1  **Hospitalized premature infants are colonized by related bacterial strains**
2  **with distinct proteomic profiles**

3  Christopher. T. Brown[1], Weili Xiong[2], Matthew R. Olm[1], Brian C. Thomas[3],
4  Robyn Baker[4], Brian Firek[5], Michael J. Morowitz[5,6], Robert L. Hettich[7],
5  and Jillian F. Banfield[3,8,9*]

6  [1]Department of Plant and Microbial Biology, University of California, Berkeley, CA, USA,
7  [2]Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA,
8  [3]Department of Earth and Planetary Science, University of California, Berkeley, CA, USA,
9  [4]Division of Newborn Medicine, Children's Hospital of Pittsburgh and Magee-Womens Hospital
10  of UPMC, Pittsburgh,  United States
11  [5]Department of Surgery, University of Pittsburgh School of Medicine, Pittsburgh, PA, USA,
12  [6]Division of Pediatric General and Thoracic Surgery, Children's Hospital of Pittsburgh of UPMC,
13  Pittsburgh, PA, USA
14  [7]Chemical Science Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA
15  [8]Department of Environmental Science, Policy, and Management, University of California,
16  Berkeley, CA, USA
17  [9]Earth Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

18  *Corresponding author
19  Email: jbanfield@berkeley.edu
20  Telephone: 510-643-2155
21  Address: McCone Hall, Berkeley, CA 94720

22
23  **Running Title**

24  Colonization of the infant gut microbiome

25  **Keywords**

26  Metagenomics, Metaproteomics, Microbial Genomics, iRep, Human Microbiome,
27  Necrotizing Enterocolitis, Neonates, Microbial Colonization

28      **Abstract**

29      During the first weeks of life, microbial colonization of the gut impacts human immune

30      system maturation and other developmental processes. In premature infants, aberrant

31      colonization has been implicated in the onset of necrotizing enterocolitis (NEC), a life-

32      threatening intestinal disease. To study the premature infant gut colonization process,

33      genome-resolved metagenomics was conducted on 343 fecal samples collected during the

34      first three months of life from 35 premature infants housed in a neonatal intensive care

35      unit, 14 of which developed NEC, and metaproteomic measurements were made on 87

36      samples. Microbial community composition and proteomic profiles remained relatively

37      stable on the time scale of a week, but the proteome was more variable. Although

38      genetically similar organisms colonized many infants, most infants were colonized by

39      distinct strains with metabolic profiles that could be distinguished using metaproteomics.

40      Microbiome composition correlated with infant, antibiotics administration, and NEC

41      diagnosis. Communities were found to cluster into seven primary types, and community

42      type switched within infants, sometimes multiple times. Interestingly, some communities

43      sampled from the same infant at subsequent time points clustered with those of other

44      infants. In some cases, switches preceded onset of NEC; however, no species or

45      community type could account for NEC across the majority of infants. In addition to a

46      correlation of protein abundances with organism replication rates, we found that

47      organism proteomes correlated with overall community composition. Thus, this genome-

48      resolved proteomics study demonstrates that the contributions of individual organisms to

49      microbiome development depend on microbial community context.

**Importance**

Humans are colonized by microbes at birth, a process that is important to health and development. However, much remains to be known about the fine-scale microbial dynamics that occur during the colonization period. We conducted a genome-resolved study of microbial community composition, replication rates, and proteomes during the first three months of life of both healthy and sick premature infants. Infants were found to be colonized by similar microbes, but each underwent a distinct colonization trajectory. Interestingly, related microbes colonizing different infants were found to have distinct proteomes, indicating that microbiome function is not only driven by which organisms are present, but also largely depends on microbial responses to the unique set of physiological conditions in the infant gut.

**Introduction**

Infants have high levels of between-individual variation in microbiome composition compared with adult humans (1, 2). Variation in the infant microbiome exists at both the species and strain level (3, 4). During the first one to two years of life, the gut microbiome of infants begins to converge upon an adult-like state (2, 5). However, aberrations in this process may contribute to diseases such as type 1 and 2 diabetes, irritable bowel disease, and necrotizing enterocolitis (NEC) in premature infants (6-11). Because establishment of the microbiome is a key driver of immune system development, changes in the process of colonization may have life-long implications, even if they do not result in drastically different microbiome composition later in life (12, 13).

71    Infants born prematurely have lower diversity microbial communities compared

72    with full term infants, and are susceptible to life-threatening diseases such as NEC (4, 14-

73    17). While it has long been thought that bacterial infection may contribute to NEC

74    pathogenesis, strain-resolved microbial community analysis has not identified a single

75    pathogen that is responsible for the disease (3). However, it is still likely that microbial

76    communities play an important role, with the context-dependent metabolism of specific

77    strains potentially critical to infant health and disease. Recent studies have applied

78    proteomics and metabolomics to premature infant gut microbiomes to measure functional

79    profiles in healthy premature infants and those that went on to develop NEC (18, 19).

80    These studies reported temporal variation in the infant proteome and identified

81    metabolites associated with NEC. However, further study is required to better understand

82    the range of functional and developmental patterns during the microbial colonization

83    process.

84    To investigate microbial community assembly, and how microbes modulate their

85    metabolism and replication rate during colonization, we conducted a combined

86    metagenomics and metaproteomics study of the microbiome of both healthy premature

87    infants and infants that went on to develop NEC. Microbiome samples were collected

88    during the first three months of life with the goal of measuring the physiological changes

89    of dominant and ubiquitous bacterial species. Genomes assembled from metagenomes

90    enabled analysis of microbial community membership, and tracking of both community

91    composition and replication rates over time. The availability of genome sequences made

92    it possible to map protein abundance measurements to bacterial species and strains.

93    Microbial communities were clustered into distinct types in order to provide context for

94    proteomics analyses. Statistical analyses showed that, while species and strain-specific

95    proteomic profiles correlated with overall community composition, the proteomes of

96    members of the same species and strain were largely infant-specific. These analyses also

97    show that bacterial proteome features are correlated with infant development, health

98    status, and antibiotics administration.

99    **Results**

100   *Metagenome sequencing and genome binning*

101   In order to study the developing gut microbiome, stool samples were collected during the

102   first three months of life for 35 infants born prematurely and housed in the neonatal

103   intensive care unit at Magee-Womens Hospital at the University of Pittsburgh Medical

104   Center. Two of the infants in the study cohort developed sepsis (N1_017 and N1_019)

105   and 14 infants developed necrotizing enterocolitis (NEC; **Table 1**). To study the gut

106   microbiome, we analyzed 1,149 Gbp of DNA sequences generated by our laboratory (3,

107   4, 20). These sequences were from 343 metagenomes (average of 3.3 Gbp sequencing per

108   sample; **Supplemental Figure 1 and Supplemental File 1a**). Metagenomes were

109   assembled into 6.79 Gbp of scaffolds ≥1 Kbp that represented 92% of all sequenced

110   DNA.

111        Scaffolds assembled from metagenomes were grouped into 3,643 bins, 1,457 of

112   which were ≥50% complete with ≤5% contamination; **Supplemental Figure 2**,

113   **Supplemental File 2**). These genomes were assigned to 270 groups approximating

114   different bacterial sub-species based on sharing ≥98% average nucleotide identity (ANI)

115   (**Supplemental File 1b**). These genomes account for 91% of the total sequencing.

5

116    Genomes suitable for iRep replication rate analysis (≥75% complete with ≤175

117    fragments/Mbp and ≤5% contamination) were available for 193 genome clusters (21).

118    ***Protein quantification by metaproteomics***

119    Across all metagenomes, 5,233,047 proteins were predicted, 897,520 of which were from

120    a non-redundant set of representative genomes clustered at 98% ANI. Proteins clustered

121    into 121,746 putative families (**Supplemental File 3**). Metaproteomics measurements

122    were conducted on 87 metagenome-matched samples that spanned 16 infants, six of

123    which developed NEC and one of which was diagnosed with sepsis (N1_019;

124    **Supplemental Figure 1**). Conducting metagenomics and metaproteomics on the same

125    samples was critical for obtaining an appropriate database for matching peptides to

126    proteins. On average, 71,676 unique bacterial spectral counts were detected per sample,

127    with an average of 33% of predicted bacterial proteins identified (**Supplemental Figure**

128    **1**, **Supplemental File 1b, and Supplemental File 4**).

129    ***Premature infants were colonized by genetically similar organisms, and microbial***

130    ***communities clustered into seven primary types***

131    The majority of infants were colonized by *Enterococcus faecalis*, *Klebsiella pneumoniae*,

132    and *Staphylococcus epidermidis* (**Figure 1a,b**). Overall, infants that developed NEC were

133    colonized by organisms genetically similar to those colonizing other infants, and most

134    genotypes were seen in only one infant. No individual species was strongly associated

135    with NEC (**Supplemental Figure 3**).

136         The premature infant microbiome was found to be highly variable. In some cases,

137    samples collected from an infant at subsequent time points were as different from earlier

138    samples as those collected from other infants (**Figure 1c**). Communities were clustered

6

139    based on species membership and abundance in order to identify microbial consortia

140    common during the colonization process. In order to account for both genomic

141    differences and organism abundance, clustering was conducted based on weighted

142    UniFrac distances, where the tree used for calculating UniFrac was constructed using

143    genome ANI. Nine distinct community types were identified, seven of which were

144    comprised of samples collected from multiple infants and were thus considered primary

145    types  (**Figure 1d, Supplemental Figure 4, and Supplemental Figure 5**). Each

146    community type is characterized by the dominance of different community members

147    (**Supplemental Figure 6**). Microbiomes from different infants clustered into the same

148    community type, and the microbiome of individual infants was found to switch types,

149    sometimes multiple times, during the colonization process (**Figure 2**). Although infants

150    shared community types, overall colonization patterns were not replicated across infants.

151    Microbiomes associated with infants that did and did not go on to develop NEC were

152    often classified in the same community type. In some cases, switches preceded onset of

153    NEC, but no type or switch could explain all cases of NEC.

154    *Microbial replication rates and proteins*

155    iRep is a newly-developed method that enables measurement of bacterial replication rates

156    based on metagenome sequencing data when high-quality draft genome sequences are

157    available (21). We applied the iRep method using genomes recovered from metagenomes

158    sequenced for each infant in the study, and quantified 1,328 iRep replication rates from

159    330 samples. Sample clustering was conducted based on community iRep profiles,

160    identifying nine distinct iRep types that were correlated with community type (Mantel

161    test p-value = $1 \times 10^{-3}$, **Figure 2a**). Likewise, analysis of protein family abundance

162    clustered samples into four distinct proteome types, which also correlated with

163    community type (Mantel test p-value = 1 x $10^{-3}$, **Figure 2b**). Interestingly, there are

164    several cases in which iRep and/or proteome type switched when community type was

165    constant, or when community type switched but iRep and/or proteome type remained

166    constant.

167    *Microbiome development*

168    Peptide spectral counts were matched to infant-specific databases containing both human

169    and microbial proteins. This allowed for the relative proportions of human and microbial

170    proteins to be determined for each time point. Samples are dominated by human proteins

171    during the first 10 days of life (DOL), and then microbial proteins become dominant

172    around DOL 18. Ratios of human versus bacterial protein abundances show that the

173    premature infant gut microbiome is established over a period of approximately two weeks

174    (**Figure 3a**).

175         The presence of multiple data types (microbial community abundance and iRep,

176    microbial community proteome composition, and human proteome composition) enabled

177    tracking of various aspects of human and microbiome development during the first

178    months of life (**Figure 3b,c**). All measurements from an infant were stable within the

179    time scale of a week, but diverged over time. Interestingly, communities from different

180    infants neither converged nor diverged over time in terms of similarity based on three of

181    these five metrics. However, we observed that human proteome measurements and

182    microbial protein family abundances from different infants became increasingly different

183    when samples with time separations of greater than three weeks were compared. Overall,

184    the microbial proteome was more variable (higher variance) than community composition

185 (**Figure 3d,e**). After approximately two weeks, both microbial community abundance

186 and proteome measurements collected from the same infant became as different from

187 each other as samples collected from other infants.

188       The majority of human and microbiome features recorded in our analyses were

189 correlated with one another (**Figure 3f**). However, an exception is that microbial

190 community abundance and iRep were not correlated with human proteome composition

191 (Mantel test p-value >0.01). This is interesting in that it shows that there is no strong

192 connection between the overall human proteome and either the composition or replication

193 activity of the microbiome.

194       As shown in **Figure 3g**, microbial features were also correlated with a variety of

195 infant factors, including infant health and development (gestational age and weight), as

196 well as antibiotics administration (Mantel or permutational analysis of variance,

197 PERMANOVA, p-value ≤0.01). Notably, whether or not an infant developed NEC

198 (condition) correlated with several microbiome factors (infant genome inventory, and

199 both community composition and iRep), but not with proteome measurements. However,

200 these correlations were in part driven by antibiotics, as only iRep was correlated with

201 infant condition when excluding samples collected during or within five days of

202 antibiotics administration. Regardless of the influence of antibiotics on the microbiome,

203 microbial responses to treatment likely impact infant health.

204 ***Different species expressed varying amount of their proteome in the infant gut***

205 Microbes present in the gut environment are not expected to express their complete

206 complement of proteins at all times. In order to investigate the extent of proteome

207 expression for different bacteria, we compared depth of proteome sampling for each

208   organism to the percent of the predicted proteome that could be detected (**Figure 4**). The

209   median proteome detection across all samples was 11%, but this was largely due to low

210   sampling depth. Higher depth of proteome sampling corresponded with detection of a

211   larger fraction of the predicted proteins. The median percent of the proteome detected for

212   organisms with the best detection in each sample was 31% (max. 48%). For several

213   frequently detected colonists, including *Klebsiella pneumoniae*, *Klebsiella oxytoca*, and

214   members of the genus *Enterobacter,* maximum proteome expression was ~50%.

215   However, *Propionibacterium sp.*, *Anaerococcus vaginalis*, and members of the genus

216   *Bifidobacterium* expressed a greater proportion of their encoded genes than other

217   organisms. We infer that these bacteria may be specifically adapted to environments and

218   resource availability within the infant gut, whereas other bacteria may maintain capacities

219   that enable adaption to other environments.

220   ***Members of the same bacterial species replicated at different rates during colonization***

221   Across all infants, *Streptococcus agalactiae*, *Pseudomonas aeruginosa*, *Klebsiella*

222   *pneumoniae*, and members of the genera *Veillonella* and *Clostridium* exhibited some of

223   the highest replication rates (**Supplemental File 1c**). iRep values for organisms sampled

224   in this cohort during or immediately after antibiotics administration were not significantly

225   different from those at other time points (**Figure 5a**). This indicates that populations

226   present after antibiotics administration are both resistant to antibiotics and are continuing

227   to replicate. Members of several species were replicating quickly during or immediately

228   following antibiotic treatment (*Veillonella sp.*, *Streptococcus agalactiae*, *Finegoldia*

229   *magna*, and others). However, we did not detect overall higher iRep values following

230   antibiotics administration, although this was reported previously (21). Most species were

10

231  found only to be replicating in the absence of antibiotics, consistent with their

232  susceptibility to the treatment.

233  ***Species-specific proteomic profiles are associated with infant and microbiome features***

234  Relative protein abundance levels were determined for each genome and tracked across

235  samples. This identified population-specific proteome profiles and enabled us to test for

236  correlations with various human and microbial properties (**Figure 3h**, **Supplemental**

237  **Figure 1**, **and Supplemental File 1d**). *Veillonella spp., Klebsiella pneumoniae,*

238  *Escherichia coli,* and *Propionibacterium sp.* all had infant-specific profiles

239  (PERMANOVA p-value ≤0.01), indicating that although similar organisms are

240  colonizing different infants, each population is expressing a different complement of

241  proteins. *K. pneumoniae* and *Veillonella spp.* proteomes also correlated with community

242  type, as did the *Bifidobacterium breve* proteome (Mantel test p-value ≤0.01), showing

243  that populations respond to overall microbial community context. Interestingly, both

244  *Enterococcus faecalis* and *Propionibacterium sp.* exhibited proteomes that were also

245  correlated with infant development, and the *K. pneumoniae* proteome correlated with

246  both iRep and infant health. Although overall microbial proteome correlated with

247  antibiotics administration, species-specific proteome profiles did not; however, this may

248  be due to a lack of available data for the same species in multiple samples with and

249  without antibiotics.

250       Because of the existence of 35 samples in which ≥10% of the *K. pneumoniae*

251  proteome could be detected (max. = 38%, median = 25%), correlations between

252  individual protein abundances and iRep could be determined. Amongst proteins

253  positively correlated with iRep were a transcriptional regulator (LysR), proteins involved

11

254    in cell wall biogenesis, and ribosomal proteins (Pearson ≥0.5, q-value ≤0.01, observed in

255    ≥15 samples; **Supplemental File 1e**).

256    ***Infants were colonized by different strains with distinct proteomes***

257    The finding that *K. pneumoniae, E. coli, Propionibacterium sp.*, and *Veillonella spp.* have

258    infant-specific proteomes raised the question of whether or not each infant was being

259    colonized by different strains. Genome sequences ≥50% complete with ≤5%

260    contamination that were assembled for each species from each infant were compared with

261    one another, and hierarchical clustering conducted on pairwise ANI values was used to

262    delineate strains (**Supplemental Figure 7**). Clustering showed that in most cases each

263    infant was indeed colonized by distinct strains, which proteomics analysis showed are

264    functionally distinct. However, there were a few notable exceptions. Twin infants

265    N2_069 and N2_070, as well as infant N1_003 were all colonized by the same strain of *K.*

266    *pneumoniae*. The proteomic profiles for the strains colonizing N2_069 and N2_070 were

267    more similar to one another than they were to profiles recovered from other strains;

268    however, they were still distinguishable (**Figure 6**). Likewise, the same strain of

269    *Propionibacterium sp.* colonized twin infants N2_038 and N2_039. As with shared

270    strains of *K. pneumoniae*, their functional profiles clustered together but were still

271    distinguishable from one another (**Supplemental Figure 8**).

272        Analysis showed that few proteins were responsible for distinguishing proteomes

273    of the same bacterial types in different infants (**Figure 6, Supplemental Figure 8, and**

274    **Supplemental File 1d**). Common amongst these were proteins involved in nucleotide,

275    amino acid, carbohydrate and lipid metabolism. Also notable were several proteins

276    produced by *K. pneumoniae* involved in central carbohydrate metabolism, and both

12

277     galactose degradation and D-galacturonate degradation, indicating different carbon

278     preferences for strains colonizing different infants (**Figure 6**). Proteins involved in

279     bacterial secretion were differentially abundant between *K. pneumonia* colonizing

280     different infants, indicating variations in secretion potential that could affect human-

281     microbe interactions. Relatedly, the abundance of proteins involved in transport of metals,

282     ions, citrate, and several sugars also differed between infants.

283     ***Low microbiome diversity was associated with both antibiotics administration and NEC***

284     Microbiome diversity was lower during or within five days of antibiotics administration

285     compared with other time points (Mann-Whitney U test, MW, p-value = $2.6 \times 10^{-9}$;

286     **Figure 7a**), and the microbiome of infants that developed NEC was typically less

287     diverse compared with healthy infants (MW p-value = $4 \times 10^{-4}$; **Figure 7b**). However, the

288     difference in diversity between healthy and NEC infants was driven by the fact that NEC

289     infants more frequently receive antibiotics (**Figure 2**). When comparing within either

290     periods with or without antibiotics, microbiome diversity for healthy and NEC infants

291     (pre-NEC diagnosis) was indistinguishable (**Figure 7c**). When excluding antibiotics

292     samples, both groups of infants had higher diversity microbial communities later in

293     development (post GA + DOL 220; **Figure 7d**).

294     ***Microbial community composition was correlated with infant health***

295     Premature infants that developed NEC had different microbial community abundance

296     profiles (PERMANOVA p-value = $3 \times 10^{-3}$; **Supplemental Figure 4g**). Interestingly,

297     there were a variety of species detected in healthy infants, but never detected in those that

298     developed NEC; however, the opposite was not true. It should be noted that species

299     unique to healthy infants were not consistently detected. No species identified five days

300    prior to NEC diagnosis showed a significant difference in abundance, or was unique to

301    NEC infants.

302        Overall community composition was also correlated with each infant, antibiotics

303    administration, birth weight, gestational age, and gestational age corrected day of life

304    (GA + DOL; PERMANOVA or Mantel test p-value ≤0.01; **Supplemental Figure 4**).

305    Several species were more abundant members of communities associated with infants

306    that developed NEC: *Enterobacter sp.*, *Propionibacterium sp.*, and *Peptostreptococcus sp.*

307    (edgeR q-value ≤0.01 after excluding samples collected within five days of antibiotics

308    administration; **Supplemental File 1f**). *Vellonella sp.* replicated faster in NEC infants,

309    while several groups of organisms were replicating faster in control infants, including

310    members of the genera *Anaerococcus, Klebsiella, Actinomyces, Eggerthella,*

311    *Streptococcus, Clostridiales,* and *Bifidobacterium*  (MW p-value ≤0.01 after excluding

312    samples collected within five days of antibiotics administration; **Supplemental File 1c**).

313    Several different species were active in control infants, but were not detected in infants

314    that went on to develop NEC. However, combined iRep values collected from infants that

315    did and did not go on to develop NEC were not statistically different, even when

316    considering only samples collected within the five days prior to NEC diagnosis (**Figure**

317    **5b**).

318    ***Microbial proteins associated with proteome type, antibiotics administration, and NEC***

319    As described above, we used protein abundance patterns to cluster microbial community

320    proteomes into functionally distinct proteome types. Statistical analysis identified 3,085

321    differentially abundant proteins distinguish proteome types (edgeR q-value ≤ 0.01;

322    **Supplemental File 1g**). Of these, 461 were found to distinguish only one proteome type

14

323    from all others. Notable amongst all of these proteins were those involved in central

324    carbohydrate metabolism and energy metabolism (**Supplemental Figure 9**). Proteome

325    types differ in terms of the amount and variety of carbon degradation enzymes, as well as

326    the propensity for aerobic versus anaerobic respiration (based on the abundance of

327    oxidases and reductases).

328         Samples collected during antibiotics treatment were enriched in 56 different

329    proteins (identified in more than one treated infant, edgeR q-value ≤0.01; **Supplemental**

330    **File 1g**). Amongst these proteins were those involved in secretion, transcription, and

331    DNA degradation. Along with iRep results, the findings indicate that a subset of

332    organisms remain active in the presence of antibiotics.

333         Although overall community proteome abundance profiles were not correlated

334    with NEC, microbial proteins from 160 different protein families, many with no known

335    function, were more abundant in samples from infants that went on to develop NEC

336    (identified in more than one NEC infant, edgeR q-value ≤0.01; **Supplemental File 1g**).

337    Annotated proteins were dominantly involved in transport of ions, metals, and other

338    substrates, iron acquisition, and both motility and chemotaxis. Among proteins

339    responsible for iron scavenging was subunit E of enterobactin synthase, a high-affinity

340    siderophore involved in iron acquisition, which is often used by pathogenic organisms.

341    Also more abundant was outer membrane receptor FepA, which is involved in

342    transporting iron bound by extracellular enterobactin. Subunit F of enterobactin synthase

343    was also identified in NEC infants, as were an iron-enterobactin ABC transporter

344    substrate-binding protein, and an enterobactin esterase. The abundance of this protein

345    suggests a possible role for iron acquisition by organisms that may contribute to disease

346  onset. Interestingly, 21 *K. pneumoniae* proteins were correlated with NEC, including a

347  ferrous iron transporter (family 2834) that was 3.9-fold more abundant in two infants that

348  developed NEC. The abundance of this protein was also correlated with infant, proteome

349  type, community type, and antibiotics administration.

350  **Discussion**

351  Most studies to date have focused on the composition of the gut microbiome, typically at

352  the low-resolution afforded by 16S rRNA gene amplicon methods. We used genome-

353  resolved time-series metagenomics in conjunction with iRep replication rate and

354  metaproteomics measurements to obtain a more comprehensive view of the colonization

355  process. The dataset included information about the gut colonization trajectories of both

356  healthy infants and infants that went on to develop NEC, enabling exploration of

357  microbiome variability, at both the community composition and organism functional

358  levels.

359          Microbial communities were classified into types based on the mixture of

360  organisms present. Interestingly, most types occurred in multiple infants, a result that

361  indicates the tendency of gut colonizing bacteria to form networks of interaction, possibly

362  based on metabolic complementarity. An important factor determining the community

363  type present may be the specific organisms that are introduced, and the extent to which

364  they are able to colonize. Other factors that may dictate the community type include

365  human genetic selection, diet, and antibiotics administration. Within a single infant,

366  community types often switched several times over the observation period. Given the

367  lack of evidence for consistent transitions from one type to another across multiple

368  infants, the high degree of variation in iRep replication rates observed for members of the

16

369    same species, and a lack of convergence of communities in different infants, we conclude

370    that colonization is a chaotic process.

371         Overall microbial physiology, as measured by whole proteome abundance

372    patterns, was more dynamic than community composition. Thus, metagenomics-enabled

373    proteomic analyses indicate functional flexibility that does not depend on addition or loss

374    of organisms. Shifts in the importance of specific pathways or metabolisms with

375    environmental conditions would not be apparent in studies that only use organism

376    identification or metabolic potential predictions.

377         It is possible that onset of NEC is due to fast growth rates of potential pathogens

378    within communities that are imbalanced due to low species richness, ultimately resulting

379    in overgrowth by a pathogen. For this reason, we compared microbial community

380    diversity and composition, growth rates, and metabolic features in infants that did and did

381    not develop NEC. A clear finding of this study, and evident from prior research (17), is

382    that microbial communities associated with infants that develop NEC are of lower

383    diversity compared with control infants. However, this was due to the frequency of

384    antibiotics administration for NEC infants. Regardless of the cause of the lower-diversity

385    communities, microbial activities throughout the colonization process, including during

386    periods of antibiotics administration, are likely important to infant health.

387         Several different species had higher relative abundance in infants that developed

388    NEC, but none of these species were consistently associated with the disease. The

389    correlation could be the consequence of the loss of other organisms from the community

390    rather than their higher absolute abundance. Interestingly, *Veillonella spp.* were  found to

17

391    replicate more quickly in NEC versus control infants. This may be medically important,

392    but additional examples are needed to establish a link between rapid growth and NEC.

393          Surprisingly, whether or not an infant developed NEC was not correlated with

394    overall proteome composition. However, there were specific proteins that were associated

395    with NEC, notably several involved in iron scavenging. Given that this is an important

396    process often associated with pathogenesis, it is possible that increased activity of iron

397    scavenging pathways could contribute to organism proliferation and onset of NEC. In

398    addition, the *Klebsiella pneumoniae* proteome was correlated with NEC, including a

399    protein involved in transport of iron. This is intriguing considering the prior finding that

400    supplementation of lactoferrin, an abundant breast milk protein involved in modulating

401    iron levels in the gut, decreases risk of developing necrotizing enterocolitis (22, 23).

402    Overall, these findings indicate that fine-scale, species-specific proteins are important for

403    understanding disease onset. Although the microbial community, and specific microbial

404    proteins were correlated with NEC, no individual organism or protein was significantly

405    more abundant in all cases. This finding supports the hypothesis that NEC is a

406    multifaceted disease with multiple routes that lead to onset.

407          Although species-specific proteome profiles were correlated with community

408    composition, they were largely infant specific. This is an interesting observation because

409    it implies a feedback between human physiological conditions in the gut, which likely

410    vary substantially from infant to infant and over time, and microbiome function.

18

411    **Methods**

412    *Sample collection and metagenome sequencing*

413    Samples were collected, processed for metagenome sequencing, and sequenced as part of

414    three prior studies (**accession numbers in Supplemental File 1a**) (3, 4, 20). Stool

415    samples were collected from infants and stored at −80°C. DNA was extracted from

416    frozen fecal samples using the MO BIO PowerSoil DNA Isolation Kit, with

417    modifications (4). DNA libraries were sequenced on an Illumina HiSeq for 100 or 150

418    cycles (Illumina, San Diego, CA). All samples were collected with parental consent.

419    *Metagenome assembly and genome binning*

420    We re-assembled and analyzed metagenomes generated as part of a prior study, referred

421    to as NIH1 (4). The data were processed in a manner consistent with the two other prior

422    studies analyzed, referred to as NIH2 (20) and NIH3 (3). All raw sequencing reads were

423    trimmed using Sickle (https://github.com/najoshi/sickle). Each metagenome was

424    assembled separately using IDBA_UD (24). Open reading frames (ORFs) were predicted

425    using Prodigal (25) with the option to run in metagenome mode. Predicted protein

426    sequences were annotated based on USEARCH (–ublast) (26, 27) searches against

427    UniProt (28), UniRef100 (29), and KEGG (30, 31). Scaffold coverage was calculated by

428    mapping reads to the assembly using Bowtie2 (32) with default parameters for paired

429    reads.

430         Scaffolds from NIH1 infants were binned to genomes using Emergent Self-

431    Organizing Maps (ESOMs) generated based on time-series abundance profiles (15, 33).

432    Reads from every sample were mapped independently to every assembly using SNAP

433    (34), and the resulting coverage data were combined. Coverage was calculated over non-

19

434 overlapping three Kbp windows. Coverage values were normalized first by sample, and

435 then the values for each scaffold fragment were normalized from zero to one. Combining

436 coverage data from scaffolds assembled from different samples prior to normalization

437 made it possible to generate a single ESOM map for binning genomes assembled

438 independently from each sample. ESOMs were trained for ten epochs using the Somoclu

439 algorithm (35) with the option to initialize the codebook using Principal Component

440 Analysis (PCA). Genomes were binned by manually selecting data points on the ESOM

441 map using Databionics ESOM Tools (36). Binning was aided by coloring scaffold

442 fragments on the map based on BLAST (37) hits to the genomes assembled in the prior

443 study.

444  As part of the NIH2 and NIH3 studies, scaffolds were binned based on their GC

445 content, DNA sequence coverage, and taxonomic affiliation using ggKbase tools

446 (ggkbase.berkeley.edu). Genome bins from all three datasets were classified based on the

447 consensus of taxonomic assignments for predicted protein sequences. Genome

448 completeness and contamination were estimated for all genomes using CheckM with the

449 taxonomy_wf option (38). Genomes with extra single copy genes, but with ≤175

450 fragments/Mbp (normalized for contamination) that were estimated to be ≥75% complete

451 were manually curated based on scaffold GC content and coverage.

452 ***Clustering genomes into sub-species groups***

453 Genomes were clustered into sub-species groups based on sharing ≥98% average

454 nucleotide identity (ANI), as estimated by MASH (39). Representative genomes were

455 selected for each cluster as the largest genome with the highest expected completeness

456    and smallest amount of contamination. Genomes were classified based on the lowest

457    possible consensus of taxonomic assignments for predicted protein sequences.

458        Taxonomic assignments for representative genomes were checked manually based

459    on hits to ribosomal protein S3, or visual inspection of protein taxonomic assignments. In

460    order to identify cases in which the same bacterial strain was present in multiple samples,

461    sub-species groups were further analyzed with the ANIm algorithm (40) implemented in

462    dRep (41).

463    *Measuring microbial community abundance and replication rates*

464    In order to achieve accurate abundance and replication rate measurements from read

465    mapping, databases of representative genomes were created for each sample. Each

466    database was constructed in order to include a representative genome from important sub-

467    species groups. Priority was given to high-quality draft genome sequences reconstructed

468    from the same sample. Genomes were classified as high-quality draft based on the

469    requirements for iRep replication rate analysis

470    (https://github.com/christophertbrown/iRep): ≥75% complete, ≤2.5% contamination, and

471    ≤175 scaffolds per Mbp of sequence (21). Genomes were selected to represent sub-

472    species groups using the following priority scheme: 1) high-quality draft genome

473    assembled from the same sample, 2) high-quality draft genome from the same infant, 3)

474    high-quality draft genome representative of sub-species group from any infant (if group

475    had ≥5 representatives), 4) best genome from infant (if a genome was available). iRep

476    was conducted using reads that mapped to genome sequences with ≤1 mismatch per read

477    sequence. In cases where iRep values were ≥3, coverage plots were inspected and values

478    were removed if there was evidence of strain variation.

21

479    We considered bacterial sub-species to be present in a sample if ≥97% of the

480    genome was covered by an average of ≥2 reads. Abundance and iRep measurements were

481    compared across samples by linking sample-specific representative genomes to sub-

482    species groups. Relative abundance measurements for each sub-species group were

483    calculated by converting DNA sequencing coverage values to a percentage. UniFrac (42)

484    analysis was conducted based on rarefied abundance data and a tree constructed based on

485    pairwise genome ANI values measured using MASH (-ms 5000000).

486    ***Metaproteomics analysis***

487    Metaproteomics sequencing was conducted on 0.3 g of stool as previously described (18).

488    Each sample was suspended in 10 mL cold phosphate buffered saline. Samples were

489    filtered through a 20 µm size filter to enrich for microbial cells and proteins. Microbial

490    cells were collected by centrifugation, boiled in 4% sodium dodecyl sulfate for 5 minutes,

491    and sonicated to lyse cells. The resulting protein extract was precipitated with 20%

492    trichloroacetic acid at -80°C overnight. The protein pellet was washed with ice-cold

493    acetone, solubilized in 8 M urea, reduced with 5 mM dithiothreitol, and cysteines were

494    blocked with 20 mM iodoacetamide. Then sequencing grade trypsin was used to digest

495    the proteins into peptides. Proteolyzed peptides were then salted and acidified by

496    adjusting the sample to 200 mM NaCl, 0.1% formic acid, followed by filtering through a

497    10 kDa cutoff spin column filter to collect tryptic peptides.

498    Peptides were quantified by BCA assay and 50 µg peptides of each sample were

499    analyzed via two-dimensional nanospray LC-MS/MS system on an LTQ-Orbitrap Elite

500    mass spectrometer (Thermo Scientific). Each peptide mixture was loaded onto a biphasic

501    back column containing both strong-cation exchange and reverse phase resins (C18). As

22

502     previously described, loaded peptides were separated and analyzed using a 11-salt-pusle

503     MudPIT protocol over a 22-h period (43). Mass spectra were acquired in a data-

504     dependent mode with following parameters: full scans were acquired at 30 k resolution (1

505     microscan) in the Orbitrap, followed by CID fragmentation of the 20 most abundant ions

506     (1 microscan). Charge state screening and monoisotopic precursor selection were enabled.

507     Unassigned charge and charge state +1 were rejected. Dynamic exclusion was enabled

508     with a mass exclusion width of 10 ppm and exclusion duration of 30 seconds. Two

509     technical replicates were conducted for each sample.

510            Protein databases were generated for each infant from protein sequences predicted

511     from assembled metagenomes. The database also included human protein sequences

512     (NCBI Refseq_2011), common contaminants, and reverse protein sequences, which were

513     used to control the false discovery rate (FDR). Collected MS/MS spectra were matched to

514     peptides using MyriMatch v2.1 (44), filtered, and assembled into proteins using IDPicker

515     v3.0 (45). All searches included the following peptide modifications: a static cysteine

516     modification (+57.02 Da), an N-terminal dynamic carbamylation modification (+43.00

517     Da), and a dynamic oxidation modification (+15.99). A maximum 2% peptide spectrum

518     match level FDR and a minimum of two distinct peptides per protein were applied to

519     achieve confident peptide identifications (FDR <1%). To alleviate the ambiguity

520     associated with shared peptides, proteins were clustered into protein groups by 100%

521     identity for microbial proteins and 90% amino acid sequence identity for human proteins

522     using USEARCH (26). Spectral counts were balanced between shared proteins, and

523     proteins were considered to be present if ≥2 unique peptides were identified.

23

*Identification of putative protein families*

Putative protein families were identified in order to track the presence and abundance of different protein types across samples. ORFs were first pre-clustered at 95% identity using USEARCH (-cluster_smallmem -target_cov 0.50 -query_cov 0.95 -id 0.95), and then all-versus-all protein searches were conducted (–ublast -evalue 10e-10 -strand both). Protein families were delineated from within the all-versus-all network graph using the MCL clustering algorithm (-I 2 -te 10) (46). The most common annotation observed across all protein sequences in the group was selected as the annotation for the putative protein family. Proteins were also grouped based on sharing 97% amino acid identity using USEARCH (-cluster_smallmem -target_cov 0.50 -query_cov 0.95 -id 0.97).

*Tracking human and bacterial protein abundances*

Human and bacterial protein abundances were normalized using the weighted trimmed mean method from EdgeR (47). Species-specific proteomic profiles were normalized as the percent of total balanced spectral counts.

*Sample clustering and statistical analyses*

Sample clustering was conducted based on microbial community abundance and iRep profiles, and bacterial protein family abundance profiles. In each case, the number of clusters was determined using the gap statistic (48), and then samples were grouped into the appropriate number of clusters using hierarchical clustering (average linkage method). Microbial community data was clustered based on weighted UniFrac distances, and protein data using Bray-Curtis distance. EdgeR was used to calculate statistically significant differences between conditions using quasi-likelihood linear modeling (glmQLFTest).

**Acknowledgements**

**Author Contributions**

MJM oversaw sample collection, RB collected all samples and managed metadata, and BF coordinated sample processing for DNA sequencing and proteomics analysis. CTB and MRO assembled and annotated the metagenome data. CTB and JFB carried out the genome binning and curation. CTB conducted the microbial community time series abundance and iRep analyses. WX and RLH generated the proteomics data, which was analyzed by CTB. CTB, MRO, and BCT provided bioinformatics support. CTB and JFB wrote the paper, and all authors provided input to the final text.

**Disclosure Declaration**

The authors declare no competing financial interests.

**Figures**

**Figure 1 | Premature infant gut microbial communities associated into seven primary types.** Genomes reconstructed from metagenomes were clustered into sub-species groups based on sharing 98% average nucleotide identity (ANI). **a,** The number of genomes assigned to each group and **b,** the number of infants with a reconstructed genome from the group. Shown are groups comprised of five or more genomes. **c,** Pairwise weighted UniFrac distances calculated between all microbiome samples based

568    on genome sequence ANI and abundance. **d,** PCoA clustering of samples based on

569    weighted UniFrac distances. Samples are colored based on community type assignment.

570    **Figure 2 | Microbial colonization patterns for preterm infants.** Samples were

571    clustered into types based on microbial community composition ("community type"),

572    bacterial iRep profiles ("iRep type"), and overall bacterial proteome composition

573    ("proteome type"). Microbial community type is shown along with iRep **(a)** and

574    proteome **(b)** types. Infants are arranged based on hierarchical clustering of unweighted

575    UniFrac distances calculated based on the set of genomes recovered from each infant

576    (**Supplemental Figure 3**). Antibiotics administration is indicated with pink bars and

577    NEC diagnoses with red bars. DOL stands for day of life.

578    **Figure 3 | Microbiome stability and correlations. a,** The relative contribution of human

579    and bacterial proteins to overall proteome composition during development of the

580    premature infant gut. **b,** Similarity measurements for microbiomes sampled either from

581    the same infant or **c,** from different infants. Comparison of similarity measurements

582    calculated between samples collected either form the same or different infants based

583    either on weighted microbial community UniFrac (**d**), or weighted microbial proteome

584    BrayCurtis (**e**) measurements. Human proteome and microbial community correlations

585    calculated between one another (**f**), with infant metadata (**g**), and determined based on

586    microbial species (**h**). Shown are PERMANOVA or Mantel test p-values (**f-h**). Microbial

587    proteome "family" refers to protein family analysis, and "group" refers to analysis of

588    proteins clustered at 97% amino acid identity.

589    **Figure 4 | Proteome detection for species colonizing premature infants.** Depth of

590    proteome sampling for organisms in each sample is compared against the percent of

26

591   predicted proteins that could be detected. Data point sizes and histograms are scaled

592   based on organism abundance as determined by metagenome sequencing.

593   **Figure 5 | Replication rates for bacteria colonizing premature infants. a,** Replication

594   rates for bacteria sampled during periods with or without antibiotics administration and **b,**

595   associated with infants that did and did not go on to develop NEC. Statistically

596   significant differences between replication rates observed for individual species under

597   different conditions are indicated with an asterisk (MW p-value ≤0.01). Shown are all

598   species with at least five observations.

599   **Figure 6 |** *Klebsiella pneumonia* **proteins with infant-specific abundance profiles.**

600   Hierarchical clustering was conducted on all *K. pneumonia* protein families, showing that

601   strains colonizing different infants have distinct proteomic profiles. Infant and species

602   metadata are shown for each sample. Metadata significantly correlated with the *K.*

603   *pneumonia* proteome are indicated with an asterisk (PERMANOVA or Mantel test p-

604   value ≤0.01). Protein families that correlated with at least one infant are shown in the

605   heatmap (edgeR q-value ≤0.01). Samples colonized by the same *K. pneumonia* strain are

606   shown with red text.

607   **Figure 7 | Microbial community diversity.** Shannon diversity measurements for

608   microbial communities associated with infants during periods with or without antibiotics

609   administration (**a**), and between infants that did and did not go on to develop NEC (**b-d**).

610   Significant differences are indicated with an asterisk (MW p-value ≤0.01). "Early"

611   samples were collected prior to GA + DOL 220. Samples collected after NEC diagnosis

612   were excluded from **c** and **d**.

613 **Tables**

614 **Table 1 | Infant medical information.**

| infant | campaign | sex | delivery | mult. gest. | gestational age (weeks) | birth weight (g) | feeding | condition | NEC diagnosis (DOL) |
|---|---|---|---|---|---|---|---|---|---|
| N1_003 | NIH1 | F | C-section | Single | 26 | 822 | Breast | control | n/a |
| N1_004 | NIH1 | F | C-section | N1_005 | 32 | 1450 | Formula | control | n/a |
| N1_008 | NIH1 | F | Vaginal | Single | 32 | 1230 | Formula | NEC | 9 |
| N1_009 | NIH1 | M | C-section | Single | 29 | 1820 | Combination | control | n/a |
| N1_011 | NIH1 | M | C-section | N1_012 | 26 | 523 | Combination | NEC | 34, 62 |
| N1_014 | NIH1 | M | Vaginal | Single | 32 | 2035 | Combination | control | n/a |
| N1_017 | NIH1 | F | Vaginal | Single | 26 | 748 | Combination | NEC | 11 |
| N1_018 | NIH1 | M | C-section | Single | 29 | 1133 | Combination | control | n/a |
| N1_019 | NIH1 | F | C-section | N1_020, N1_021 | 24 | 731 | Combination | control | n/a |
| N1_021 | NIH1 | F | C-section | N1_019, N1_020 | 24 | 697 | Breast | NEC | 32 |
| N1_023 | NIH1 | F | Vaginal | Single | 27 | 875 | Breast | control | n/a |
| N2_031 | NIH2 | M | C-section | Single | 26 | 773 | Formula | control | n/a |
| N2_035 | NIH2 | M | Vaginal | Single | 25 | 795 | Breast | control | n/a |
| N2_038 | NIH2 | F | C-section | N2_039 | 30 | 1381 | Combination | control | n/a |
| N2_039 | NIH2 | F | C-section | N2_038 | 30 | 1470 | Combination | NEC | 24 |
| N2_060 | NIH2 | M | C-section | Single | 30 | 1878 | Combination | control | n/a |
| N2_061 | NIH2 | M | Vaginal | Single | 28 | 1184 | Combination | NEC | 9, 34 |
| N2_064 | NIH2 | M | Vaginal | Single | 28 | 1100 | Combination | control | n/a |
| N2_065 | NIH2 | F | Vaginal | Single | 25 | 841 | Combination | control | n/a |
| N2_066 | NIH2 | F | Vaginal | Single | 28 | 1028 | Breast | control | n/a |
| N2_069 | NIH2 | M | C-section | N2_070 | 26 | 637 | Breast | NEC | 32 |
| N2_070 | NIH2 | F | C-section | N2_069 | 26 | 633 | Combination | control | n/a |
| N2_071 | NIH2 | M | C-section | Single | 25 | 754 | Combination | NEC | 31 |
| N2_088 | NIH2 | F | C-section | N2_089 | 28 | 1057 | Formula | control | n/a |
| N2_093 | NIH2 | M | C-section | Single | 26 | 924 | Breast | NEC | 12 |
| N3_172 | NIH3 | M | C-section | Single | 28 | 1250 | Breast | NEC | 37, 54 |
| N3_173 | NIH3 | M | C-section | Single | 29 | 1530 | Breast | NEC | 25 |
| N3_174 | NIH3 | F | C-section | Single | 30 | 980 | Breast | control | n/a |
| N3_175 | NIH3 | M | Vaginal | Single | 29 | 1480 | Combination | control | n/a |
| N3_176 | NIH3 | M | C-section | Single | 28 | 990 | Combination | control | n/a |
| N3_177 | NIH3 | F | Vaginal | Single | 28 | 900 | Combination | control | n/a |
| N3_178 | NIH3 | M | Vaginal | Single | 32 | 2050 | Combination | NEC | 16 |
| N3_182 | NIH3 | M | C-section | Single | 39 | 3010 | Combination | NEC | 6 |
| N3_183 | NIH3 | M | Vaginal | Single | 32 | 2410 | Combination | NEC | 11 |
| S2_010 | NIH3 | M | C-section | Single | 32 | 1810 | Combination | control | n/a |

615

616 **Supplemental Materials**

617 *Supplemental Figures*

618 **Supplemental Figure 1 | Metagenome sequencing and metaproteomics conducted on**

619 **microbiome samples collected from premature infants.** Frequency of sample

620 collection for metagenomics (**a**) and metaproteomics (**b**) based on infant day of life

28

621  (DOL). **c,** Metagenome sequencing, and **d**, the percentage of each metagenome

622  represented by assembled genome sequences ≥50% complete with ≤5% contamination. **e,**

623  The number of proteomics spectral counts that could be uniquely assigned to human or

624  bacterial proteins. **f,** The percent of predicted proteins that could be detected in each

625  sample. **g,** The percent of species-specific proteomes that could be detected for species

626  where ≥10% of the proteome could be detected in at least one sample. **h,** Histogram

627  showing the distribution of the maximum percent of the proteome detected for all species

628  present in each sample.

629  **Supplemental Figure 2 | ESOM genome binning.** Genome binning was conducted

630  based on Emergent Self-Organizing Map (ESOM) clustering of scaffolds assembled from

631  individual metagenomes. Data points represent 3 Kbp fragments of assembled scaffolds.

632  Coloring is based on the species-level assignment of reconstructed genomes. The map is

633  periodic, and red boxes indicate a single period.

634  **Supplemental Figure 3 | Infants that developed NEC and healthy controls are**

635  **colonized by genetically similar bacteria.** Presence (dark boxes) and absence (white

636  boxes) of members of bacterial sub-species in microbial communities from different

637  infants. Sub-species were identified based on sharing ≥98% genome average nucleotide

638  identity (ANI), and were determined to be present if ≥97% of the genome was covered by

639  an average of ≥2 reads. Hierarchical clustering was conducted based on unweighted

640  UniFrac distances calculated between infant genome inventories.

641  **Supplemental Figure 4 | Studied infant gut microbial communities associate into**

642  **seven primary community types. a**, Hierarchical clustering was conducted based on the

643  abundance of bacterial sub-species using weighted UniFrac distances. Microbial

29

644 community types are identified by colored boxes. Metadata are shown for each sample,

645 and indicated with an asterisk if significantly correlated with microbial community

646 abundance data (PERMANOVA or Mantel test p-value ≤0.01). **b-i**, PCoA clustering of

647 microbial communities with associated metadata: antibiotics administration (**b**), infant (**c**),

648 developmental age (**d**; number of days since conception: gestational age + day of life, GA

649 + DOL), proteome type (**e**), iRep type (**f**), infant health (**g**), days prior to NEC diagnosis

650 (**h;** DOL – NEC diagnosis), and human proteome type (**i**).

651 **Supplemental Figure 5 | Microbial community abundance and replication rate**

652 **profiles.** Relative abundance (bars) and iRep replication rate (scatter plot) values for

653 bacterial sub-species colonizing studied premature infants. The five days following

654 antibiotics administration are indicated with a color gradient.

655 **Supplemental Figure 6 | Microbial community types are distinguished by their**

656 **abundant members.** Rank abundance curves showing the average and range (95%

657 confidence interval) of relative abundance values for sub-species groups associated with

658 each community type.

659 **Supplemental Figure 7 | Hierarchical clustering of genomes for members of the**

660 **same sub-species group.** dRep results show ANI clustering of assembled genomes.

661 Genome names indicate the metagenome that each genome was assembled from (**see**

662 **Supplemental File 1b**). Clustering dendrograms show that most infants are colonized by

663 different strains.

664 **Supplemental Figure 8 | Multiple species have infant-specific proteome profiles. a**,

665 Analysis of *Veillonella spp.* genomes shows the presence of four different species. **b-e**,

666 Proteome profiles for different species colonizing premature infants. Hierarchical

30

667   clustering was conducted based on all detected protein families, and shows that strains

668   colonizing different infants typically have distinct proteomic profiles. Infant and species

669   metadata are shown for each sample. Metadata significantly correlated with the species

670   proteome are indicated with an asterisk (PERMANOVA or Mantel test p-value ≤0.01).

671   Protein families that correlated with at least one infant are shown in the heatmap (edgeR

672   q-value ≤0.01). Samples colonized by the same strain are shown with colored text.

673   **Supplemental Figure 9 | Proteome types are distinguished by the abundance of**

674   **proteins from different KEGG modules.** Hierarchical clustering of proteome types was

675   conducted based on the abundance of proteins associated with KEGG modules. The

676   relative abundance of proteins associated with each module was summed for each sample,

677   and then the average was taken across all samples associated with each proteome type.

678   *Supplemental Files*

679   **Supplemental File 1a | DNA sequencing and metaproteomics statistics.**

680   **Supplemental File 1b | Genomes reconstructed from metagenomes.**

681   **Supplemental File 1c | Species iRep replication rates and statistical analysis.**

682   **Supplemental File 1d | Species-specific microbial protein family abundance and**

683   **statistical analysis.**

684   **Supplemental File 1e | Correlation of species-specific protein family abundances**

685   **with iRep replication rates and gestational age corrected day of life (GA + DOL).**

686   **Supplemental File 1f | Species relative abundance and statistical analysis.**

687   **Supplemental File 1g | Microbial protein family abundance and statistical analysis.**

688   **Supplemental File 2 | Scaffolds binned to reconstructed genomes.**

689   **Supplemental File 3 | Proteins assigned to putative families.**

690     **Supplemental File 4 | Metaproteomics spectral counts.**

691     **References**

692     1.     **Costello EK**, **Lauber CL**, **Hamady M**, **Fierer N**, **Gordon JI**, **Knight R**. 2009.
693            Bacterial community variation in human body habitats across space and time.
694            Science **326**:1694–1697.

695     2.     **Palmer C**, **Bik EM**, **DiGiulio DB**, **Relman DA**, **Brown PO**. 2007. Development
696            of the human infant intestinal microbiota. PLoS Biol **5**:e177.

697     3.     **Raveh-Sadka T**, **Thomas BC**, **Singh A**, **Firek B**, **Brooks B**, **Castelle CJ**,
698            **Sharon I**, **Baker R**, **Good M**, **Morowitz MJ**, **Banfield JF**. 2015. Gut bacteria are
699            rarely shared by co-hospitalized premature infants, regardless of necrotizing
700            enterocolitis development. Elife **4**:–.

701     4.     **Raveh-Sadka T**, **Firek B**, **Sharon I**, **Baker R**, **Brown CT**, **Thomas BC**,
702            **Morowitz MJ**, **Banfield JF**. 2016. Evidence for persistent and shared bacterial
703            strains against a background of largely unique gut colonization in hospitalized
704            premature infants. ISME J.

705     5.     **Bokulich NA**, **Chung J**, **Battaglia T**, **Henderson N**, **Jay M**, **Li H**, **Lieber AD**,
706            **Wu F**, **Perez-Perez GI**, **Chen Y**, **Schweizer W**, **Zheng X**, **Contreras M**,
707            **Dominguez-Bello MG**, **Blaser MJ**. 2016. Antibiotics, birth mode, and diet shape
708            microbiome maturation during early life. Sci Transl Med **8**.

709     6.     **Xavier RJ**, **Podolsky DK**. 2007. Unravelling the pathogenesis of inflammatory
710            bowel disease. Nature **448**:427–434.

711     7.     **Brown CT**, **Davis-Richardson AG**, **Giongo A**, **Gano KA**, **Crabb DB**,
712            **Mukherjee N**, **Casella G**, **Drew JC**, **Ilonen J**, **Knip M**, **Hyöty H**, **Veijola R**,
713            **Simell T**, **Simell O**, **Neu J**, **Wasserfall CH**, **Schatz D**, **Atkinson MA**, **Triplett**
714            **EW**. 2011. Gut microbiome metagenomics analysis suggests a functional model
715            for the development of autoimmunity for type 1 diabetes. PLoS ONE **6**:e25792.

716     8.     **Qin J**, **Li Y**, **Cai Z**, **Li S**, **Zhu J**, **Zhang F**, **Liang S**, **Zhang W**, **Guan Y**, **Shen D**,
717            **Peng Y**, **Zhang D**, **Jie Z**, **Wu W**, **Qin Y**, **Xue W**, **Li J**, **Han L**, **Lu D**, **Wu P**, **Dai**
718            **Y**, **Sun X**, **Li Z**, **Tang A**, **Zhong S**, **Li X**, **Chen W**, **Xu R**, **Wang M**, **Feng Q**,
719            **Gong M**, **Yu J**, **Zhang Y**, **Zhang M**, **Hansen T**, **Sanchez G**, **Raes J**, **Falony G**,
720            **Okuda S**, **Almeida M**, **LeChatelier E**, **Renault P**, **Pons N**, **Batto J-M**, **Zhang Z**,
721            **Chen H**, **Yang R**, **Zheng W**, **Li S**, **Yang H**, **Wang J**, **Ehrlich SD**, **Nielsen R**,
722            **Pedersen O**, **Kristiansen K**, **Wang J**. 2012. A metagenome-wide association
723            study of gut microbiota in type 2 diabetes. Nature **490**:55–60.

724    9.    **Mshvildadze M**, **Neu J**, **Shuster J**, **Theriaque D**, **Li N**, **Mai V**. 2010. Intestinal
725          Microbial Ecology in Premature Infants Assessed with Non–Culture-Based
726          Techniques. J Pediatr **156**:20–25.

727    10.   **Mai V**, **Young CM**, **Ukhanova M**, **Wang X**, **Sun Y**, **Casella G**, **Theriaque D**, **Li
728          N**, **Sharma R**, **Hudak M**, **Neu J**. 2011. Fecal microbiota in premature infants
729          prior to necrotizing enterocolitis. PLoS ONE **6**:e20647–e20647.

730    11.   **Morrow AL**, **Lagomarcino AJ**, **Schibler KR**, **Taft DH**. 2013. Early microbial
731          and metabolomic signatures predict later onset of necrotizing enterocolitis in
732          preterm infants. Microbiome **1**:13.

733    12.   **Maslowski KM**, **Vieira AT**, **Ng A**, **Kranich J**, **Sierro F**, **Yu D**, **Schilter HC**,
734          **Rolph MS**, **Mackay F**, **Artis D**, **Xavier RJ**, **Teixeira MM**, **Mackay CR**. 2009.
735          Regulation of inflammatory responses by gut microbiota and chemoattractant
736          receptor GPR43. Nature **461**:1282–1286.

737    13.   **Lathrop SK**, **Bloom SM**, **Rao SM**, **Nutsch K**, **Lio C-W**, **Santacruz N**, **Peterson
738          DA**, **Stappenbeck TS**, **Hsieh C-S**. 2011. Peripheral education of the immune
739          system by colonic commensal microbiota. Nature **478**:250–254.

740    14.   **Neu J**, **Walker WA**. 2011. Necrotizing enterocolitis. N Engl J Med **364**:255–264.

741    15.   **Sharon I**, **Morowitz MJ**, **Thomas BC**, **Costello EK**, **Relman DA**, **Banfield JF**.
742          2012. Time series community genomics analysis reveals rapid shifts in bacterial
743          species, strains, and phage during infant gut colonization. Genome Res **23**:111–
744          120.

745    16.   **Brown CT**, **Sharon I**, **Thomas BC**, **Castelle CJ**, **Morowitz MJ**, **Banfield JF**.
746          2013. Genome resolved analysis of a premature infant gut microbial community
747          reveals a *Varibaculum cambriense* genome and a shift towards fermentation-based
748          metabolism during the third week of life. Microbiome **1**:30.

749    17.   **Pammi M**, **Cope J**, **Tarr PI**, **Warner BB**, **Morrow AL**, **Mai V**, **Gregory KE**,
750          **Kroll JS**, **McMurtry V**, **Ferris MJ**, **Engstrand L**, **Lilja HE**, **Hollister EB**,
751          **Versalovic J**, **Neu J**. 2017. Intestinal dysbiosis in preterm infants preceding
752          necrotizing enterocolitis: a systematic review and meta-analysis. Microbiome **5**:31.

753    18.   **Xiong W**, **Brown CT**, **Morowitz MJ**, **Banfield JF**, **Hettich RL**. 2017. Genome-
754          resolved metaproteomic characterization of preterm infant gut microbiota
755          development reveals species-specific metabolic shifts and variabilities during early
756          life. Microbiome **5**:72.

757    19.   **Stewart CJ**, **Embleton ND**, **Marrs ECL**, **Smith DP**, **Nelson A**, **Abdulkadir B**,
758          **Skeath T**, **Petrosino JF**, **Perry JD**, **Berrington JE**, **Cummings SP**. 2016.
759          Temporal bacterial and metabolic development of the preterm gut reveals specific
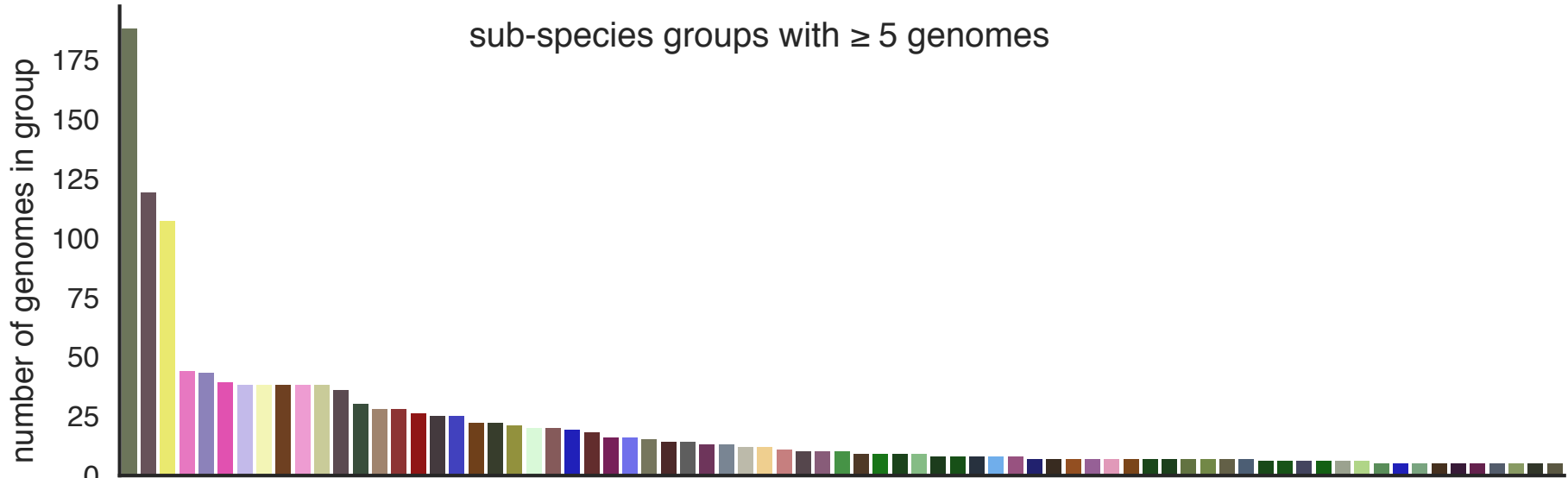760          signatures in health and disease. Microbiome **4**:67.

761  20.  **Brooks B**, **Olm MR**, **Firek BA**, **Baker R**, **Thomas BC**, **Morowitz MJ**, **Banfield**
762      **JF**. 2017. Strain-resolved analysis of hospital rooms and infants reveals overlap
763      between the human and room microbiome. Nat Commun **8**:1814.

764  21.  **Brown CT**, **Olm MR**, **Thomas BC**, **Banfield JF**. 2016. Measurement of bacterial
765      replication rates in microbial communities. Nat Biotechnol **34**:1256–1263.

766  22.  **Manzoni P**. 2016. Clinical Benefits of Lactoferrin for Infants and Children. J
767      Pediatr **173**:S43–S52.

768  23.  **Raghuveer TS**, **McGuire EM**, **Martin SM**, **Wagner BA**, **Rebouché CJ**,
769      **Buettner GR**, **Widness JA**. 2002. Lactoferrin in the Preterm Infants' Diet
770      Attenuates Iron-Induced Oxidation Products. Pediatric Research 2002 52:6
771      **52**:pr2002280–972.

772  24.  **Peng Y**, **Leung HCM**, **Yiu SM**, **Chin FYL**. 2012. IDBA-UD: a de novo
773      assembler for single-cell and metagenomic sequencing data with highly uneven
774      depth. Bioinformatics **28**:1420–1428.

775  25.  **Hyatt D**, **Chen G-L**, **LoCascio PF**, **Land ML**, **Larimer FW**, **Hauser LJ**. 2010.
776      Prodigal: prokaryotic gene recognition and translation initiation site identification.
777      BMC Bioinformatics **11**:119.

778  26.  **Edgar RC**. 2010. Search and clustering orders of magnitude faster than BLAST.
779      Bioinformatics **26**:2460–2461.

780  27.  **Edgar RC**. UBLAST. drive5com.

781  28.  **The UniProt Consortium**. 2015. UniProt: a hub for protein information. Nucleic
782      Acids Res **43**:D204–D212.

783  29.  **Suzek BE**, **Huang H**, **McGarvey P**, **Mazumder R**, **Wu CH**. 2007. UniRef:
784      comprehensive and non-redundant UniProt reference clusters. Bioinformatics
785      **23**:1282–1288.

786  30.  **Kanehisa M**, **Goto S**, **Sato Y**, **Furumichi M**, **Tanabe M**. 2012. KEGG for
787      integration and interpretation of large-scale molecular data sets. Nucleic Acids Res
788      **40**:D109–14.

789  31.  **Minoru Kanehisa SG**. 2000. KEGG: Kyoto Encyclopedia of Genes and
790      Genomes. Nucleic Acids Res **28**:27.

791  32.  **Langmead B**, **Salzberg SL**. 2012. Fast gapped-read alignment with Bowtie 2. Nat
792      Meth **9**:357–359.

793   33.   **Dick GJ**, **Andersson AF**, **Baker BJ**, **Simmons SL**, **Thomas BC**, **Yelton AP**,
794          **Banfield JF**. 2009. Community-wide analysis of microbial genome sequence
795          signatures. Genome Biol **10**:R85–.

796   34.   **Zaharia M**, **Bolosky WJ**, **Curtis K**, **Fox A**, **Patterson D**, **Shenker S**, **Stoica I**,
797          **Karp RM**, **Sittler T**. 2011. Faster and More Accurate Sequence Alignment with
798          SNAP.

799   35.   **Wittek P**, **Gao SC**, **Lim IS**, **Zhao L**. 2013. Somoclu: An Efficient Parallel Library
800          for Self-Organizing Maps.

801   36.   **Ultsch A**. 2005. ESOM-Maps: tools for clustering, visualization, and classification
802          with Emergent SOM.

803   37.   **Altschul SF**, **Gish W**, **Miller W**, **Meyers EW**, **Lipman DJ**. 1990. Basic Local
804          Alignment Search Tool. J Mol Biol **215**:403–410.

805   38.   **Parks DH**, **Imelfort M**, **Skennerton CT**, **Hugenholtz P**, **Tyson GW**. 2015.
806          CheckM: assessing the quality of microbial genomes recovered from isolates,
807          single cells, and metagenomes. Genome Res **25**:gr.186072.114–1055.

808   39.   **Ondov BD**, **Treangen TJ**, **Melsted P**, **Mallonee AB**, **Bergman NH**, **Koren S**,
809          **Phillippy AM**. 2016. Mash: fast genome and metagenome distance estimation
810          using MinHash. Genome Biol **17**:1.

811   40.   **Richter M**, **Rossello-Mora R**. 2009. Shifting the genomic gold standard for the
812          prokaryotic species definition. PNAS **106**:19126–19131.

813   41.   **Olm MR**, **Brown CT**, **Brooks B**, **Banfield JF**. 2017. dRep: a tool for fast and
814          accurate genomic comparisons that enables improved genome recovery from
815          metagenomes through de-replication. ISME J.

816   42.   **Lozupone C**, **Knight R**. 2005. UniFrac: a New Phylogenetic Method for
817          Comparing Microbial Communities. Appl Environ Microbiol **71**:8228–8235.

818   43.   **Xiong W**, **Abraham PE**, **Li Z**, **Pan C**, **Hettich RL**. 2015. Microbial
819          metaproteomics for characterizing the range of metabolic functions and activities
820          of human gut microbiota. Proteomics **15**:3424–3438.

821   44.   **Tabb DL**, **Fernando CG**, **Chambers MC**. 2007. MyriMatch:  Highly Accurate
822          Tandem Mass Spectral Peptide Identification by Multivariate Hypergeometric
823          Analysis. J Proteome Res **6**:654–661.

824   45.   **Ma Z-Q**, **Dasari S**, **Chambers MC**, **Litton MD**, **Sobecki SM**, **Zimmerman LJ**,
825          **Halvey PJ**, **Schilling B**, **Drake PM**, **Gibson BW**, **Tabb DL**. 2009. IDPicker 2.0:
826          Improved Protein Assembly with High Discrimination Peptide Identification
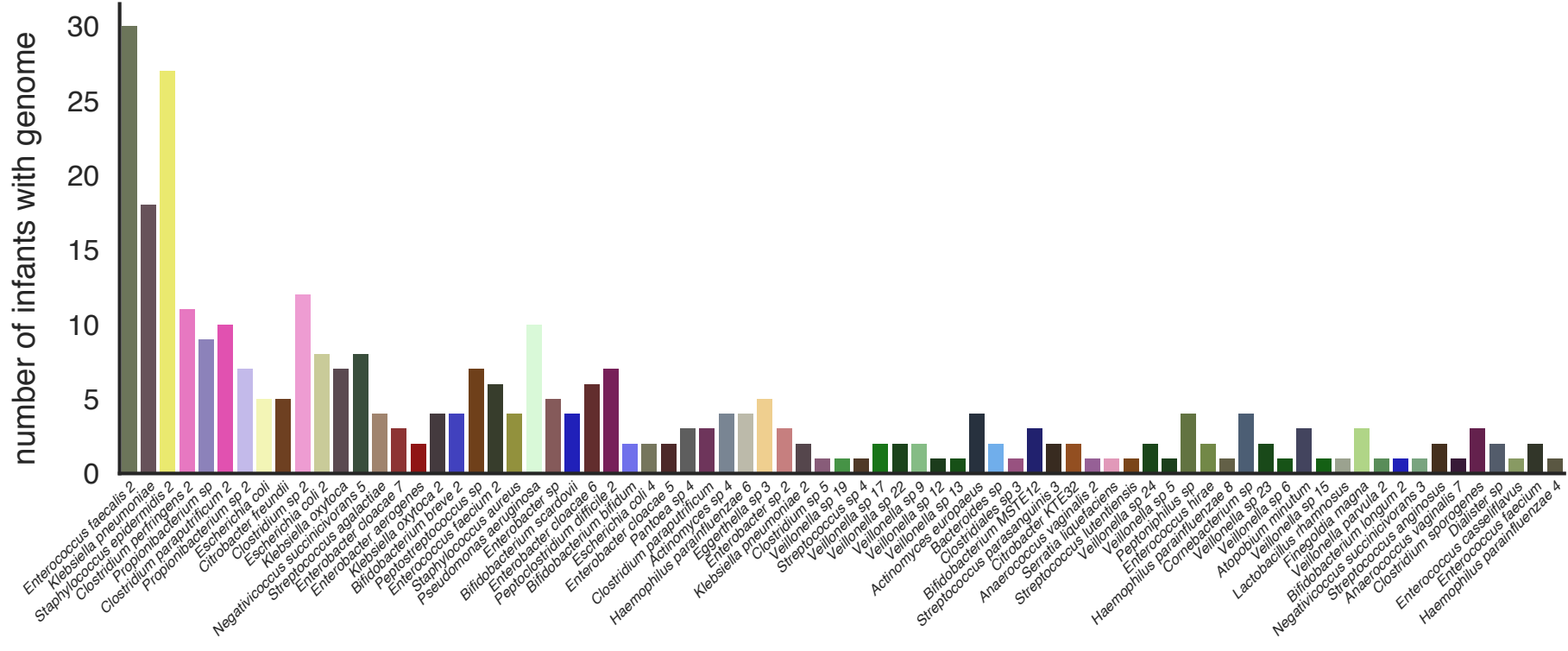827          Filtering. J Proteome Res **8**:3872–3881.

35

828   46.   **Enright AJ**, **Van Dongen S**, **Ouzounis CA**. 2002. An efficient algorithm for
829         large-scale detection of protein families. Nucleic Acids Res **30**:1575–1584.

830   47.   **Robinson MD**, **McCarthy DJ**, **Smyth GK**. 2009. edgeR: a Bioconductor package
831         for differential expression analysis of digital gene expression data. Bioinformatics
832         **26**:139–140.

833   48.   **Tibshirani R**, **Walther G**, **Hastie T**. 2001. Estimating the number of clusters in a
834         data set via the gap statistic. Journal of the Royal Statistical Society: Series B
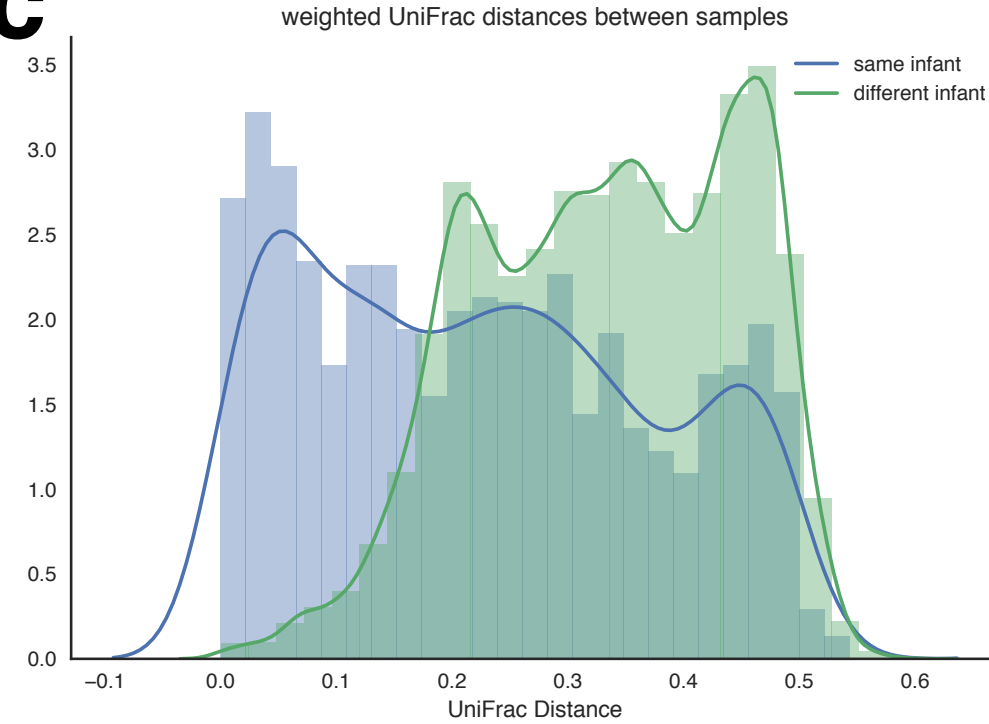835         (Statistical Methodology) **63**:411–423.

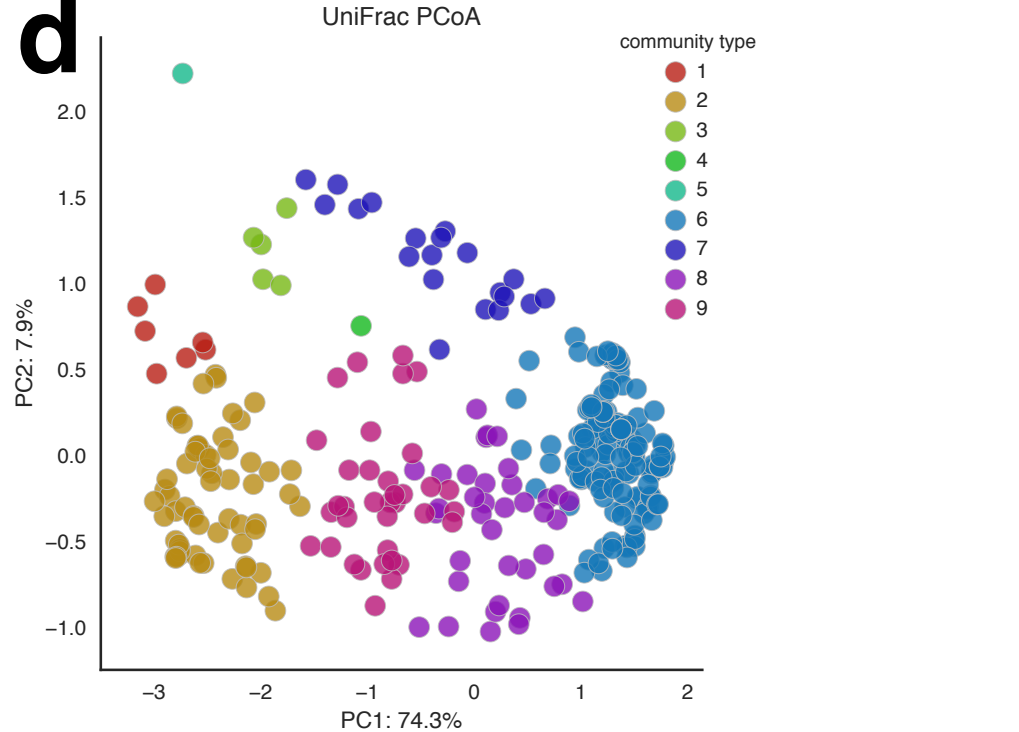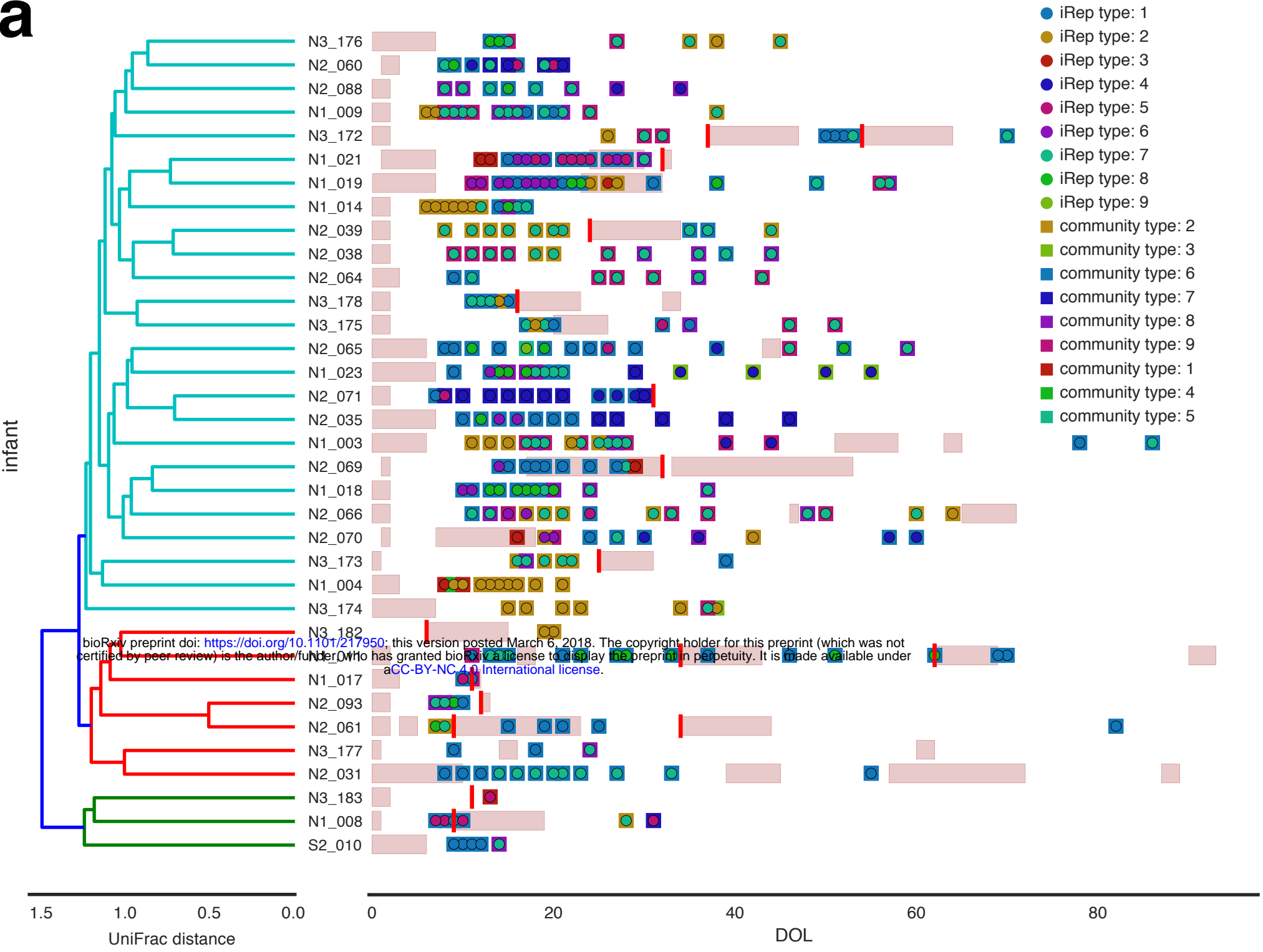**Figure 1**

# Figure 2
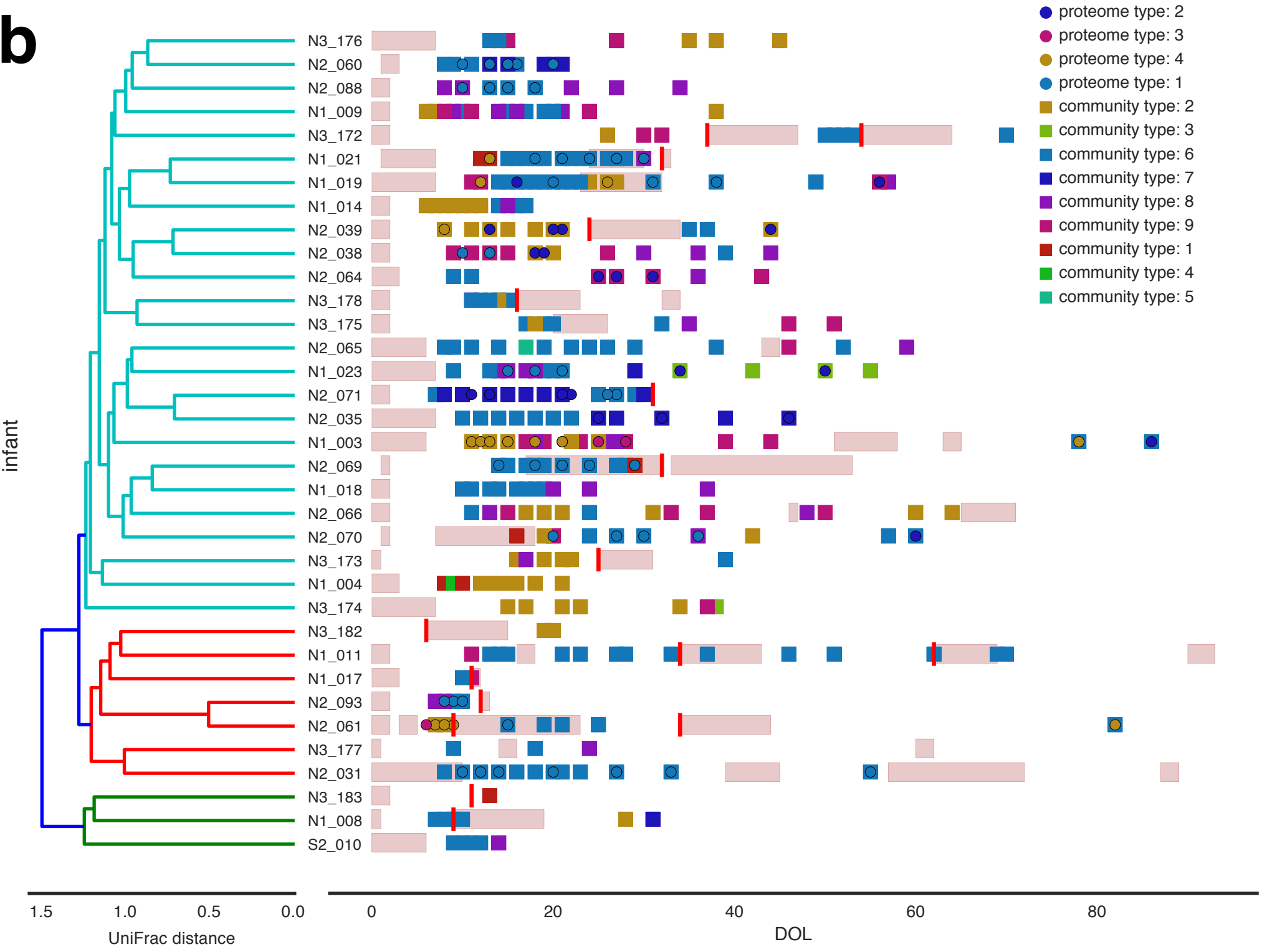
# Figure 3

# Figure 4



A scatter plot titled "Figure 4" showing percent proteome detected versus depth of proteome sampling, with marginal histograms showing count distributions along both axes.

Labeled data points include:
- *Propionibacterium sp*
- *Bifidobacterium breve*
- *Bifidobacterium bifidum*
- *Anaerococcus vaginalis*
- *Enterobacter sp*
- *Enterobacter cloacae*
- *Escherichia coli*
- *Veillonella sp*
- *Klebsiella pneumoniae*
- *Klebsiella oxytoca*
- *Enterococcus faecalis*

Axes:
- y-axis: percent proteome detected (0, 10, 20, 30, 40)
- x-axis: depth of proteome sampling (0.0, 0.1, 0.2, 0.3, 0.4)

Legend for point size:
- min. abundance 0
- med. abundance 23
- max. abundance 100

Color legend (species):
- *Peptoclostridium difficile*
- *Escherichia coli*
- *Lachnospiraceae oral*
- *Bifidobacterium bifidum*
- *Klebsiella oxytoca*
- *Anaerococcus sp*
- *Staphylococcus epidermidis*
- *Streptococcus agalactiae*
- *Bifidobacterium longum*
- *Finegoldia magna*
- *Citrobacter sp*
- *Clostridium sp*
- *Acinetobacter baumannii*
- *Anaerococcus vaginalis*
- *Peptoniphilus harei*
- *Clostridiales sp*
- *Dialister sp*
- *Enterococcus faecium*
- *Peptostreptococcus sp*
- *Bifidobacterium breve*
- *Negativicoccus succinicivorans*
- *Varibaculum cambriense*
- *Propionibacterium sp*
- *Streptococcus lutentiensis*
- *Peptoniphilus sp*
- *Bifidobacterium scardovii*
- *Citrobacter freundii*
- *Pantoea sp*
- *Eggerthella sp*
- *Lactobacillus rhamnosus*
- *Streptococcus anginosus*
- *Enterococcus gallinarum*
- *Enterobacter sp*
- *Clostridium paraputrificum*
- *Veillonella sp*
- *Enterococcus faecalis*
- *Cornebacterium sp*
- *Enterobacter cloacae*
- *Bacteroides sp*
- *Staphylococcus warneri*
- *Staphylococcus aureus*
- *Streptococcus salivarius*
- *Lactobacillus fermentum*
- *Dermabacter HFH0086*
- *Streptococcus parasanguinis*
- *Actinomyces urogenitalis*
- *Actinomyces sp*
- *Klebsiella pneumoniae*
- *Haemophilus parainfluenzae*

# Figure 5

# Figure 6



Klebsiella pneumoniae (min. 10% proteome)

# Figure 7



**a**

MW p-value = 2.6 x 10⁻⁹*

shannon

- no antibiotics
- during or within 5 days antibiotics

**b**

MW p-value = 4.0 x 10⁻⁴*

shannon diversity

control    NEC

**c**

- no antibiotics
- during or within 5 days antibiotics

control    NEC

**d**

- control
- NEC

early    late    early    late