

1 **Title:**

2 Strain-level characterization of health-associated bacterial consortia that colonize the human  
3 gut during infancy

4

5 **Authors:**

- 6 • Samuel S. Minot<sup>1</sup> ([sminot@fredhutch.org](mailto:sminot@fredhutch.org))
- 7 • Koshlan Mayer-Blackwell<sup>2</sup> ([kmayerbl@scharp.org](mailto:kmayerbl@scharp.org))
- 8 • Andrew Fiore-Gartland<sup>2</sup> ([agartlan@fredhutch.org](mailto:agartlan@fredhutch.org))
- 9 • Andrew Johnson<sup>2</sup> ([amjohns3@fredhutch.org](mailto:amjohns3@fredhutch.org))
- 10 • Steven Self<sup>2</sup> ([sgself@uw.edu](mailto:sgself@uw.edu))
- 11 • Parveen Bhatti<sup>4,5,6</sup> ([pbhatti@bccrc.ca](mailto:pbhatti@bccrc.ca))
- 12 • Lena Yao<sup>2</sup> ([lenayao07@gmail.com](mailto:lenayao07@gmail.com))
- 13 • Lili Liu<sup>7</sup> ([lililiu@noihp.chinacdc.cn](mailto:lililiu@noihp.chinacdc.cn))
- 14 • Xin Sun<sup>8</sup> ([sunxin@niohp.chinacdc.cn](mailto:sunxin@niohp.chinacdc.cn))
- 15 • Yi Jinfan<sup>9</sup> ([13927733228@139.com](mailto:13927733228@139.com))
- 16 • \*James Kublin<sup>2,3</sup> ([jkublin@fredhutch.org](mailto:jkublin@fredhutch.org))

17 \* Corresponding author

18

19 **Affiliations:**

- 20 1. Data Core, Fred Hutchinson Cancer Center, Seattle, USA
- 21 2. Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Center, Seattle, USA
- 22 3. HIV Vaccine Trials Network, Fred Hutchinson Cancer Center, Seattle, USA

- 23 4. Cancer Control Research, BC Cancer Research Institute, Vancouver, BC, Canada
- 24 5. Epidemiology Program, Public Health Sciences Division, Fred Hutchinson Cancer Center,  
25 Seattle, USA
- 26 6. School of Population and Public Health, University of British Columbia, Vancouver, BC,  
27 Canada
- 28 7. Key Laboratory of Occupational Disease Prevention and Treatment, Guangdong Province  
29 Hospital for Occupational Disease Prevention and Treatment, Guangzhou, China
- 30 8. National Institute of Occupational Health and Poison Control, Chinese Center for Disease  
31 Control and Prevention, Beijing, China
- 32 9. Nanhai Maternity and Child Healthcare Hospital of Foshan, Foshan, China

33

34 **Abstract (250 words):**

35 Background: The human gut microbiome develops rapidly during infancy, a key window of  
36 development coinciding with maturation of the adaptive immune system. However, little is  
37 known of the microbiome growth dynamics over the first few months of life and whether there  
38 are any generalizable patterns across human populations. We performed metagenomic  
39 sequencing on stool samples (n=94) from a cohort of infants (n=15) at monthly intervals in the  
40 first six months of life, augmenting our dataset with seven published studies for a total of 4,441  
41 metagenomes from 1,162 infants.

42 Results: Strain-level *de novo* analysis was used to identify 592 of the most abundant organisms  
43 in the infant gut microbiome. Previously unrecognized consortia were identified which  
44 exhibited highly correlated abundances across samples and were composed of diverse species

45 spanning multiple genera. Analysis of a cohort of infants with cystic fibrosis identified one such  
46 novel consortium of diverse *Enterobacterales* which was positively correlated with weight gain.  
47 While all studies showed an increased community stability during the first year of life, microbial  
48 dynamics varied widely in the first few months of life, both by study and by individual.

49 Conclusion: By augmenting published metagenomic datasets with data from a newly  
50 established cohort we were able to identify novel groups of organisms that are correlated with  
51 measures of robust human development. We hypothesize that the presence of these groups  
52 may impact human health in aggregate in ways that individual species may not in isolation.

53

54 **Keywords (3-10):**

55 Microbiome, Metagenomics, Human Development, Bacterial Consortia

56

57 **Background**

58 Early-life colonization of the human gut by microorganisms can have long-term implications for  
59 physiology and disease[1-3]. Species- and strain-level analyses suggest that most taxa can be  
60 inherited from the mother during vaginal birth, and microbial transfer is likely reduced in  
61 infants born by Caesarean delivery or by those treated with antibiotics[4-6]. Disruptions to  
62 natural bacterial exposures and microbiome development (e.g., by Caesarian section delivery,  
63 excessively sterile environment, or antibiotic-treatment) are associated with increased  
64 susceptibility to inflammatory and metabolic diseases, and intervention studies in animal  
65 models have defined key pre- and post-natal developmental windows during which the  
66 developing microbiome affects important immune processes, such as tolerance induction.

67

68 Key knowledge gaps remain concerning the immune phenotypes of at-risk infant populations  
69 and how early-life complications, such as microbiome disruption, malnutrition, and pathogen  
70 exposures, alter immune ontogeny and lead to vaccine response deficiencies in some children.  
71 Emerging evidence suggests that individual variation in response to infection or vaccination  
72 may be influenced by past viral and bacterial exposures, which shape the immune system and  
73 can establish pre-existing immune-reactivity[7-10].

74

75 Murine systems and longitudinal human birth cohorts have defined critical neonatal windows in  
76 which the intestinal microbiome stimulates immune maturation and provides colonization  
77 resistance to protect against infectious and immune-mediated disease[3, 11-28]. While  
78 neonatal taxa-immune pathways remain to be fully elucidated, the acquisition of Clostridiales  
79 taxa in early-life is clearly vital [2, 29-44]. Clostridiales provide colonization resistance[13],  
80 stimulate immune-regulatory responses[18, 26, 45-47], and activate IFN-mediated lung  
81 protection[48]. A failure to acquire Clostridiales taxa, especially *Ruminococcaceae*,  
82 *Lachnospiraceae* and Clostridium Cluster XIVa, represents the major deficiency of the CF infant  
83 microbiome, a finding that is highly reproducible across multiple independent cohorts,  
84 including the most extensively characterized BONUS cohort[30, 31, 35, 37, 44].

85

86 Longitudinal studies of birth cohorts – so far conducted predominantly in North America and  
87 Europe – have begun to characterize compositional changes to the gut microbiome that occur  
88 in the first years of life. These studies have relied primarily on amplification of the bacterial 16S

89 ribosomal RNA gene or, more seldomly, whole genome sequencing[4, 6, 49-51]. These  
90 longitudinal studies, along with one major cross-sectional study[52], have demonstrated that  
91 there is considerable inter-individual and temporal variation in the neonatal and infant  
92 microbiome community starting from birth and extending to approximately three years of age.  
93 During this period the microbiome gains richness and stability to form a microbial community  
94 that is more reflective of the adult microbiome[4, 6, 49-51, 53]. This represents a general  
95 transition where bacteria specialized to the aerobic neonatal gut (e.g., *E. coli*) or for growth on  
96 complex sugars in breastmilk (e.g., Bifidobacterium and Veillonella) and are outcompeted by  
97 organisms found more commonly in the adult gut microbiome, such as Bacteroidaceae and  
98 Ruminococcaceae[4, 6, 50]. This is reflected in the metagenomic composition of the bacterial  
99 communities, with genes involved in milk oligosaccharide metabolism giving way to those  
100 better suited to solid foods, such as fiber degradation[4, 6, 49].

101

102 While these studies have elucidated general trends in infant microbiome development, most  
103 prior studies are limited by low density of fecal sampling in the first 6 months of life when  
104 temporal intraindividual variation in the microbiome is highest and exposures to the immune  
105 system are particularly impactful. Furthermore, bioinformatic approaches have focused  
106 predominantly on identifying microbiome community states that are reflective of specific ages  
107 and which are generalizable across individuals. In contrast, these approaches offer more limited  
108 insight into the growth dynamics of individual taxa or clusters of interacting consortia. It is now  
109 evident that the path of individual microbiome development is highly variable across infants.  
110 For example, a recent analysis of transitions between ten different microbiome community

111 states in early life observed great diversity in the patterns of transitions between states. In fact,  
112 the most common transition pattern was only observed in 20% percent of infants, with the  
113 remaining 80% of infants displaying unique maturation transition patterns.[54]

114

115 Recognizing that the microbiome may not conform to consistent community states, an  
116 alternate ontological approach is to identify the subsets of microbial organisms which are found  
117 together in concert, as would be expected from a group of Bifidobacteria which jointly  
118 metabolize human milk oligosaccharides. From an informatic point of view, we approach the  
119 identification of microbial consortia by testing for organisms with correlated abundances across  
120 large numbers of microbiome samples[55]. When organisms are more likely to be found  
121 together than would be expected by random assortment, we may hypothesize that there is a  
122 shared underlying biological process which is jointly driving their growth and survival.

123

124 In human microbiome research, detection of well-studied bacterial species and genera can  
125 reveal considerable information about environmental conditions present during health and  
126 disease conditions of the host; however, there are limitations in taxonomic-centered  
127 approaches that orient analysis around finding associations with host characteristics and  
128 relative abundance of bacterial groups agglomerated by phylogenetic clade. Critically,  
129 aggregating organismal relative abundances within phylogenetic clades (i.e., summarizing  
130 microbiome features to the genus, family, or order level) becomes less informative as  
131 physiology and metabolism of bacteria within taxonomically derived grouping can vary greatly.  
132 However, strain-level analysis suffers from high-dimensionality and high sparsity of features

133 between samples. Thus, a major challenge in analysis of microbiome is finding a flexible unit of  
134 analysis that permits detection of consistent and interpretable ecological changes in the host  
135 via phylogenetic-independent agglomeration of co-abundant organisms.

136

137 Thus, to gain a better understanding of dynamics of microbiome development in the critical  
138 development period between 1 to 6 months of age, we conducted such a gene-level  
139 microbiome analyses on stool samples collected monthly in a longitudinal mother-infant birth  
140 cohort[56]. Because the aggregate gene content of the gut microbiome is comprised of tens or  
141 hundreds of millions of genes[57], a meaningful embedding in lower dimensional space is  
142 helpful for comparisons across samples. For this purpose, we use Co-Abundant Gene Groups  
143 (CAGs)[58, 59] which represent sets of genes that are expected to be found together in the  
144 same genetic element (chromosome, plasmid, virus, etc.) across all of the samples in the  
145 collection. To increase the total amount of biological information used for CAG construction, we  
146 augmented the data from our own cohort with additional metagenomic data from seven  
147 published infant microbiome datasets[4, 6, 49-53], for a total of 4,441 biological samples in the  
148 combined dataset. We used a reproducible pipeline, *geneshot* [60], for constructing and  
149 quantifying CAGs to describe groups of organisms that colonize the gut according to patterns of  
150 correlated abundance which are reproduced across cohorts. We describe these groups of  
151 microbes which are present or absent in concert as previously-unrecognized microbial  
152 consortia, which may help researchers more succinctly describe the patterns of rapid turnover  
153 which are observed during the first six months of life. Moreover, parallel analysis of a published  
154 cross-sectional cohort[53] identified one such consortium whose presence was strongly

155 correlated with infant growth rate. We propose that this dynamic growth variation may  
156 underlie altered immune development between individuals and associated susceptibility to  
157 immune-mediated disease in later life, and therefore that CAG-based analysis of microbial  
158 consortia would be a useful approach for the analysis of existing and future longitudinal birth  
159 cohorts.

160

## 161 **Results**

162

### 163 *De novo* metagenomic analysis identifies 592 bacterial genomospecies in the infant gut

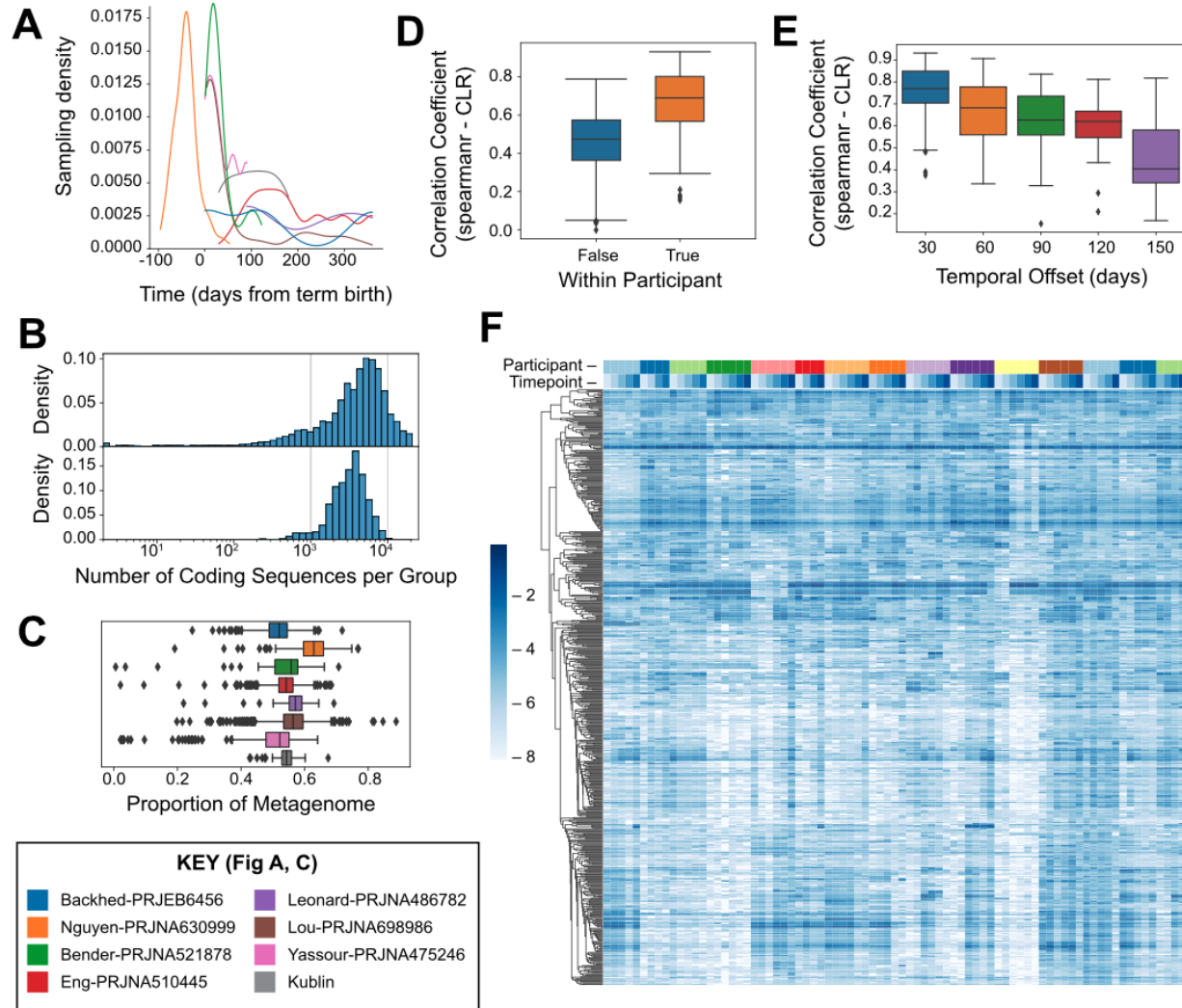
164 To quantify the relative abundance of microorganisms present in the gut during early human  
165 life, metagenomic whole-genome shotgun sequencing (WGS) data were generated from the  
166 total DNA isolated from stool samples collected from cohort of healthy infants (n = 15) at  
167 monthly intervals from birth until 6 months. To mitigate the inherent stochasticity of  
168 metagenomic sequencing, we sought to increase generalizability of gene-level analysis by  
169 including other publicly available deeply sequenced microbiome samples from other infant  
170 cohorts . the most that are In all, we analyzed metagenomic data from the seven largest  
171 published infant microbiome datasets[4, 6, 49-53], for a total of 4,441 biological samples in the  
172 combined dataset (Fig. 1A, Table 1, Table S1, S2). Gene-level metagenomic analysis was  
173 performed by: (1) generating a gene catalog via *de novo* assembly and centroid clustering  
174 (based on 90% amino acid identity); (2) estimating the relative abundance of organisms  
175 encoding each individual gene via short-read alignment; and (3) grouping genes with correlated  
176 abundances into co-abundant gene groups (CAGs) via iterative greedy single linkage



177 clustering[60]. Thus, the primary unit of measurement for downstream analyses was the  
178 relative abundance of each CAG in a particular sample, which is an estimate of the relative  
179 abundance of the organisms contained within that sample encoding the genes in that CAG.  
180 Because the groups of genes contained in each CAG have highly correlated abundances, they  
181 are predicted to be contained within the same genomic context with the metagenome ,  
182 representing the complete or partial core genome of a set of closely related isolates or  
183 strains[58].

184  
185 To focus our analysis on CAGs most likely to represent species-level groupings of genes, we  
186 subset our analysis to those CAGs containing between 1,000 and 10,000 genes ( $n = 592$ ), a  
187 range which encompasses most representative bacterial genomes in the NCBI RefSeq database  
188 (Fig. 1B). We conceptualized these CAGs as “genomospecies” because they define a group of  
189 organisms at the species or strain level based on a high degree of shared genomic content[61].  
190 The filtered set of 592 appropriately sized CAGs, i.e., genomospecies, account for over half of  
191 the raw sequence fragments recovered from infant stool samples with no clear bias by study or  
192 timepoint (Fig. 1C). While the organisms contributing the remaining sequence fragments may  
193 also have a meaningful influence on human health, they were not observed consistently at high  
194 enough abundance across multiple samples to enable gene-level analysis in this study. Using  
195 these genomospecies-level abundances as the basis of characterizing microbiome composition  
196 in our to -cohort, we compared pairs of samples from the same or different individuals using  
197 rank correlations across the CAGs. Sample pairs from the same individual were more correlated  
198 (Fig. 1D,  $p=1.16E-46$  Mann-Whitney U), as were samples collected from the same individual at

199 shorter time intervals compared to longer intervals (Fig. 1E,  $p=4.12E-10$  Spearman); together  
200 these show that there is some degree of temporal stability in community composition. A  
201 graphical summary of microbial abundances across samples is shown in Figure 1F, with each  
202 genomospecies shown in a row and each sample shown in a column. Samples are ordered by  
203 participant and timepoint, and the genomospecies are grouped by linkage clustering based on  
204 the similarity of abundance patterns across samples. The presence of genomospecies with very  
205 similar patterns of abundance in this dataset suggests that organisms are not distributed  
206 randomly across individuals, but that there may be groups of genomospecies whose relative  
207 abundance are correlated when comparing across specimens (Fig. 1F).



208

209 Figure 1. Quantification of 592 microbial genomes from 0-6 months across 8 studies.

210 A) Density of sample collection per study over time relative to term birth (40 weeks of

211 gestational age).

212 B) Distribution of genomespecies genome sizes (number of coding sequences, top) in

213 comparison to NCBI prokaryotic reference genomes (bottom).

214 C) Proportion of metagenomic sequence data from each sample which can be

215 unambiguously assigned to any one of the 592 bacterial genomes, compared

216 across studies.

217 D) Within- versus intra-participant variation in metagenome similarity estimated using the  
218 Spearman correlation coefficient of Centered-Log Ratio (CLR) abundances.

219 E) Comparison of within-participant variation in metagenome similarity between samples  
220 at varying time intervals estimated using the Spearman correlation coefficient of CLR  
221 abundances.

222 F) Comparison of microbiome composition within the samples collected for this study,  
223 with participants (n=15) and timepoints (n=6) indicated on the top marginal axis.  
224 Dendrogram indicates hierarchical clustering of the 592 genomospecies based on  
225 similarity of abundance profiles across samples. Color scale indicates log-scaled relative  
226 abundances.

227

### 228 Bacterial strains are observed in tightly correlated consortia across populations

229 To identify microbial species with correlated abundances, we calculated the Kendall rank  
230 correlation coefficient[62] for every pair of genomospecies across all of the samples from both  
231 published and newly-generated microbiome samples. Ordination of genomospecies based on  
232 similarity of abundance profiles across samples suggested a correlation within taxonomic  
233 groups (Fig. 2A; Supp. Data 1). While correlation coefficients were higher overall within  
234 taxonomic groups than between taxonomic groups (ANOSIM R=0.43 p=0.001) this was not  
235 observed universally across taxonomic groups, with many CAG pairs from the same taxa  
236 showing a complete lack of correlation (Fig. 2B, S1). Having observed that taxonomic similarity  
237 was not the primary driver of correlated abundances, we compared all pairs of genomospecies  
238 in a taxonomically-agnostic analysis. Out of all pairwise comparisons of genomospecies

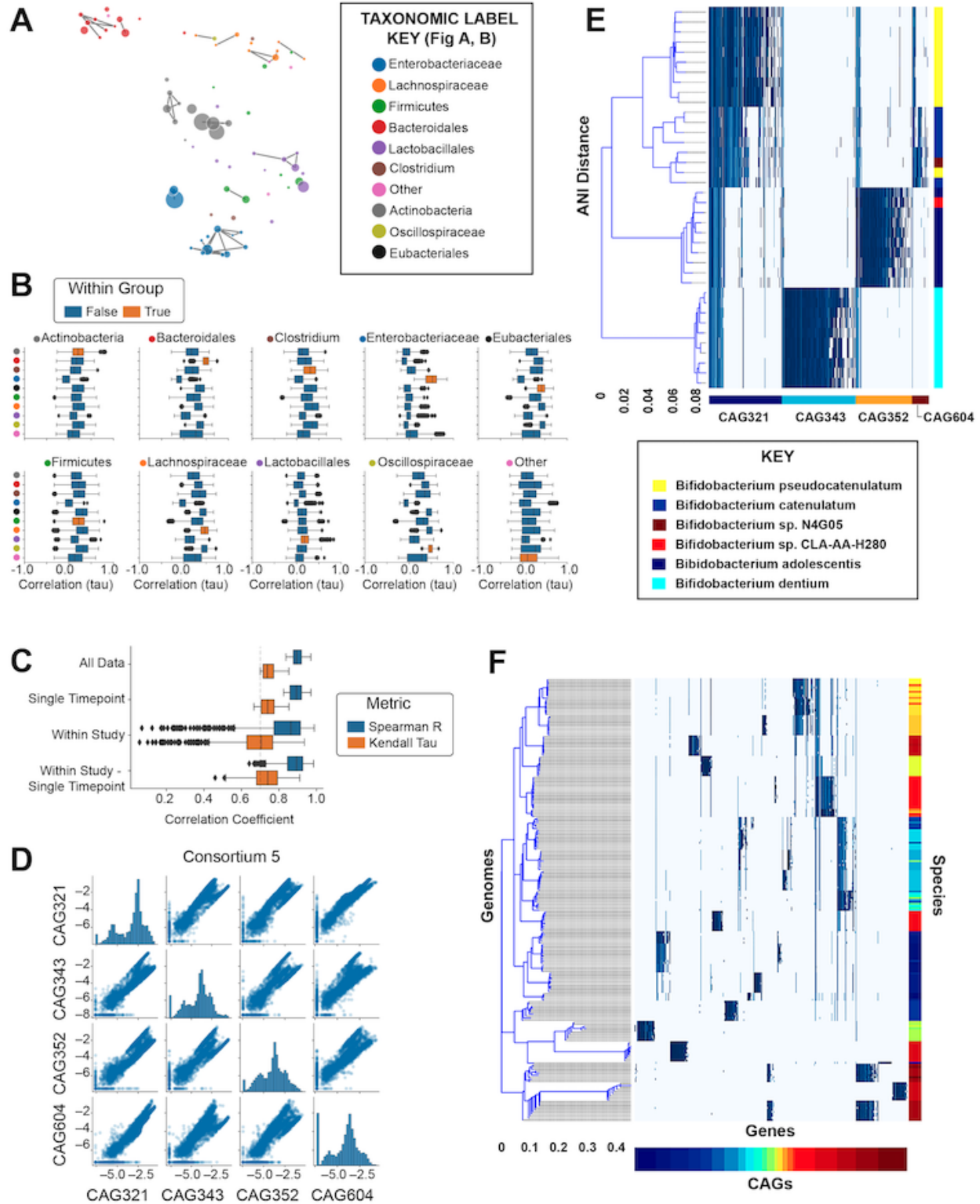
239 (n=174,936) only 143 had a Kendall's tau value of 0.7 or greater (shown with connecting lines in  
240 Fig. 2A). The correlation of this subset of CAG pairs remained strong even after filtering the data  
241 to only include a single sample from each participant or calculating the correlation  
242 independently for each study (Fig. 2C), and so it is not likely to be driven by the confounding  
243 effect of intra-individual or inter-population differences in community composition.

244

245 Next, groups of microbes were identified by single-linkage clustering using these highly  
246 correlated genomospecies, which we conceptualized as "consortia" because of their high  
247 degree of co-abundance in the infant gut microbiome (Table S3). We retained all  
248 genomospecies in this analysis, with those that did not have any correlated match being  
249 treated as individual single-genomospecies "consortia." The most abundant consortia  
250 accounted for 1-5% of all predicted genome copies on average across all specimens in the  
251 meta-analysis (Table S4, S5). To test our hypothesis that these highly correlated genomospecies  
252 represent multiple organisms (in comparison to the null hypothesis that that a single organism  
253 encodes all of the observed co-abundant genes), we compared these metagenome-derived  
254 genomospecies to the reference genomes of bacterial isolates. To identify the bacterial  
255 reference genomes which are most similar to each genomospecies we searched the NCBI  
256 RefSeq collection of bacterial genomes (n=113,938; downloaded June 6<sup>th</sup>, 2022) by amino acid  
257 sequence alignment. To better understand genomospecies relationships to conventional  
258 phylogenetic based metagenome interpretation, we closely examined the two largest groups of  
259 genomospecies with highly correlated abundances (Fig. 2D, S2). We made two observations.  
260 First, the genes contained within each individual genomospecies generally mapped to a

261 consistent set of genomes. Second, genomospecies within those consortia often mapped to  
262 different strains and species within a genus (Fig. 2A,2E) or even different orders within a class  
263 (Fig. 2F, Supp. Data 2). Finding no single genome with the complete genetic content present in  
264 these correlated genomospecies, it is likely that they represent groups of distinct organisms  
265 that are present at correlated relative abundances in the human gut microbiome during early  
266 life.

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



267

268 Figure 2. Groups of microbial genomespecies are reproducibly observed at correlated

269 abundances across studies.

- 270 A) UMAP-based ordination of genomospecies (filtered at a minimum threshold of 0.1%  
271 average abundance) based on correlation of relative abundances across samples  
272 (Kendall's tau).
- 273 B) Comparison of relative abundance-based correlation coefficients for genomospecies  
274 pairs based on order-level taxonomic annotations.
- 275 C) Considering only those pairs of genomospecies with a correlation coefficient greater  
276 than 0.7 (Kendall's tau) using all available data, correlation metrics were recalculated  
277 using a single timepoint per participant across all studies ("Single Timepoint"); while  
278 calculating an independent correlation metric for each individual study ("Within Study");  
279 or using a single timepoint per participant while also calculating an independent  
280 correlation metric for each individual study ("Within Study – Single Timepoint"). The  
281 Spearman correlation coefficient was also calculated for all comparisons (blue) in  
282 addition to Kendall's tau (orange).
- 283 D) Comparison of relative abundances (CLR) across all samples for each pair of  
284 genomospecies within Consortium 5.
- 285 E) Bacterial reference genome similarity for each of the genes within the 4 genomospecies  
286 which make up Consortium 5. Each column represents a single gene reconstructed from  
287 the metagenomic analysis. The bottom color bar indicates the genomospecies (CAG)  
288 assignment for each gene. Blue marks indicate reference genomes (each shown in a  
289 distinct row) in which that gene was detected by sequence alignment. The right-hand  
290 color bar indicates the species-level assignment for each reference genome. Hierarchical



291 clustering of reference genomes is based on the average nucleotide identity-based  
292 dissimilarity matrix.

293 F) Bacterial reference genome similarity for Consortium 3 (following E), with the full set of  
294 reference genomes available for inspection in Supplementary Data 2.

295

### 296 Relative abundance of microbial consortia changes rapidly during human infancy

297 Considering the human gut microbiome as a collection of microbial consortia, we wanted to  
298 better understand how this complex community evolves during early life. Individual consortia  
299 vary widely in relative abundance both as a function of host age as well as study population  
300 (Fig. 3A-B, Supp. Data 3). Because each study included samples from a single population, it was  
301 not possible to distinguish between study population differences and batch effects of sampling.

302 Ordinating samples based on consortium abundances across all studies shows a complex  
303 pattern, suggesting that samples at earlier timepoints are more varied in community  
304 composition, and samples at later timepoints converge on a smaller number of community  
305 types (Fig. 3C, Supp. Data 4). Consistent with this hypothesis, we found that the composition of  
306 microbial consortia was more similar at premature-birth and later timepoints within each study  
307 (Fig. 3D,  $p=0.008$  Spearman).

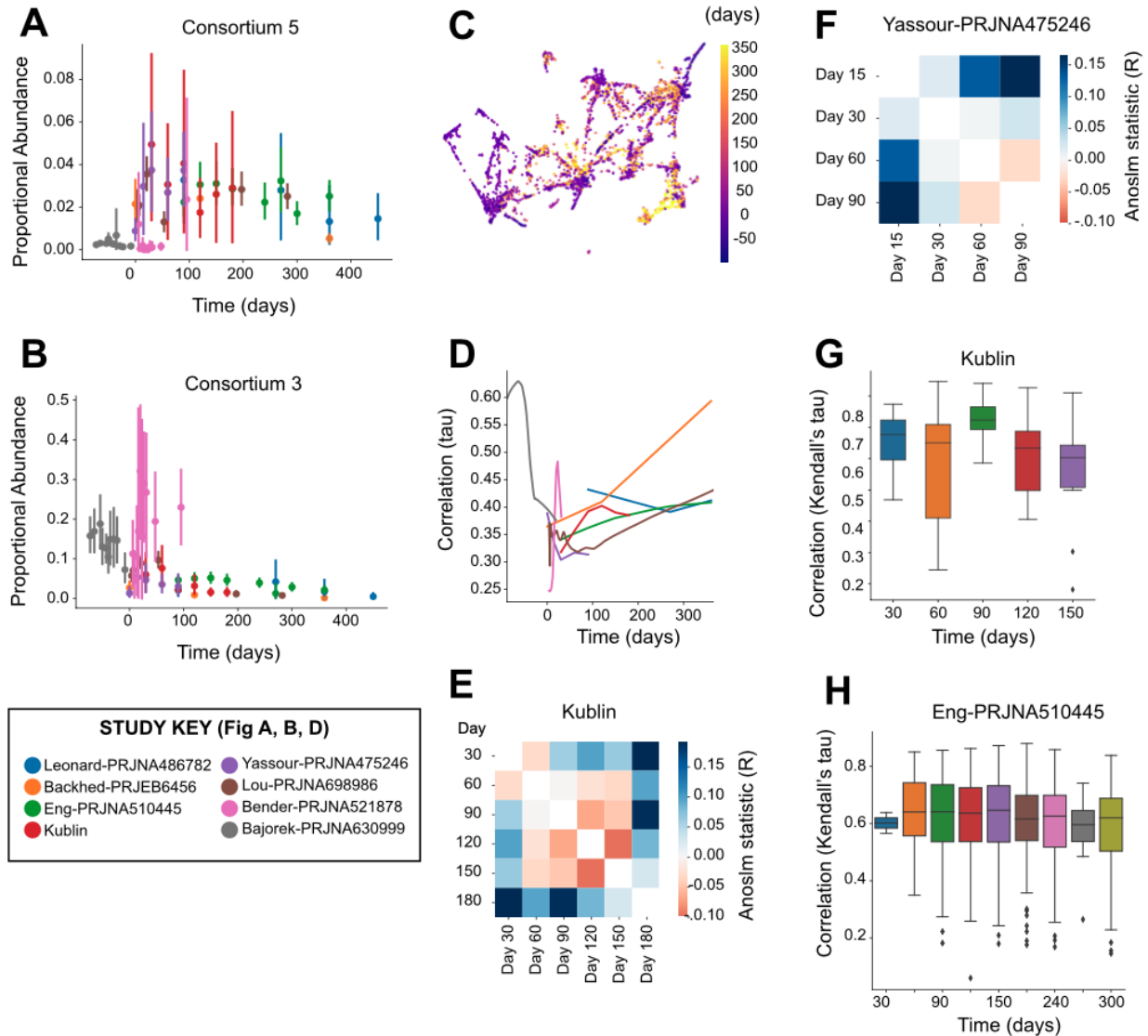
308

309 To better understand the dynamics of microbial communities as a function of host age, we used  
310 ANOSIM to compare all the samples collected at each pair of timepoints within each study. In  
311 our data, we noted a similarity of days 30-60 as well as days 90-150, with day 180 as the most  
312 distinct (Fig. 3E). In contrast, the Yassour et al. dataset showed a greater similarity of days 60-90

313 than 30-60 (Fig. 3F). Looking entirely at the similarity of samples from the same individual over  
314 time, our data showed a greater degree of change (lower correlation) from days 60-90 than 30-  
315 60 or 90-120 (Fig. 3G), while the Eng et al. data showed a greater degree of change from 30-60  
316 than 60-90 or 90-120 (Fig. 3H). The combined analysis across cohorts emphasizes the high  
317 degree of interpersonal heterogeneity and temporal transience in the human gut microbiome  
318 during early life.

319

320



321

322 Figure 3. Rapid changes in relative abundance of microbial consortia during early human life.

323 A) Summary of the relative abundance of Consortium 5 (vertical axis) across stool samples  
 324 as a function of time since term birth (horizontal axis), comparing samples obtained  
 325 from different studies (indicated by color).

326 B) Summary of the relative abundance for Consortium 3, as in (A).

- 327 C) UMAP-based ordination of microbiome samples based on similarity of microbiome  
328 composition as measured by the relative abundance of microbial consortia. Colors  
329 indicate the timepoint of sample collection relative to term birth.
- 330 D) Similarity of sample composition was compared for pairs of samples collected at similar  
331 timepoints from different individuals within each study using Kendall's tau. The  
332 horizontal axis indicates the time of sampling, and the vertical axis indicates the  
333 similarity of microbial abundances observed between different individuals.
- 334 E) Similarity of microbial relative abundances were compared between pairs of samples  
335 collected at different timepoints from different individuals within the samples collected  
336 for this study (Kublin). The pairwise comparison of each timepoint using the ANOSIM R  
337 metric is shown in a heatmap, with positive values indicating more distinct microbial  
338 compositions within each of the pair of timepoints and negative values indicating more  
339 similar microbial compositions within the pair of timepoints.
- 340 F) Similarity of microbial relative abundances for the samples from the Yassour study (as in  
341 E).
- 342 G) Similarity of samples collected from the same individual at adjacent timepoints within  
343 the samples collected for this study (Kublin). The horizontal axis indicates the timepoint  
344 which was compared to samples from the immediately preceding timepoint. Higher  
345 values on the vertical axis indicate a greater similarity of samples based on the relative  
346 abundance of microbial consortia.
- 347 H) Similarity of samples collected from the same individual from the Eng study (as in G).
- 348

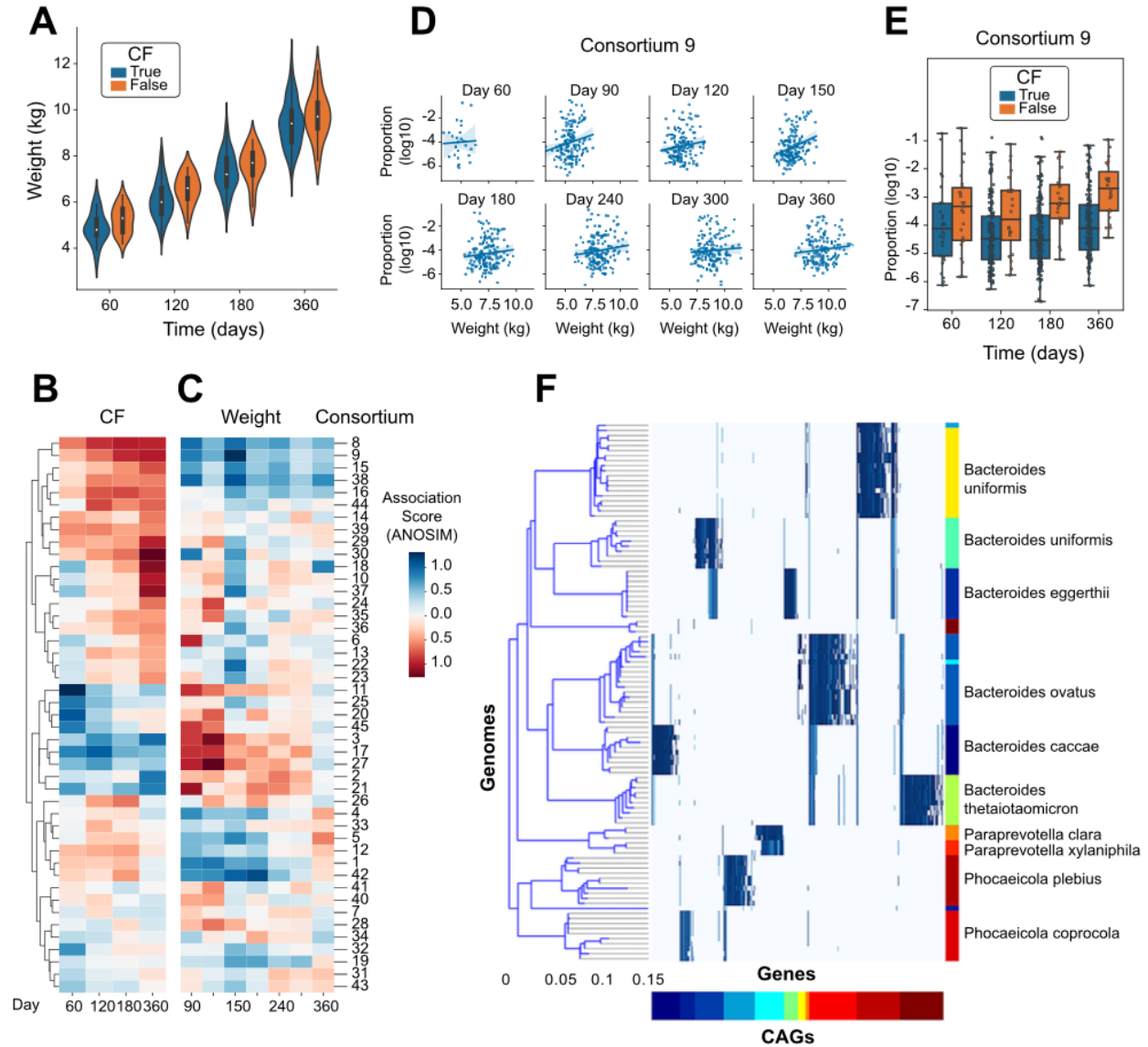
349 Taxonomically diverse microbial consortia are less abundant in the gut of infants with cystic  
350 fibrosis

351 An important application of detailed microbiome analysis is to identify microbes which may  
352 influence human health and disease. A study published by Eng, et al.[53] paired metagenomic  
353 sequencing with infant health indicators and compared the metabolic pathways encoded by the  
354 microbiome with inflammation and nutritional failure in We accessed a rich dataset that paired  
355 data from infants with cystic fibrosis (CF; n=207) to healthy controls (n=25). As previously  
356 observed[63], infants with CF had lower weight than healthy controls at each timepoint (Fig.  
357 4A, Wilcoxon  $p=0.0017$ ). To identify organisms with relative abundances that are correlated  
358 with CF status and/or weight, we performed independent linear modeling at each timepoint.  
359 The weight-association analysis was performed using only samples from participants with CF.  
360 Because of uneven sampling between groups, the CF- and weight-association analyses were  
361 performed over an overlapping but distinct set of days. When comparing the strength of  
362 association with these two clinical features across the consortia, we observed that the  
363 organisms with positive weight-associations (observed at higher abundances in CF infants with  
364 greater weights) generally had negative CF-associations (observed at lower abundances in CF  
365 infants compared to healthy controls), and vice versa (Fig. 4B-C).  
366 The microbial consortium showing the strongest association with weight in the CF infants was  
367 #9 (Fig. 4D), which was also found at lower abundance in CF infants compared to healthy  
368 controls (Fig. 4E). Alignment of the genomic markers of consortium #9 against the NCBI RefSeq  
369 catalog of microbial genomes identified species spanning *Bacteroides* (*B. uniformis*, *B. stercoris*,  
370 *B. eggerthii*, *B. ovatus*, *B. caccae*, and *B. thetaiotaomicron*), *Paraprevotella clara/xylaniphila*,

371 and *Phocaeicola* (*P. plebius* and *P. coprocola*) (Fig 4F). While these genera have been identified  
372 previously as being altered in the gut microbiome of infants with CF, these results: (a) identify  
373 the species that are most likely involved with a specific set of genomic markers (Supp. Data 5),  
374 (ii) indicate that those species are generally found together rather than individually in the gut  
375 microbiome, and (iii) suggest that the combined metabolism of a multi-species consortium may  
376 collectively mediate weight gain in infants with CF .

377

378



379

380 Figure 4. Association of specific microbial consortia with infant CF status and weight.

381 A) Measured weight of each infant at each timepoint, distinguishing infants diagnosed with

382 CF from healthy controls.

383 B) Estimated coefficient of association for the relative abundance of each microbial

384 consortium with CF status, calculated independently at each timepoint.

- 385 C) Estimated coefficient of association for the relative abundance of each microbial  
386 consortium with infant weight using only those participants diagnosed with CF,  
387 calculated independently at each timepoint.
- 388 D) Comparison of the relative abundance of Consortium 9 with weight within the group of  
389 participants diagnosed with CF, shown independently at each timepoint.
- 390 E) Comparison of the relative abundance of Consortium 9 between participants  
391 distinguished by CF diagnosis, shown independently at each timepoint.
- 392 F) Comparison of the genomic content of Consortium 9 to a reference genome collection,  
393 as in Figure 2E.

394

## 395 **Discussion**

### 396 Quantification of microbes sampled from the human gut

397 By using the reference-free analysis approach to microbiome analysis implemented in the  
398 *geneshot* analysis pipeline[60], our analysis aimed to expand our understanding of the infant  
399 gut microbiome. The advantage of this approach, which quantifies organisms on the basis of  
400 the genes encoded in their genome, is that it is not dependent on the composition of existing  
401 genome databases to detect and quantify specific organisms. While the primary drawback of  
402 this approach is a lack of sensitivity for the detection of organisms that are not sequenced to a  
403 depth sufficient for *de novo* reconstruction, approximately 50% of the raw metagenomic data  
404 was successfully assigned to just 592 distinct genomospecies representing the most abundant  
405 organisms (Fig. 1C). Based on previous work, we expected that microbial composition would  
406 reflect some degree of individuality and temporal stability[55, 64, 65]. This expectation was



407 borne out using the genomospecies-level abundance data, with samples more similar within-  
408 than between-participants (Fig. 1D) and more similar between samples collected across shorter  
409 time intervals (Fig. 1E). Based on these high-level metrics, we gained confidence that our  
410 genomospecies-level analysis is capturing a biologically meaningful profile of the most  
411 abundant organisms in the infant gut microbiome.

412

#### 413 Observation of taxonomically distinct microbial consortia

414 In addition to the detection of previously unsequenced organisms, an advantage of *de novo*  
415 metagenomic analysis is the ability to precisely identify organisms with correlated abundances  
416 that are taxonomically similar or diverse. While marker-gene or *k*-mer based analyses run the  
417 risk of confounding taxa that share a subset of genomic content, our *de novo* gene-level  
418 analysis assigns each raw sequence read unambiguously to a single genomospecies reference  
419 (using an expectation maximization approach to resolve duplicate alignments). Moreover, by  
420 limiting to the 592 organisms, evaluating all possible pairwise correlations among microbes  
421 became computationally tractable. Using this approach, we found only 143 pairs of microbes  
422 (out of the 174,936 total pairwise comparisons) with a Kendall's tau correlation coefficient  $\geq 0.7$ .  
423 Noting the inter-participant individuality of microbiome composition (Fig. 1D), we were  
424 encouraged that this high degree of pairwise correlation between individual genomospecies  
425 was observed after downsampling to a single sample per participant and after controlling for  
426 batch effects (Fig. 1C).  
427 While it is possible that genomospecies with correlated abundances may represent a single  
428 species which was inappropriately split due to noise in the metagenomic sequencing process, it

429 is more likely that correlated genomospecies represent different species with correlated  
430 abundances. One biological concept used to describe such multi-species groups would be  
431 “consortia” of distinct organisms formed by cross-feeding or syntrophic interaction[66] or by  
432 stable niche partitioning of a common source of energy (such as the degradation of diverse  
433 human milk oligosaccharides by related *Bifidobacteria*[67]). By comparing each genomospecies’  
434 genetic content to the extensive NCBI RefSeq genome collection, we observed candidate  
435 consortia containing genetically distinct organisms from the same genus (e.g., Consortium 5,  
436 Fig. 2E, Supp. Data 2), as well as single consortia containing organisms spanning multiple  
437 diverse genera (*Klebsiella*, *Enterobacter*, *Leclercia*, *Citrobacter*, *Cronobacter*, *Proteus*, *Serratia*,  
438 and *Pseudomonas*) (e.g, Consortium 3, Fig. 2F, Supp. Data 2). The robust correlation of relative  
439 abundances between these genetically distinct organisms is highly unlikely to be caused by  
440 technical artifacts, and strongly suggests that these groups of organisms are present or absent  
441 in the microbiome as a correlated group.

442

#### 443 Complex, rapid temporal dynamics of the infant gut microbiome

444 Using the aggregate abundances of microbial consortia to measure the composition of the gut  
445 microbiome, we sought to better understand the complex temporal dynamics of the developing  
446 human microbiome during infancy. While there were some consistent patterns across datasets  
447 – consortium #5 of *Bifidobacteria* was observed at higher abundance during later timepoints  
448 (Fig. 3A) and consortium #3 of diverse *Enterobacterales* observed at higher abundance during  
449 earlier timepoints (Fig. 3B) – those patterns were not consistent across all individuals or all  
450 studies. Clustering of samples by total community composition did not reveal any single

451 community state associated with earlier or later timepoints (Fig. 3C). The most consistent  
452 pattern we observed was that microbial communities from later timepoints were more similar  
453 across individuals than the communities from later timepoints, an effect which was observed  
454 across multiple independent studies (Fig. 3D) and which is consistent with the previous  
455 observations made within single cohorts [11, 54, 55].

456 The development of the human microbiome during the earliest days of life is a highly dynamic  
457 process which has not been measured at consistent, dense intervals across previous studies.

458 We augmented the published set of microbiome studies by collecting stool samples at 30-day  
459 intervals from birth to day 180 in a cohort of 15 infants. While the stool microbiome in this  
460 cohort was more similar between days 30-60 and 90-150 (Fig. 3E), a previous study has shown a  
461 stronger signal of similarity between days 60-90 (Fig. 3F). To identify the patterns of  
462 microbiome development that are consistent across populations, the field will need to collect  
463 considerably more metagenomic data at higher temporal frequency from this early time period  
464 across multiple geographically diverse study sites.

465

466 Identification of a multispecies microbial consortia associated with health outcomes in infants  
467 with cystic fibrosis

468 To assess whether any of the newly-identified microbial consortia were correlated with human  
469 health outcomes, we focused on a study of infants with cystic fibrosis [53], performing linear  
470 modeling of consortium abundances with both CF status and weight among infants with CF  
471 (those analyses being performed independently). The strongest association was observed with  
472 a multispecies consortium (#9) containing species of *Bacteroides*, *Paraprevotella*, and

473 *Phocaeicola* (Fig. 4F). This group of organisms was found at lower abundance in infants with CF  
474 compared to healthy controls; and critically was also found at lower abundance in those infants  
475 with CF who weighed less; in particular at day 90 of life (Fig. 4B-E). Our biological interpretation,  
476 which is heavily influenced by the identification of this microbial consortium, is that there is a  
477 mechanistic link between the presence of this group of microbes with the factors influencing  
478 weight gain during early life. Similar to the joint metabolism of human milk oligosaccharides  
479 distributed across related *Bifidobacteria*[67], we hypothesize that there is a consequence from  
480 the presence of this group of organisms which may not be recapitulated by any single member.  
481 When translating these findings to a controlled experimental setting, our results would imply  
482 that the administration of any single species may not be sufficient to reproduce the same  
483 biological effect, but instead the full or partial set of the multi-species community may be  
484 required.

485

## 486 **Methods**

### 487 **Study sites and enrollment of cohort.**

488 We obtained data from a collaborative mother-infant cohort that enrolled pregnant women in  
489 the Guangdong and Zhejiang Provinces of China from and followed their newborn off-spring  
490 from birth up to two years of age[56, 68]. Samples were collected for this analysis from  
491 12/19/2017 to 08/21/2018. Pregnant women provided written informed consent and were  
492 screened and enrolled between 14 and 20 weeks of gestation. The study completed enrollment  
493 into the cohort in January, 2020.

494

495 **Specimen and data collection and processing.**

496 Stool samples were initially collected at the hospital (10 grams of fecal matter collected from  
497 diapers, placed in plastic containers and stored at -80°C). Subsequent stool samples were  
498 collected monthly by parents in the home. On the morning of sample collection, a cooler with  
499 ice packs and sample collection materials was delivered to the home of each participant. In the  
500 early evening, the coolers were collected and returned to the laboratory where samples were  
501 aliquoted, labelled and stored at -80°C. If an infant did not provide a stool sample on the  
502 collection day, collection was rescheduled for the following day and a new cooler was provided.

503

504 **Sequencing**

505 DNA was extracted from each stool sample (n=94) and prepared for sequencing using the  
506 TruPrep DNA Library Prep Kit V2 for Illumina. Libraries were clustered and sequenced on an  
507 Illumina HiSeq2000 instrument and sequenced to an average depth of 61.7 million paired-end  
508 reads per sample.

509

510 Sequencing data from published datasets were obtained using the SRA Toolkit, downloading all  
511 paired-end FASTQ data available from the BioProject accessions PRJNA521878 (Bender),  
512 PRJEB6456 (Backhed), PRJNA630999 (Bajorek), PRJNA475246 (Yassour), PRJNA698986 (Lou),  
513 PRJNA486782 (Leonard), and PRJNA510445 (Eng).

514

515 **Analysis**

516 Identifying and quantifying CAGs from metagenomes

517 While previous reports have described bacteria in broad taxonomic groups (e.g.  
518 Bifidobacteriaceae, Lactobacillus, Enterococcus, Bacteroides, Streptococcus), we used gene-  
519 level metagenomics to increase this level of resolution to the species- and strain-level, while  
520 also identifying horizontally transferred genetic elements which play a role in microbiome  
521 development. Analysis of raw FASTQ datasets was performed using the geneshot analysis  
522 pipeline, available at <https://github.com/Golob-Minot/geneshot>. The exact version of that  
523 software used was v0.9 with the commit hash 4d700993660ed8fdf4df6432d2c7cb2ddd8ce85f.  
524 The geneshot pipeline (described previously [60]) performs the following bioinformatics  
525 analysis steps:

- 526 1. De novo assembly of each sample independently (using megahit v1.2.9);
- 527 2. Identification of protein-coding sequences in each assembly (using prodigal v2.6.3);
- 528 3. Deduplication of protein-coding sequences at 90% sequencing identity and 50%  
529 coverage (using linclust/MMseqs2 release 12-113e3);
- 530 4. Alignment of conceptually translated sequence reads against that deduplicated gene  
531 catalog (using DIAMOND v0.9.10);
- 532 5. Clustering of protein-coding sequences into CAGs using a maximum cosine distance  
533 threshold of 0.35.

534 To effectively process the large number of metagenomic samples in this project, a subset was  
535 selected for *de novo* assembly and gene identification which included only a single  
536 representative per participant across all projects, while genes and CAGs were quantified across  
537 the full set of samples. The computational resources required for this analysis were  
538 considerable, with ~104,000 CPU hours required for the *de novo* assembly and gene

539 identification and ~364,000 CPU hours required to align the full dataset against that gene  
540 catalog.

541

#### 542 Comparing the composition of microbial communities

543 The similarity of organisms present in different microbial communities was estimated using the  
544 non-parametric Spearman correlation of CLR-transformed abundances. The Spearman R value  
545 was used when comparing pairs of samples in terms of their microbial composition. When  
546 comparing pairs of CAGs to find organisms with correlated abundances, the more conservative  
547 Kendall Tau metric was also calculated using the CLR-transformed abundances.

548

#### 549 Identifying genomospecies associated with health status

550 Statistical analysis for the association of genomospecies relative abundance with the health  
551 status of human hosts (CF status and weight) was performed using Generalized Estimating  
552 Equations as implemented in the statsmodels package (Python). All GEE models were  
553 constructed using an exchangeable covariance structure and Gaussian family. Adjustment for  
554 multiple hypothesis testing was performing with the FDR-BH protocol as implemented in  
555 statsmodels.

556

#### 557 Comparing genetic content of genomospecies to reference genomes

558 The reference genomes most closely resembling the organisms reconstructed *de novo* from this  
559 metagenomic dataset were identified by alignment against the NCBI RefSeq database  
560 (downloaded June 6, 2022). The protein-coding sequences from each CAG (“genomospecies”)

561 was aligned against that genome collection using the gig-map workflow (available at  
562 <https://github.com/FredHutch/gig-map/>), which employs the DIAMOND aligner for rapid  
563 alignment of conceptually-translated genomes in amino acid space, at a minimum alignment  
564 threshold of 90% sequence identity and 90% alignment coverage (of the *de novo* assembled  
565 gene sequence). The similarity of reference genomes (used for the dendrogram display in CAG-  
566 genome heatmaps) was estimated by gig-map with Average Nucleotide Identity (ANI)  
567 calculated using the MASH software [69].

568

569

## 570 **Tables**

571

### 572 Table 1

573

| <b>Dataset</b>                 | <b>Citation</b> | <b>Samples<br/>(#)</b> | <b>Participants<br/>(#)</b> | <b>Mean Reads per<br/>Sample</b> |
|--------------------------------|-----------------|------------------------|-----------------------------|----------------------------------|
| <b>Kublin</b>                  | This<br>study   | 94                     | 15                          | 61,722,891.79                    |
| <b>Bender-<br/>PRJNA521878</b> | <sup>1</sup>    | 62                     | 29                          | 4,166,371.97                     |
| <b>Backhed-PRJEB6456</b>       | <sup>2</sup>    | 400                    | 100                         | 39,764,708.24                    |
| <b>Nguyen-<br/>PRJNA630999</b> | <sup>3</sup>    | 292                    | 77                          | 64,487,422.86                    |



|                                 |   |       |     |               |
|---------------------------------|---|-------|-----|---------------|
| <b>Yassour-<br/>PRJNA475246</b> | 4 | 169   | 43  | 28,045,378.54 |
| <b>Lou-PRJNA698986</b>          | 5 | 2,049 | 642 | 26,306,197.79 |
| <b>Leonard-<br/>PRJNA486782</b> | 6 | 96    | 24  | 43,482,751.56 |
| <b>Eng-PRJNA510445</b>          | 7 | 1,279 | 232 | 27,790,671.30 |

574 Summary of metagenomic WGS gut microbiome datasets included in meta-analysis.

575

## 576 **Supplementary Tables**

577

578 • Supplementary Table 1: Manifest with metadata for all specimens, including number of  
579 reads, number of genes detected, etc.

580 • Supplementary Table 2: Relative abundance of all genomes across all samples

581 • Supplementary Table 3: Annotation of genomes by consortium

582 • Supplementary Table 4: Annotation of consortia, mean relative abundance

583 • Supplementary Table 5: Relative abundance of all consortia across all samples

584

## 585 **Supplementary Figures**

586

587 • Supplementary Figure 1: Kendall's tau within the most frequent taxonomic groups at  
588 various levels

589

## 590 **Supplementary Data**

591

592 • Supplementary Data 1: Interactive display showing ordination of CAGs with taxonomic  
593 annotations

594 • Supplementary Data 2: gig-map displays for all consortia

595 • Supplementary Data 3: Relative abundance displays over time and across studies for all  
596 consortia

597 • Supplementary Data 4: Interactive UMAP of samples by consortium abundance

598 • Supplementary Data 5: Sequences of genes in all CAGs

599

## 600 **Declarations**

601 Ethics approval and consent to participate: The study was approved by the ethical committee of  
602 the Chinese Center for Disease Control and Prevention.

603 Consent for publication: Not applicable

604 Availability of data and materials: The data produced by this study has been made available in

605 the Cirro Data Platform (<https://cirro.bio>) using anonymous login information – username:

606 infant\_microbiome\_2023@cirro.bio; password: Public\_Manuscript\_Data123. To ensure the

607 privacy of all participants, all sequences from the metagenomic data which align to the human

608 genome have been masked from the FASTQ files which are provided publicly. Note that

609 Supplementary Data 5 is only available via the Cirro platform due to the large file size (870MB).

610 The datasets supporting the conclusions of this article are included within the article and its

611 additional files.

612 Competing interests: The authors declare that they have no competing interests.

613 Funding: SM was supported by funding from the Microbiome Research Initiative (Fred Hutch  
614 Cancer Center, PI: David Fredricks M.D.). SM, KMB, AFG, AJ, and JK were supported by funding  
615 from NIH NIAID R01AI127100.

616 Authors' contributions: JK, XS, PB, and SS conceived and designed the study for the newly  
617 collected cohort. XS, LL, and YJ analyzed and generated data from physical specimens. SM  
618 performed bioinformatics analysis and implemented statistical analyses. SM, KMB, AFG, AJ, and  
619 JK collaboratively developed the statistical analysis approach. SM, PB, KMB, AFG, and JK  
620 collaboratively wrote the manuscript. All authors read and approved the final manuscript.

621 Acknowledgements: We would like to acknowledge the significant contribution of all study  
622 participants.

623

## 624 **References**

- 625 1. Wang S, Ryan CA, Boyaval P, Dempsey EM, Ross RP, Stanton C: **Maternal Vertical**  
626 **Transmission Affecting Early-life Microbiota Development.** *Trends Microbiol* 2020,  
627 **28**:28-45.
- 628 2. Hayden HS, Eng A, Pope CE, Brittnacher MJ, Vo AT, Weiss EJ, Hager KR, Martin BD, Leung  
629 DH, Heltshe SL, et al: **Fecal dysbiosis in infants with cystic fibrosis is associated with**  
630 **early linear growth failure.** *Nat Med* 2020, **26**:215-221.
- 631 3. Ennamorati M, Vasudevan C, Clerkin K, Halvorsen S, Verma S, Ibrahim S, Prosper S,  
632 Porter C, Yeliseyev V, Kim M, et al: **Intestinal microbes influence development of**  
633 **thymic lymphocytes in early life.** *Proc Natl Acad Sci U S A* 2020, **117**:2570-2578.
- 634 4. Backhed F, Roswall J, Peng Y, Feng Q, Jia H, Kovatcheva-Datchary P, Li Y, Xia Y, Xie H,  
635 Zhong H, et al: **Dynamics and Stabilization of the Human Gut Microbiome during the**  
636 **First Year of Life.** *Cell Host Microbe* 2015, **17**:690-703.
- 637 5. Turnbaugh PJ, Hamady M, Yatsunencko T, Cantarel BL, Duncan A, Ley RE, Sogin ML, Jones  
638 WJ, Roe BA, Affourtit JP, et al: **A core gut microbiome in obese and lean twins.** *Nature*  
639 2009, **457**:480-484.
- 640 6. Yassour M, Jason E, Hogstrom LJ, Arthur TD, Tripathi S, Siljander H, Selvenius J,  
641 Oikarinen S, Hyoty H, Virtanen SM, et al: **Strain-Level Analysis of Mother-to-Child**

- 642 **Bacterial Transmission during the First Few Months of Life.** *Cell Host Microbe* 2018,  
643 **24:**146-154 e144.
- 644 7. Edwards KM: **Maternal antibodies and infant immune responses to vaccines.** *Vaccine*  
645 2015, **33:**6469-6472.
- 646 8. Voysey M, Kelly DF, Fanshawe TR, Sadarangani M, O'Brien KL, Perera R, Pollard AJ: **The**  
647 **Influence of Maternally Derived Antibody and Infant Age at Vaccination on Infant**  
648 **Vaccine Responses : An Individual Participant Meta-analysis.** *JAMA Pediatr* 2017,  
649 **171:**637-646.
- 650 9. Zimmermann P, Curtis N: **Factors That Influence the Immune Response to Vaccination.**  
651 *Clin Microbiol Rev* 2019, **32.**
- 652 10. Tsang JS, Dobano C, VanDamme P, Moncunill G, Marchant A, Othman RB, Sadarangani  
653 M, Koff WC, Kollmann TR: **Improving Vaccine-Induced Immunity: Can Baseline Predict**  
654 **Outcome?** *Trends Immunol* 2020, **41:**457-465.
- 655 11. Kostic AD, Gevers D, Siljander H, Vatanen T, Hyotylainen T, Hamalainen AM, Peet A,  
656 Tillmann V, Poho P, Mattila I, et al: **The dynamics of the human infant gut microbiome**  
657 **in development and in progression toward type 1 diabetes.** *Cell Host Microbe* 2015,  
658 **17:**260-273.
- 659 12. Gray J, Oehrle K, Worthen G, Alenghat T, Whitsett J, Deshmukh H: **Intestinal commensal**  
660 **bacteria mediate lung mucosal immunity and promote resistance of newborn mice to**  
661 **infection.** *Sci Transl Med* 2017, **9.**
- 662 13. Kim YG, Sakamoto K, Seo SU, Pickard JM, Gilliland MG, 3rd, Pudlo NA, Hoostal M, Li X,  
663 Wang TD, Feehley T, et al: **Neonatal acquisition of Clostridia species protects against**  
664 **colonization by bacterial pathogens.** *Science* 2017, **356:**315-319.
- 665 14. Fujimura KE, Sitarik AR, Havstad S, Lin DL, Levan S, Fadrosch D, Panzer AR, LaMere B,  
666 Rackaityte E, Lukacs NW, et al: **Neonatal gut microbiota associates with childhood**  
667 **multisensitized atopy and T cell differentiation.** *Nat Med* 2016, **22:**1187-1191.
- 668 15. Vatanen T, Franzosa EA, Schwager R, Tripathi S, Arthur TD, Vehik K, Lernmark A,  
669 Hagopian WA, Rewers MJ, She JX, et al: **The human gut microbiome in early-onset type**  
670 **1 diabetes from the TEDDY study.** *Nature* 2018, **562:**589-594.
- 671 16. Vatanen T, Kostic AD, d'Hennezel E, Siljander H, Franzosa EA, Yassour M, Kolde R,  
672 Vlamakis H, Arthur TD, Hamalainen AM, et al: **Variation in Microbiome LPS**  
673 **Immunogenicity Contributes to Autoimmunity in Humans.** *Cell* 2016, **165:**842-853.
- 674 17. Arrieta MC, Stiemsma LT, Dimitriu PA, Thorson L, Russell S, Yurist-Doutsch S, Kuzeljevic  
675 B, Gold MJ, Britton HM, Lefebvre DL, et al: **Early infancy microbial and metabolic**  
676 **alterations affect risk of childhood asthma.** *Sci Transl Med* 2015, **7:**307ra152.
- 677 18. Abdel-Gadir A, Stephen-Victor E, Gerber GK, Noval Rivas M, Wang S, Harb H, Wang L, Li  
678 N, Crestani E, Spielman S, et al: **Microbiota therapy acts via a regulatory T cell**  
679 **MyD88/RORgammat pathway to suppress food allergy.** *Nat Med* 2019, **25:**1164-1174.
- 680 19. Cahenzli J, Koller Y, Wyss M, Geuking MB, McCoy KD: **Intestinal microbial diversity**  
681 **during early-life colonization shapes long-term IgE levels.** *Cell Host Microbe* 2013,  
682 **14:**559-570.
- 683 20. Brand S, Teich R, Dicke T, Harb H, Yildirim AO, Tost J, Schneider-Stock R, Waterland RA,  
684 Bauer UM, von Mutius E, et al: **Epigenetic regulation in murine offspring as a novel**

- 685            **mechanism for transmaternal asthma protection induced by microbes. *J Allergy Clin*  
686            *Immunol* 2011, **128**:618-625 e611-617.**
- 687    21.    Herbst T, Sichelstiel A, Schar C, Yadava K, Burki K, Cahenzli J, McCoy K, Marsland BJ,  
688            Harris NL: **Dysregulation of allergic airway inflammation in the absence of microbial**  
689            **colonization.** *Am J Respir Crit Care Med* 2011, **184**:198-205.
- 690    22.    Olszak T, An D, Zeissig S, Vera MP, Richter J, Franke A, Glickman JN, Siebert R, Baron RM,  
691            Kasper DL, Blumberg RS: **Microbial exposure during early life has persistent effects on**  
692            **natural killer T cell function.** *Science* 2012, **336**:489-493.
- 693    23.    Trompette A, Gollwitzer ES, Yadava K, Sichelstiel AK, Sprenger N, Ngom-Bru C, Blanchard  
694            C, Junt T, Nicod LP, Harris NL, Marsland BJ: **Gut microbiota metabolism of dietary fiber**  
695            **influences allergic airway disease and hematopoiesis.** *Nat Med* 2014, **20**:159-166.
- 696    24.    Deshmukh HS, Liu Y, Menkiti OR, Mei J, Dai N, O'Leary CE, Oliver PM, Kolls JK, Weiser JN,  
697            Worthen GS: **The microbiota regulates neutrophil homeostasis and host resistance to**  
698            **Escherichia coli K1 sepsis in neonatal mice.** *Nat Med* 2014, **20**:524-530.
- 699    25.    Constantinides MG, Link VM, Tamoutounour S, Wong AC, Perez-Chaparro PJ, Han SJ,  
700            Chen YE, Li K, Farhat S, Weckel A, et al: **MAIT cells are imprinted by the microbiota in**  
701            **early life and promote tissue repair.** *Science* 2019, **366**.
- 702    26.    Al Nabhani Z, Dulauroy S, Marques R, Cousu C, Al Bounny S, Dejardin F, Sparwasser T,  
703            Berard M, Cerf-Bensussan N, Eberl G: **A Weaning Reaction to Microbiota Is Required**  
704            **for Resistance to Immunopathologies in the Adult.** *Immunity* 2019, **50**:1276-1288  
705            e1275.
- 706    27.    Pronovost GN, Hsiao EY: **Perinatal Interactions between the Microbiome, Immunity,**  
707            **and Neurodevelopment.** *Immunity* 2019, **50**:18-36.
- 708    28.    Knoop KA, Gustafsson JK, McDonald KG, Kulkarni DH, Coughlin PE, McCrate S, Kim D,  
709            Hsieh CS, Hogan SP, Elson CO, et al: **Microbial antigen encounter during a preweaning**  
710            **interval is critical for tolerance to gut bacteria.** *Sci Immunol* 2017, **2**.
- 711    29.    Scanlan PD, Buckling A, Kong W, Wild Y, Lynch SV, Harrison F: **Gut dysbiosis in cystic**  
712            **fibrosis.** *J Cyst Fibros* 2012, **11**:454-455.
- 713    30.    Nielsen S, Needham B, Leach ST, Day AS, Jaffe A, Thomas T, Ooi CY: **Disrupted**  
714            **progression of the intestinal microbiota with age in children with cystic fibrosis.** *Sci*  
715            *Rep* 2016, **6**:24857.
- 716    31.    Manor O, Levy R, Pope CE, Hayden HS, Brittnacher MJ, Carr R, Radey MC, Hager KR,  
717            Heltshe SL, Ramsey BW, et al: **Metagenomic evidence for taxonomic dysbiosis and**  
718            **functional imbalance in the gastrointestinal tracts of children with cystic fibrosis.** *Sci*  
719            *Rep* 2016, **6**:22493.
- 720    32.    Madan JC, Koestler DC, Stanton BA, Davidson L, Moulton LA, Housman ML, Moore JH,  
721            Guill MF, Morrison HG, Sogin ML, et al: **Serial analysis of the gut and respiratory**  
722            **microbiome in cystic fibrosis in infancy: interaction between intestinal and respiratory**  
723            **tracts and impact of nutritional exposures.** *mBio* 2012, **3**.
- 724    33.    Hoffman LR, Pope CE, Hayden HS, Heltshe S, Levy R, McNamara S, Jacobs MA, Rohmer L,  
725            Radey M, Ramsey BW, et al: **Escherichia coli dysbiosis correlates with gastrointestinal**  
726            **dysfunction in children with cystic fibrosis.** *Clin Infect Dis* 2014, **58**:396-399.
- 727    34.    Hoen AG, Li J, Moulton LA, O'Toole GA, Housman ML, Koestler DC, Guill MF, Moore JH,  
728            Hibberd PL, Morrison HG, et al: **Associations between Gut Microbial Colonization in**

- 729 **Early Life and Respiratory Outcomes in Cystic Fibrosis.** *J Pediatr* 2015, **167**:138-147  
730 e131-133.
- 731 35. Vernocchi P, Del Chierico F, Russo A, Majo F, Rossitto M, Valerio M, Casadei L, La Storia  
732 A, De Filippis F, Rizzo C, et al: **Gut microbiota signatures in cystic fibrosis: Loss of host**  
733 **CFTR function drives the microbiota enterophenotype.** *PLoS One* 2018, **13**:e0208171.
- 734 36. Duytschaever G, Huys G, Bekaert M, Boulanger L, De Boeck K, Vandamme P: **Cross-**  
735 **sectional and longitudinal comparisons of the predominant fecal microbiota**  
736 **compositions of a group of pediatric patients with cystic fibrosis and their healthy**  
737 **siblings.** *Appl Environ Microbiol* 2011, **77**:8015-8024.
- 738 37. Duytschaever G, Huys G, Bekaert M, Boulanger L, De Boeck K, Vandamme P: **Dysbiosis**  
739 **of bifidobacteria and Clostridium cluster XIVa in the cystic fibrosis fecal microbiota.** *J*  
740 *Cyst Fibros* 2013, **12**:206-215.
- 741 38. Schippa S, Iebba V, Santangelo F, Gagliardi A, De Biase RV, Stamato A, Bertasi S, Lucarelli  
742 M, Conte MP, Quattrucci S: **Cystic fibrosis transmembrane conductance regulator**  
743 **(CFTR) allelic variants relate to shifts in faecal microbiota of cystic fibrosis patients.**  
744 *PLoS One* 2013, **8**:e61176.
- 745 39. Antosca KM, Chernikova DA, Price CE, Ruoff KL, Li K, Guill MF, Sontag NR, Morrison HG,  
746 Hao S, Drumm ML, et al: **Altered Stool Microbiota of Infants with Cystic Fibrosis Shows**  
747 **a Reduction in Genera Associated with Immune Programming from Birth.** *J Bacteriol*  
748 2019, **201**.
- 749 40. Dorsey J, Gonska T: **Bacterial overgrowth, dysbiosis, inflammation, and dysmotility in**  
750 **the Cystic Fibrosis intestine.** *J Cyst Fibros* 2017, **16 Suppl 2**:S14-S23.
- 751 41. Fridge JL, Conrad C, Gerson L, Castillo RO, Cox K: **Risk factors for small bowel bacterial**  
752 **overgrowth in cystic fibrosis.** *J Pediatr Gastroenterol Nutr* 2007, **44**:212-218.
- 753 42. Lisowska A, Madry E, Pogorzelski A, Szydlowski J, Radzikowski A, Walkowiak J: **Small**  
754 **intestine bacterial overgrowth does not correspond to intestinal inflammation in**  
755 **cystic fibrosis.** *Scand J Clin Lab Invest* 2010, **70**:322-326.
- 756 43. Lisowska A, Wojtowicz J, Walkowiak J: **Small intestine bacterial overgrowth is frequent**  
757 **in cystic fibrosis: combined hydrogen and methane measurements are required for its**  
758 **detection.** *Acta Biochim Pol* 2009, **56**:631-634.
- 759 44. Coffey MJ, Nielsen S, Wemheuer B, Kaakoush NO, Garg M, Needham B, Pickford R, Jaffe  
760 A, Thomas T, Ooi CY: **Gut Microbiota in Children With Cystic Fibrosis: A Taxonomic and**  
761 **Functional Dysbiosis.** *Sci Rep* 2019, **9**:18593.
- 762 45. Stefka AT, Feehley T, Tripathi P, Qiu J, McCoy K, Mazmanian SK, Tjota MY, Seo GY, Cao S,  
763 Theriault BR, et al: **Commensal bacteria protect against food allergen sensitization.**  
764 *Proc Natl Acad Sci U S A* 2014, **111**:13145-13150.
- 765 46. Atarashi K, Tanoue T, Oshima K, Suda W, Nagano Y, Nishikawa H, Fukuda S, Saito T,  
766 Narushima S, Hase K, et al: **Treg induction by a rationally selected mixture of Clostridia**  
767 **strains from the human microbiota.** *Nature* 2013, **500**:232-236.
- 768 47. Atarashi K, Tanoue T, Shima T, Imaoka A, Kuwahara T, Momose Y, Cheng G, Yamasaki S,  
769 Saito T, Ohba Y, et al: **Induction of colonic regulatory T cells by indigenous Clostridium**  
770 **species.** *Science* 2011, **331**:337-341.

- 771 48. Steed AL, Christophi GP, Kaiko GE, Sun L, Goodwin VM, Jain U, Esaulova E, Artyomov  
772 MN, Morales DJ, Holtzman MJ, et al: **The microbial metabolite desaminotyrosine**  
773 **protects from influenza through type I interferon.** *Science* 2017, **357**:498-502.
- 774 49. Bender JM, Li F, Purswani H, Capretz T, Cerini C, Zabih S, Hung L, Francis N, Chin S,  
775 Pannaraj PS, Aldrovandi G: **Early exposure to antibiotics in the neonatal intensive care**  
776 **unit alters the taxonomic and functional infant gut microbiome.** *J Matern Fetal*  
777 *Neonatal Med* 2021, **34**:3335-3343.
- 778 50. Nguyen M, Holdbrooks H, Mishra P, Abrantes MA, Eskew S, Garma M, Oca CG,  
779 McGuckin C, Hein CB, Mitchell RD, et al: **Impact of Probiotic B. infantis EVC001 Feeding**  
780 **in Premature Infants on the Gut Microbiome, Nosocomially Acquired Antibiotic**  
781 **Resistance, and Enteric Inflammation.** *Front Pediatr* 2021, **9**:618009.
- 782 51. Lou YC, Olm MR, Diamond S, Crits-Christoph A, Firek BA, Baker R, Morowitz MJ, Banfield  
783 JF: **Infant gut strain persistence is associated with maternal origin, phylogeny, and**  
784 **traits including surface adhesion and iron acquisition.** *Cell Rep Med* 2021, **2**:100393.
- 785 52. Leonard MM, Karathia H, Pujolassos M, Troisi J, Valitutti F, Subramanian P, Camhi S,  
786 Kenyon V, Colucci A, Serena G, et al: **Multi-omics analysis reveals the influence of**  
787 **genetic and environmental risk factors on developing gut microbiota in infants at risk**  
788 **of celiac disease.** *Microbiome* 2020, **8**:130.
- 789 53. Eng A, Hayden HS, Pope CE, Brittnacher MJ, Vo AT, Weiss EJ, Hager KR, Leung DH,  
790 Heltshe SL, Raftery D, et al: **Infants with cystic fibrosis have altered fecal functional**  
791 **capacities with potential clinical and metabolic consequences.** *BMC Microbiol* 2021,  
792 **21**:247.
- 793 54. Stewart CJ, Ajami NJ, O'Brien JL, Hutchinson DS, Smith DP, Wong MC, Ross MC, Lloyd RE,  
794 Doddapaneni H, Metcalf GA, et al: **Temporal development of the gut microbiome in**  
795 **early childhood from the TEDDY study.** *Nature* 2018, **562**:583-588.
- 796 55. Koenig JE, Spor A, Scalfone N, Fricker AD, Stombaugh J, Knight R, Angenent LT, Ley RE:  
797 **Succession of microbial consortia in the developing infant gut microbiome.** *Proc Natl*  
798 *Acad Sci U S A* 2011, **108 Suppl 1**:4578-4585.
- 799 56. Yao L, Liu L, Dong M, Yang J, Zhao Z, Chen J, Lv L, Wu Z, Wang J, Sun X, et al: **Trimester-**  
800 **specific prenatal heavy metal exposures and sex-specific postpartum size and growth.**  
801 *J Expo Sci Environ Epidemiol* 2022.
- 802 57. Coelho LP, Alves R, Del Rio AR, Myers PN, Cantalapiedra CP, Giner-Lamia J, Schmidt TS,  
803 Mende DR, Orakov A, Letunic I, et al: **Towards the biogeography of prokaryotic genes.**  
804 *Nature* 2022, **601**:252-256.
- 805 58. Nielsen HB, Almeida M, Juncker AS, Rasmussen S, Li J, Sunagawa S, Plichta DR, Gautier L,  
806 Pedersen AG, Le Chatelier E, et al: **Identification and assembly of genomes and genetic**  
807 **elements in complex metagenomic samples without using reference genomes.** *Nat*  
808 *Biotechnol* 2014, **32**:822-828.
- 809 59. Minot SS, Willis AD: **Clustering co-abundant genes identifies components of the gut**  
810 **microbiome that are reproducibly associated with colorectal cancer and inflammatory**  
811 **bowel disease.** *Microbiome* 2019, **7**:110.
- 812 60. Minot SS, Barry KC, Kasman C, Golob JL, Willis AD: **geneshot: gene-level metagenomics**  
813 **identifies genome islands associated with immunotherapy response.** *Genome Biol*  
814 2021, **22**:135.

- 815 61. Brenner DJ, Grimont PA, Steigerwalt AG, Fanning GR, Ageron E, Riddle CF: **Classification**  
816 **of citrobacteria by DNA hybridization: designation of *Citrobacter farmeri* sp. nov.,**  
817 ***Citrobacter youngae* sp. nov., *Citrobacter braakii* sp. nov., *Citrobacter werkmanii* sp.**  
818 **nov., *Citrobacter sedlakii* sp. nov., and three unnamed *Citrobacter* genomospecies.** *Int*  
819 *J Syst Bacteriol* 1993, **43**:645-658.
- 820 62. Kendall MG: **A New Measure of Rank Correlation.** *Biometrika* 1938, **30**:81-93.
- 821 63. Patterson KD, Kyriacou T, Desai M, Carroll WD, Gilchrist FJ: **Factors affecting the growth**  
822 **of infants diagnosed with cystic fibrosis by newborn screening.** *BMC Pediatr* 2019,  
823 **19**:356.
- 824 64. Yassour M, Vatanen T, Siljander H, Hamalainen AM, Harkonen T, Ryhanen SJ, Franzosa  
825 EA, Vlamakis H, Huttenhower C, Gevers D, et al: **Natural history of the infant gut**  
826 **microbiome and impact of antibiotic treatment on bacterial strain diversity and**  
827 **stability.** *Sci Transl Med* 2016, **8**:343ra381.
- 828 65. Vatanen T, Plichta DR, Somani J, Munch PC, Arthur TD, Hall AB, Rudolf S, Oakeley EJ, Ke  
829 X, Young RA, et al: **Genomic variation and strain-specific functional adaptation in the**  
830 **human gut microbiome during early life.** *Nat Microbiol* 2019, **4**:470-479.
- 831 66. Ferry JG, Wolfe RS: **Anaerobic degradation of benzoate to methane by a microbial**  
832 **consortium.** *Arch Microbiol* 1976, **107**:33-40.
- 833 67. Lawson MAE, O'Neill IJ, Kujawska M, Gowrinadh Javvadi S, Wijeyesekera A, Flegg Z,  
834 Chalklen L, Hall LJ: **Breast milk-derived human milk oligosaccharides promote**  
835 ***Bifidobacterium* interactions within a single ecosystem.** *ISME J* 2020, **14**:635-648.
- 836 68. Liu L, Yao L, Dong M, Liu T, Lai W, Yin X, Zhou S, Lv L, Li L, Wang J, et al: **Maternal urinary**  
837 **cadmium concentrations in early pregnancy in relation to prenatal and postpartum**  
838 **size of offspring.** *J Trace Elem Med Biol* 2021, **68**:126823.
- 839 69. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM:  
840 **Mash: fast genome and metagenome distance estimation using MinHash.** *Genome Biol*  
841 2016, **17**:132.
- 842