.* Bile salt hydrolase-mediated inhibitory effect of Bacteroides ovatus on
648          growth of Clostridium difficile. *J Microbiol* **55**, 892-899, doi:10.1007/s12275-017-7340-4
649          (2017).
650    66    Amrane, S. *et al.* Metagenomic and culturomic analysis of gut microbiota dysbiosis
651          during Clostridium difficile infection. *Sci Rep* **9**, 12807, doi:10.1038/s41598-019-49189-8
652          (2019).
653    67    Hourigan, S. K. *et al.* Fecal Transplant in Children With Clostridioides difficile Gives
654          Sustained Reduction in Antimicrobial Resistance and Potential Pathogen Burden. *Open
655          Forum Infect Dis* **6**, ofz379, doi:10.1093/ofid/ofz379 (2019).
656    68    Douchant, K. *et al.* Defined microbial communities and their soluble products protect
657          mice from Clostridioides difficile infection. *Commun Biol* **7**, 135, doi:10.1038/s42003-
658          024-05778-6 (2024).
659    69    Francisco, D. M. A. *et al.* Risk Factors Associated with Severe Clostridioides difficile
660          Infection in Patients with Cancer. *Infect Dis Ther* **12**, 209-225, doi:10.1007/s40121-022-
661          00722-9 (2023).
662    70    Moelling, K. & Broecker, F. Fecal microbiota transplantation to fight Clostridium difficile
663          infections and other intestinal diseases. *Bacteriophage* **6**, e1251380,
664          doi:10.1080/21597081.2016.1251380 (2016).
665    71    Roychowdhury, S. *et al.* Faecalibacterium prausnitzii and a Prebiotic Protect Intestinal
666          Health in a Mouse Model of Antibiotic and Clostridium difficile Exposure. *JPEN J
667          Parenter Enteral Nutr* **42**, 1156-1167, doi:10.1002/jpen.1053 (2018).
668    72    Vakili, B., Fateh, A., Asadzadeh Aghdaei, H., Sotoodehnejadnematalahi, F. & Siadat, S. D.
669          Intestinal Microbiota in Elderly Inpatients with Clostridioides difficile Infection. *Infect
670          Drug Resist* **13**, 2723-2731, doi:10.2147/IDR.S262019 (2020).

671   73    Vakili, B., Fateh, A., Asadzadeh Aghdaei, H., Sotoodehnejadnematalahi, F. & Siadat, S. D.
672             Characterization of Gut Microbiota in Hospitalized Patients with Clostridioides difficile
673             Infection. *Curr Microbiol* **77**, 1673-1680, doi:10.1007/s00284-020-01980-x (2020).
674   74    Shoaei, P. *et al.* Gut microbiota in burned patients with Clostridioides difficile infection.
675             *Burns* **48**, 1120-1129, doi:10.1016/j.burns.2021.11.023 (2022).
676   75    Gu, X. *et al.* Gut Ruminococcaceae levels at baseline correlate with risk of antibiotic-
677             associated diarrhea. *iScience* **25**, 103644, doi:10.1016/j.isci.2021.103644 (2022).
678   76    Dudik, B., Kinova Sepova, H., Greifova, G., Bilka, F. & Bilkova, A. Next generation
679             probiotics: an overview of the most promising candidates. *Epidemiol Mikrobiol Imunol*
680             **71**, 48-56 (2022).
681   77    Lee, Y. J., Yu, W. K. & Heo, T. R. Identification and screening for antimicrobial activity
682             against Clostridium difficile of Bifidobacterium and Lactobacillus species isolated from
683             healthy infant faeces. *Int J Antimicrob Agents* **21**, 340-346, doi:10.1016/s0924-
684             8579(02)00389-8 (2003).
685   78    Solbach, P. *et al.* Microbiota-associated Risk Factors for Clostridioides difficile
686             Acquisition in Hospitalized Patients: A Prospective, Multicentric Study. *Clin Infect Dis* **73**,
687             e2625-e2634, doi:10.1093/cid/ciaa871 (2021).
688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715
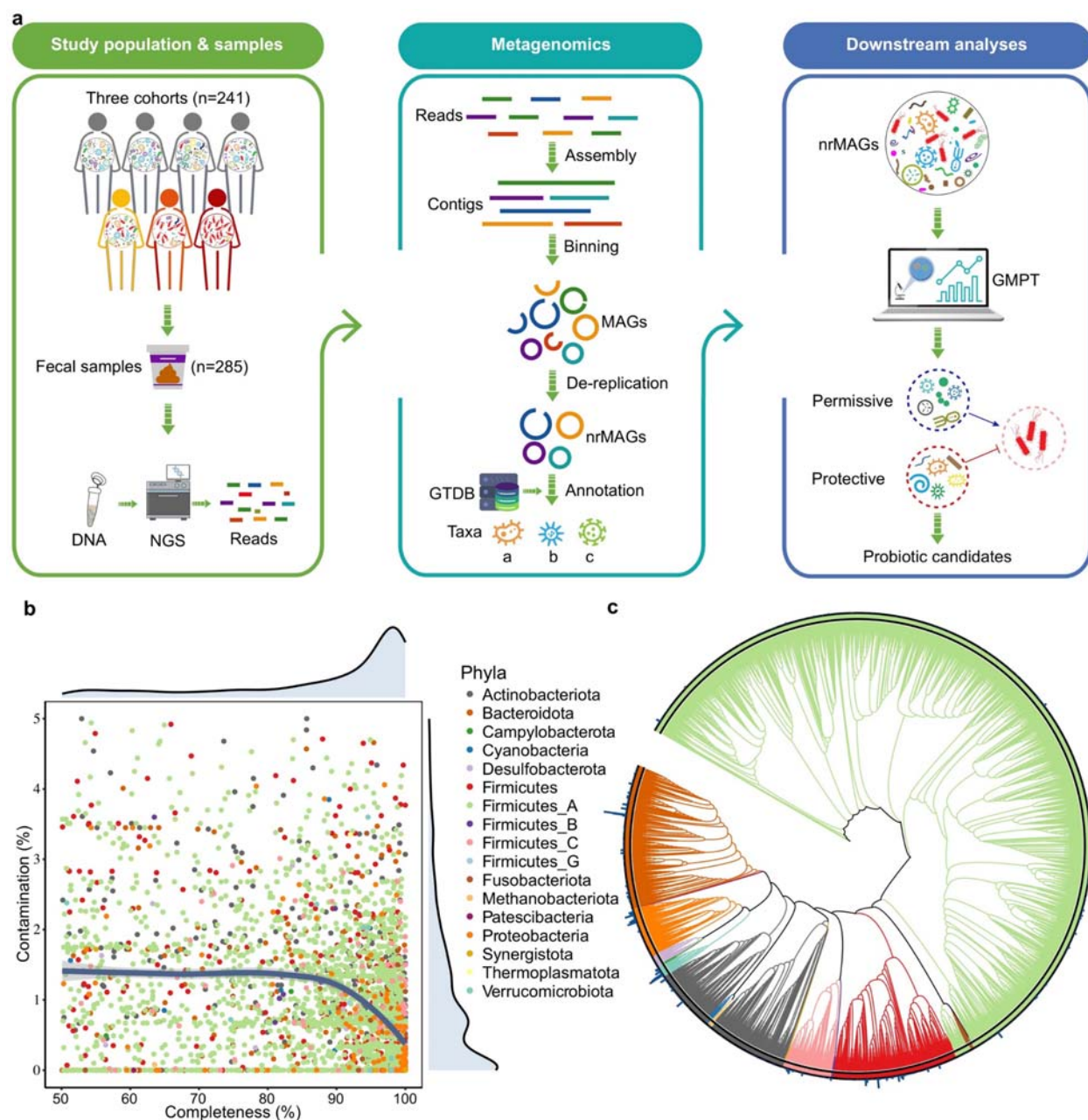
716 **Figures and Legends**

717

**Fig. 1 Study workflow and the reconstruction of the microbial genome catalog. a.** To rationally design microbial consortia against *C. difficile*, we sought to infer species that may potentially inhibit *C. difficile* from various metagenomic data. We collected a total of 285 shotgun metagenomic sequencing data from three independent cohorts. A total of 7,769 MAGs (≥50% completeness and ≤5% contamination) were constructed from all metagenomic sequencing data. The MAGs were then dereplicated to 3,741 non-redundant MAGs (nrMAGs, strain level) based on 99% of ANI. The taxonomy annotation and abundance estimation of nrMAGs were then conducted. We then

727 applied the generalized microbe-phenotype triangulation (GMPT) method to identify

728 candidate strains for the development of microbiota probiotics. **b.** The distribution of

729 completeness and contamination of nrMAGs is depicted, with the color of each point

730 representing the respective phylum. Additionally, the size of each point corresponds to

731 the genome size of the nrMAGs. **c.** A phylogenetic tree of nrMAGs was constructed

732 using PhyloPhlAn. In this representation, the color of the outer cycle and clades

733 signifies the phylum, while the bar plot within the cycle illustrates the average

734 abundance across all microbiome samples.

735

736

737

738

739

740

741

742

743

744

745

746

747

748
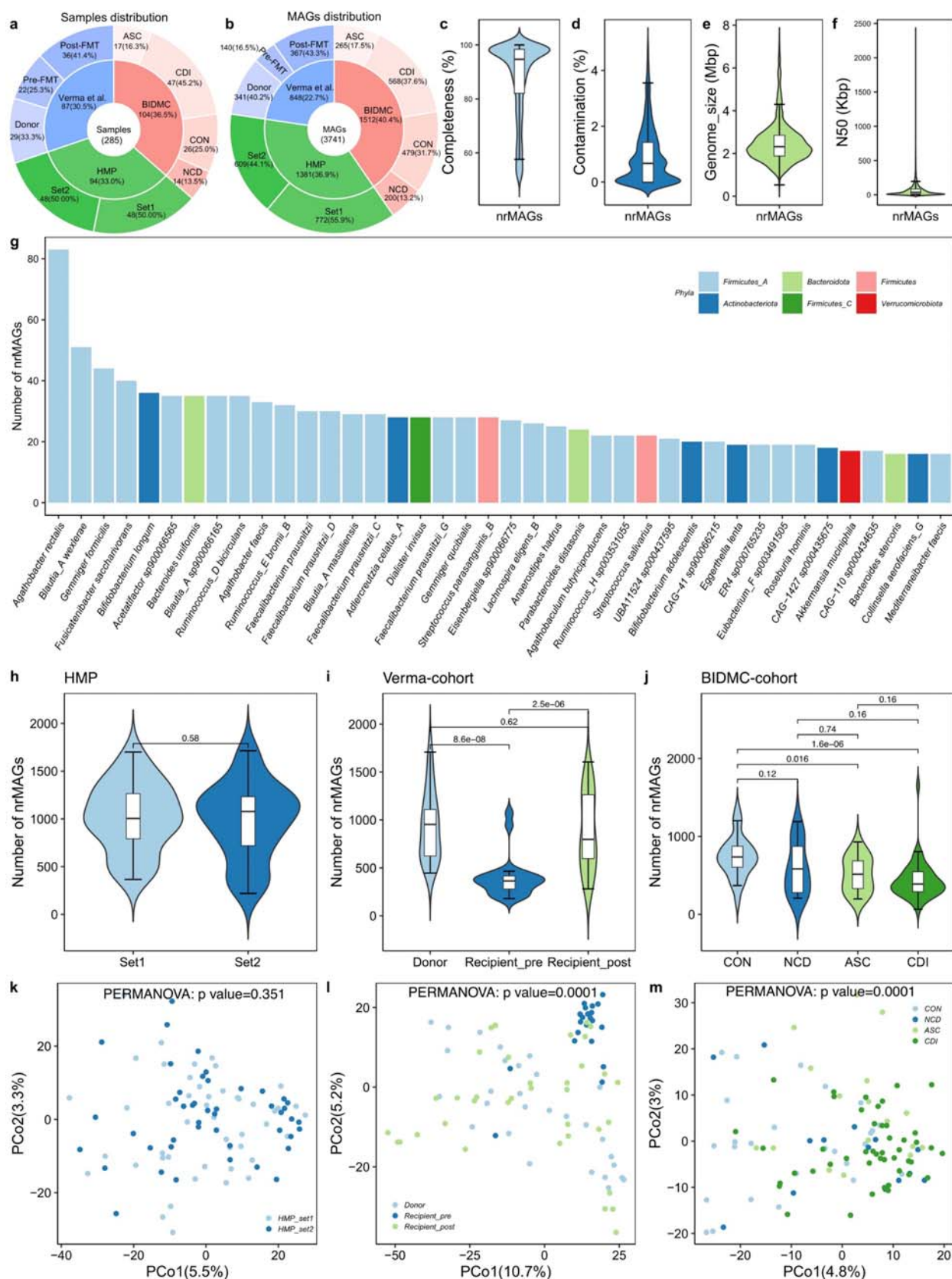
749

750

751

752

753

754

756  **Fig. 2 The microbial genome catalog and microbial diversity. a.** Sample distribution
757  among different datasets and clinical groups. **b**. Number of MAGs recovered from
758  different datasets and clinical groups. Violin plot of basic characteristics of nrMAGs on
759  completeness (**c**), contamination (**d**), genome size (**e**), and N50 (**f**). **g**. The top-40
760  species with the highest strain-richness (i.e., number of nrMAGs) identified from the
761  microbial genome catalog. The color of each bar signifies the phylum. Richness
762  (number of identified nrMAGs) of the gut microbiome from HMP (**h**), Verma-cohort. (**i**),
763  and BIDMC-cohort (**j**). Principal Coordinates Analysis (PCoA) plot based on robust
764  Aitchison distance from HMP (**k**), Verma-cohort. (**l**), and BIDMC-cohort (**m**). All
765  PERMANOVA tests were performed with 9999 permutations based on robust Aitchison
766  distance, two-sided.

767

768

769

770

771

772

773

774

775

776

777

778
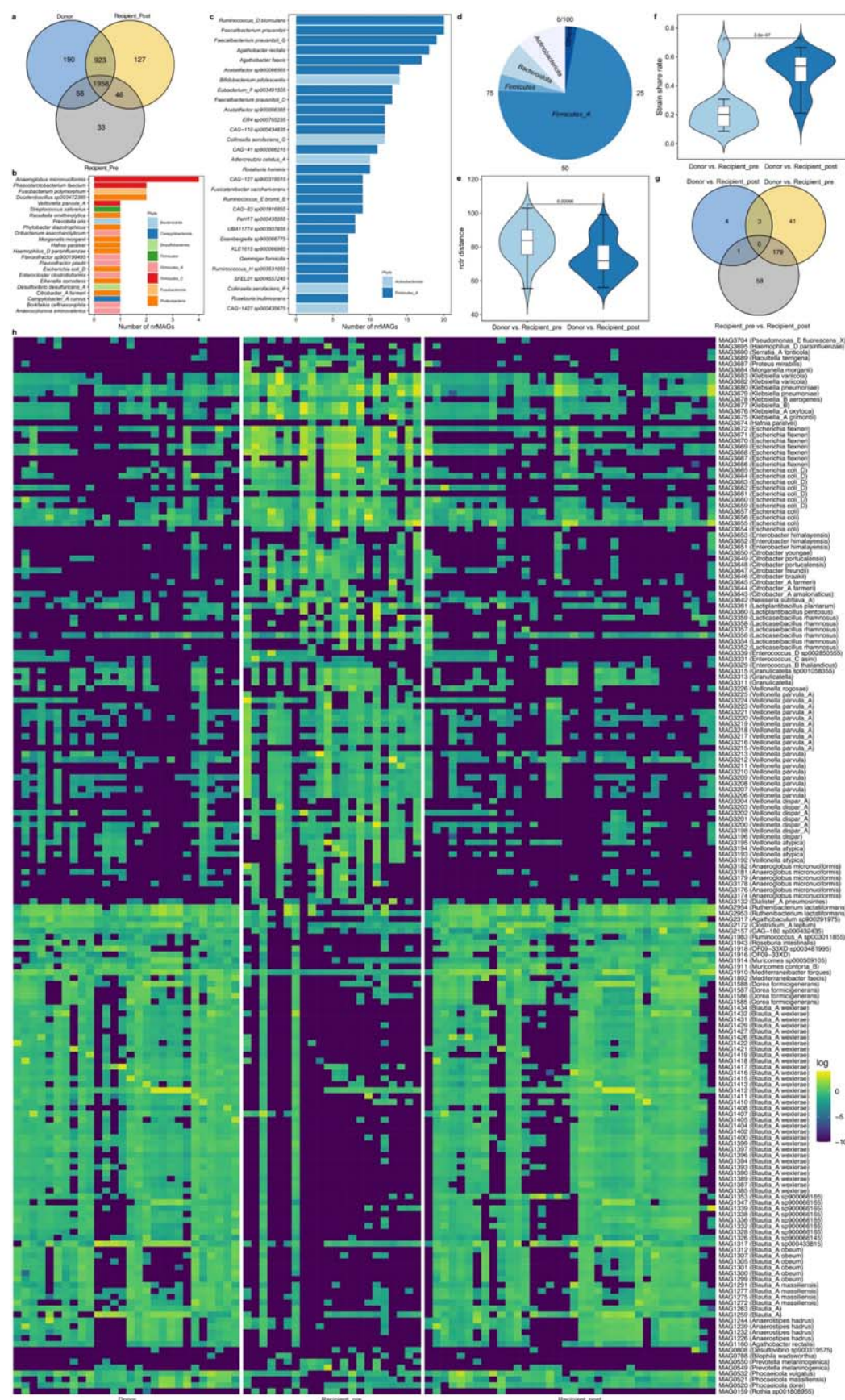
779

780

781

782

783

784

785

786  **Fig. 3 The changes in recipients' microbiome after FMT. a**. The distribution of
787  nrMAG among donors, pre-FMT recipients, and post-FMT recipients. **b**, The distribution
788  of lost nrMAGs after FMT at the species level, and the color of each bar represents the
789  phylum. **c**. The distribution of engrafted nrMAGs after FMT at the species level, and the
790  color of each bar represents the phylum. **d**. The distribution of engrafted nrMAGs after
791  FMT at the phylum level. **e**, The robust Aitchison distance between donor and recipient
792  pairs before and after FMT. **f**, The nrMAG share rate between donor and recipient pairs
793  before and after FMT. **g**. The differential abundant nrMAG distribution among three pair-
794  wise comparisons between donors, pre-FMT recipients, and post-FMT recipients. **h**.
795  The heat map showed the abundance distribution of overlapped nrMAGs identified from
796  the comparisons of donor vs. pre-FMT recipients and pre-FMT recipients vs. post-FMT
797  recipients. These nrMAGs were taxonomically annotated using GTDB-Tk based on the
798  Genome Taxonomy Database.

799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816

817 **Table. 1 Summary of the literature evidence regarding the potential role of**
818 **protective species in CDI identified through our computational pipeline.** The top
819 25 potential protective species were selected based on the overlapped protective
820 strains identified from the GMPT results with two sets of HMP microbiome data.

| Family | Rank | Species | Taxonomy ID alternate names | PMID Author reference Sponsor | Study size and Groups | Effect |
|---|---|---|---|---|---|---|
| Lachnospiraceae | 1 | *Dorea formicigenerans* | 39486 *Eubacterium formicigenerans* | 34252073[π] Verma[30] | 22 rCDI (pre/post FMT) Healthy donors | *D. formicigerans* is top 5 engrafter in rCDI patients after FMT from donor |
| | 2 | *GCA-900066135 sp900066135* | 2830660 *Lachnospiraceae bacterium Marseille-Q4251* | | | |
| | 6 | *Anaerobutyricum hallii* | 39488 *Eubacterium hallii* | 32384826 Crobach[58] | 41 CDI, 41 colonized, 43 controls | ↑ *A. hallii*: controls vs. colonized (↑6.6x AB+, ↑ 3.8x AB-) |
| | 8 | *Dorea longicatena* | 88431 *Eubacterium sp. III-35* | 37060545 Louie[39] *Vedanta* | Phase 2 RCT VE303 (8-organism consortium) 79 rCDI 1:1:1 placebo (PBO): low (LD): high-dose (HD) VE303 | *D. longicatena* is part of consortium. 8-week CDI recurrence: PBO 45.5%, LD 37.0%, HD 13.8% (p=0.006 HD vs. PBO) |
| | 9 | *Blautia obeum* | 40520 *Ruminococcus obeum* | 30816855 Mullish[59] *Finch* | 14 rCDI receiving FMT (pre, post: 1, 4, 12 weeks) 5 healthy donors | Pre-FMT: ↓ bile salt hydrolase (BSH) activity and genes (*bsh*/*baiCD*), ↑ primary bile acids (taurocholic) Post-FMT: ↑ *bsh*/*baiCD*, ↑*B. obeum* (BSH producer) Culture supernatant of *B. obeum* (& 3 other BSH+ species) attenuates CDI in mouse model |
| | | | | 33854066 Berkell[60] | AB+: 14 CDI, 64 AAD, 669 no diarrhea (ND) | *B. obeum* is 21.5% of OTU30 (oligotyping) ↑OTU30 3.8x: (AAD+ND) vs. CDI ↑OTU30 3.7x: AAD vs. CDI |
| | 10 | *Bariatricus comes* | GBIF 10828568 | | | |
| | 15 | *Mediterraneibacter faecis* | 592978 *Ruminococcus faecis* | | | |
| | 17 | *Anaerostipes hadrus* | 649756 *Eubacterium hadrum* | | | |
| | 18 | *Lachnospira sp900316325* | | | | |
| | 20 | *Agathobacter rectalis* | 39491 *Pseudobacterium / Roseburia / Eubacterium rectale Bacteroides rectalis* | | | |
| | 21 | *CAG-81 sp900066535* | | | | |
| | 23 | UBA7182 sp003480725 | 1952150 | | | |
| | 25 | *Blautia wexlerae* | 418240 | 33854066 Berkell[60] | See *B. obeum* | *B. wexlerae* is 57.9% of OTU30 (oligotyping) |
| | Family level effects (human studies only) | | | 27166072 Milani[61] | 29 Non-CDI, AB+ 30 Non-CDI, AB- 25 CDI | Lachnospira relative abundance: CDI: 0.31%, AB+: 1.22%, AB-: 3.28% (p<0.005 CDI vs. AB-) |

30

| | | | | 30785932 Han[62] | *Cd tcdB*: 79 NAAT+ 20 NAAT- | Lachnospira relative abundance: NAAT + 9.00%, NAAT – 16.51% p=0.003 |
|---|---|---|---|---|---|---|
| | | | | 35045228 Feuerstadt[63] *Seres* | Phase 3 RCT 182 rCDI 1:1 PBO: SER-109 | SER-109 contains, among others, the following genera of Lachnospiraceae: *Anaerobutyricum*, *Anaerostipes*, *Bariatricus*, *Blautia*, *Dorea*, *Lachnospira*, *Mediterraneibacter* 8-week CDI recurrence: PBO 40%, SER-109 12% (p<0.001) |
| Bacteroidaceae | 5 | *Phocaeicola dorei* | 357276 *Bacteroides dorei* | | | |
| | 13 | *Phocaeicola vulgatus* | 821 *Bacteroides vulgatus* | 2566734[π] Tvede[64] | RCT 6 rCDI, rectally instilled: 2 donor feces vs. 4 bacterial strain mix | All patients had clinical resolution at 24h. No recurrence within 1 year. *P. vulgatus*: one of 10-strain in mix, inhibited *Cd* growth in vitro |
| | 14 | *Bacteroides ovatus* | 28116 *Pasteurella ovata, Pseudobacterium ovatum Bacteroides fragilis subsp. ovatus* | 2566734[π] Tvede[64] | See *P. vulgatus* | *B. ovatus*: one of 10-strain in mix, inhibited *Cd* growth in vitro |
| | | | | 29076071[π] Yoon[65] *KoBio Labs* | In vitro study of susceptibility of *Cd* cultures to supernatants of different bacterial organisms | *B. ovatus* SNUG40239 supernatant inhibited *Cd* growth in a bile acid dependent manner. |
| | | | | 30816855 Mullish[59] *Finch* | See *B. obeum* | Post-FMT: ↑*B. ovatus* (BSH producer) Culture supernatant of *B. ovatus* (& 3 other BSH+ species) attenuates CDI in mouse model |
| | | | | 31488869 Amrane[66] | 11 CDI 8 healthy donors | *B. ovatus* among top 3 cultivable organisms absent in CDI and present in >75% of controls. |
| | | | | 31660343 Hourigan[67] | 9 children with CDI/rCDI receiving FMT (pre/post) Donor stool | ↑ *B. ovatus*: recipients post vs. pre FMT (p=0.03) donors vs. recipients pre-FMT (p=0.04) |
| | | | | 38280981 Douchant[68] | Murine model of CDI 18 & 4-strain synthetic microbial communities | *B. ovatus* part of synthetic microbial community: Protects mice from CDI (RT 027, 078) Effect persists with bacteria-free supernatant. *In vitro Cd* toxin proteolysis |
| | 24 | *Bacteroides uniformis* | 820 | 36443547 Francisco[69] | 200 CDI with cancer 42 severe/fulminant 158 non-severe | *B. uniformis* top species associated with non-severe CDI by IDSA/SHEA criteria (effect size 2.5, p<0.05) |
| Oscillospiraceae | 3 | *Oscillibacter welbionis* | -- | | | |
| | 4 | *Faecalibacterium prausnitzii* | 853 *Fusobacterium prausnitzii* | 28090385 Moelling[70] | N=1, rCDI cured with FMT Followed for 4.5 years | *F. prausnitzii* and *A. municiphila* only two species engrafting at end of follow-up |
| | | | | 29385239[π] Roychowdhury[71] | CDI mouse model, 13 mice/group, orally received: *F. prausntizii* (Fp), Fp + Potato Starch (PS), PS, Supernatant of Fp + PS, and saline | *F. prausnitzii* + PS vs. saline: ↓ *Cd* DNA (day 3) ↓ IL-1β & IL-8 mRNA (day 5 vs. 1) ↑ IL-10 mRNA (day 5 vs. 1) |
| | | | | 31660343 | See *B. ovatus* | ↑ *F. prausnitzii*: |

31

| | | | Hourigan[67] | | donors vs. recipients pre-FMT (p=0.008) |
|---|---|---|---|---|---|
| | | | 32296918 Vakili[72] | 28 CDI 56 non-CDI | ↑ *F. prausnitzii:* Non-CDI vs. CDI (p=0.015) |
| | | | 32801806 Vakili[73] | 50 CDI 50 healthy controls | ↑ *F. prausnitzii:* Controls vs. CDI (p<0.05) |
| | | | 33836037 Björkqvist[38] | 15 rCDI receiving FMT (pre, post: 2 weeks, 2-4 months) 9 healthy donors | ↑ *F. prausnitzii:* Donors vs. rCDI pre-FMT (p<0.01) rCDI post-FMT vs. pre-FMT (p<0.001) |
| | | | 34924229 Shoaei[74] | 69 Burn unit patients (23 CDI, 46 non-CDI) 46 healthy controls | ↑ *F. prausnitzii:* Non-CDI vs. CDI (p<0.001) Controls vs. CDI (p=0.003) |
| | | | 35005566 Gu[75] *Finch* | 30 healthy adults, amoxicillin-clavulanate: (13 AAD, 17 non-AAD) Days: 0, 1, 2, 3, 7, 14, 28 | ↑ *F. prausnitzii:* On day 0, ↓2.33x risk of AAD Non-AAD vs. AAD (Days 1-7, p<0.05) |
| | | | 35477270 Dudik[76] | Narrative review | *F. prausnitzii* in top 2 of most promising 'next-generation probiotics' |
| | 22 | *Lawsonibacter sp900066645* | | | |
| | Family level effects (human studies only) | | 35045228 Feuerstadt[63] *Seres* | See Lachnospiraceae | SER-109 contains, among others, the following genera of Oscillospiraceae ~[Ruminococcaceae]: *Faecalibacterium*, *Lawsonibacter* |
| Streptococc aceae | 7 | *Streptococcus salivarius* | 1304 *Lactobacillus salivarius* | 12672580 Lee[77] | In vitro study: 102 lactic acid producing bacteria from 32 healthy infants. | *S. salivarius* among 12/32 strains with anti-*Cd* activity in co-culture |
| Rikenellac eae | 11 | *Alistipes putredinis* | 28117 | | | |
| unclassified Bacillota | 12 | UBA1191 | 1947933 Firmicutes bacterium UBA1191 | | | |
| Odoribacter aceae | 16 | *Odoribacter splanchnicus* | 28118 *Bacteroides splanchnicus* | 32589701 Solbach[78] | 1506 hospitalized patients 139 colonized on admission 16 new *Cd* through admission | ↑ *O. splanchnicus* on admission associated with absence of Cd colonization through admission (LDA 3.4, FDR<0.05) |
| Peptostrepto co-ccaceae | 19 | *Romboutsia timonensis* | 1776391 *Romboutsia sp. DR1* | | | |

821
822 PMID: Pubmed ID, *Cd*: *Clostridioides difficile*, CDI: *Cd* Infection, rCDI: recurrent CDI, AB+: antibiotic-exposed, AB-: non-antibiotic-
823 exposed, FMT: Fecal Microbiota Transplantation, NAAT: Nucleic Acid Amplification Test, *tcdB*: Cd Toxin B gene, RCT: Randomized
824 Controlled Trial, OUT: Operational Taxonomic Unit, RT: Ribotype, AAD: Antibiotic Associated Diarrhea, LDA score: Linear
825 Discriminant Analysis score, FDR: False Discovery Rate
826 Ⴈ denotes strain-level information available in the reference (*Dorea formicigenerans* ATCC 27755, *Phocaeicola vulgatus* A33-14,
827 *Bacteroides ovatus* A40-4, *Bacteroides ovatus* SNUG40239, *Faecalibacterium prausnitzii* ATCC 27766)
828