

Running title: Human infant gut microbiota ecology

**Major fecal microbiota shifts in composition and diversity with age
in a geographically restricted cohort of mothers and their children**

Ekaterina Avershina^{1*}, Ola. Storrø², Torbjørn Øien², Roar Johnsen², Phil Pope¹ and Knut Rudi^{1*}

¹Department of Chemistry, Biotechnology and Food Science, University of Life Sciences, Ås, Norway, ²Department of Public Health and General Practice, Norwegian University of Science and Technology, Trondheim, Norway

* Corresponding authors: e-mail: ekaterina.avershina@umb.no and knut.rudi@umb.no, phone: +47 64 96 59 00, Fax: +47 64 96 59 01

Keywords: 16S rRNA gene, infant gut microbiota

1 ABSTRACT

2 Despite the importance, the diversity of the human infant gut microbiota still remains
3 poorly characterized at the regional scale. Here we investigated the fecal microbiota
4 diversity in a large 16S rRNA gene dataset from a healthy cohort of 86 mothers and their
5 children from the Trondheim region in Norway. Samples were collected from mothers
6 during early and late pregnancy, as well as their children at 3 days, 10 days, 4 months, 1
7 year and 2 years of age. Using a combination of Sanger sequencing of amplicon mixtures
8 (without cloning), real-time quantitative PCR and deep pyrosequencing we observed a
9 clear age related colonization pattern in children that was surprisingly evident between 3
10 and 10 days samples. In contrast, we did not observe any shifts in microbial composition
11 during pregnancy. We found that alpha-diversity was highest at 2 years and lowest at 4
12 months, whereas beta-diversity estimates indicated highest inter-individual variation in
13 newborns. Variation significantly decreased by the age of 10 days and was observed to be
14 convergent over time; however, there were still major differences between 2 years and
15 adults whom exhibited the lowest inter-individual diversity. Taken together, the major
16 age-affiliated population shift within gut microbiota suggests that there are important
17 mechanisms for transmission and persistence of gut bacteria that remain unknown.

18

19 INTRODUCTION

20 Whilst it is widely accepted that the human gut is one of the most densely populated bacterial
21 communities on Earth (Whitman, *et al.*, 1998), the general mechanisms for host-bacterial
22 interactions are not yet completely described (Avershina & Rudi, 2013). Previously, the
23 scientific community unanimously assumed that humans are born sterile (Ley, *et al.*, 2006,
24 Marques, *et al.*, 2010), although evidence now exists for pre-natal colonization (Jimenez, *et al.*,
25 2008, Satokari, *et al.*, 2009). Regardless of the required time for initial colonization, it is
26 absolute that development of this unique and intricate community takes several years to reach
27 its maturity (Marchesi, 2011). There are many factors which supposedly play a role in
28 development of gut microbiota; initial inoculation occurs via the mother's birth canal when a
29 child is born vaginally, subsequently an infant will frequently receive bacteria via breast milk
30 (Martin, *et al.*, 2007) and the surrounding environment also exerts a constant influence. Existing
31 reports have addressed various environmental influences towards gut microbiota such as age
32 (Palmer, *et al.*, 2007, Claesson, *et al.*, 2011), geography and diet (De Filippo, *et al.*, 2010,
33 Yatsunenکو, *et al.*, 2012). There are also recent suggestions of immunological modulations of
34 the microbiota during pregnancy (Koren, *et al.*, 2012). However, much less is known about
35 transmission and persistence of gut bacteria in a population during the host's first years of life.
36 We have previously described transmission of some particular gut bacteria from mother to child
37 (Bjerke, *et al.*, 2011, de Muinck, *et al.*, 2011, Avershina, *et al.*, 2013), while we have not yet
38 addressed general patterns of bacterial persistence and diversity in a healthy randomly selected
39 population of children and their mothers.

40 The aim of this study was therefore to address longitudinal fecal microbiota shifts in
41 composition and diversity in children and their mothers in a geographically restricted cohort.
42 We analyzed stool samples from 86 mother/child pairs, collected two times during the mothers
43 pregnancy (15.0 ± 4.2 and 37.5 ± 1.8 gestation weeks) and five times from infants (ages 3 and 10

44 days, 4 months, 1 year and 2 years). We used a polyphasic analytical approach consisting of
45 direct mixed 16S rRNA gene Sanger sequencing (analysis of electropherograms containing
46 information on all amplicon variants) (Zimonja, *et al.*, 2008), real-time quantitative PCR
47 (Ginzinger, 2002) and 454-sequencing (Ronaghi, 2001). We present results suggesting highly
48 age-dependent bacterial persistence and diversity patterns within the population. Furthermore,
49 we also present support for mother to child transmission of adult associated gut bacteria –
50 surprisingly not during the birth process but at a later stage.

51 **MATERIALS AND METHODS**

52 **Study material and sample preparation**

53 Fecal samples were collected from the IMPACT cohort study among small children and
54 mothers in Trondheim, which is a nested cohort within the PACT study (Prevention of Allergy
55 among Children in Trondheim) (Storro, *et al.*, 2010). Most of the children were delivered
56 vaginally (90 %), and at term (90 %). There was a high frequency of breast feeding, 97 % of
57 infants were breast-fed during the first six weeks of life. By the age of 4 months, 66.7 % of
58 infants were exclusively breast-fed, 23.8 % were receiving either formula or solid food (fruits,
59 vegetables, wheat, bread, corn, rice) complementary to breast milk, and 9.5 % of infants were
60 receiving only formula and/or solid food. More details about the cohort characteristics are given
61 by Storro *et al.* (Storro, *et al.*, 2011).

62 Fecal specimens were stored in sterile Cary Blair transport and holding medium (BD
63 Diagnostics Sparks, MD 21152 USA). Each specimen was frozen at -20°C within 2 hours after
64 defecation and transported to the laboratory for further storage at -80°C within 1 day (for
65 children) or 4 weeks (for pregnant women). Details about the IMPACT fecal material is given
66 by (Oien, *et al.*, 2006). The dataset analyzed contained samples from both early (first to second

67 trimester) and late pregnancy (third trimester) from the mothers, and 3 days, 10 days, 4 months,
68 1 year and 2 years from the children.

69 We purified fecal DNA with paramagnetic beads in accordance with an optimized and
70 automated protocol (Skanseng, *et al.*, 2006). Briefly, this protocol involved mechanical lysis
71 with glass beads, and DNA purification with silica particles. Mechanical lysis was chosen since
72 the compositions of the gut bacteria cell walls are largely unknown.

73 **Direct mixed sequence analysis**

74 The V3 – V4 region of 16S rRNA gene was PCR amplified using the primers targeting
75 universally conserved gene regions (Nadkarni, *et al.*, 2002). Subsequently the V4 region (198
76 bp) was targeted for sequencing using a mixed Sanger approach. The resulting sequence spectra
77 contained information for the 16S rRNA genes representative of all the bacteria in a given
78 sample.

79 The alpha- and beta- diversity of each spectrum was assessed by means of modified Simpson's
80 diversity index c_{mixed} (Eq. 1) and modified Bray-Curtis dissimilarity index (Eq. 2) respectively.
81 Calculations were based on the fluorescence intensity fractions of each nucleotide position. The
82 rationale is that these intensity fractions will reflect diversity. In case there is only one bacteria
83 in a sample, there will be only one nucleotide in every position of the sequence spectrum, and
84 therefore nucleotide fractions in every position will equal 1:0:0:0. In the case of a mixture of a
85 range of different bacteria, though, the fractions will converge towards 0.25:0.25:0.25:0.25.
86 Based on these fractions, one could estimate diversity in a sample which is independent of
87 operational taxonomic units (OTUs).

$$88 \quad 1/c_{mixed} = \frac{\sum_{i=1}^n (G_i)^2 + \sum_{i=1}^n (A_i)^2 + \sum_{i=1}^n (T_i)^2 + \sum_{i=1}^n (C_i)^2}{n} \text{ (Eq. 1);}$$

$$BC_{ij} = \frac{\sum_{k=1}^n |G_{ki} - G_{kj}| + \sum_{k=1}^n |A_{ki} - A_{kj}| + \sum_{k=1}^n |T_{ki} - T_{kj}| + \sum_{k=1}^n |C_{ki} - C_{kj}|}{\sum_{k=1}^n (G_{ki} + G_{kj}) + \sum_{k=1}^n (A_{ki} + A_{kj}) + \sum_{k=1}^n (T_{ki} + T_{kj}) + \sum_{k=1}^n (C_{ki} + C_{kj})} \quad (\text{Eq. 2});$$

89 Detailed description of the diversity indices calculations is given in Avershina et al. (Avershina,
 90 *et al.*, 2013). Beta-diversity was assessed both between samples belonging to the same age
 91 group, as well as between samples belonging to the same mother-child pair but at different time
 92 points. Significant difference between indices at various time points was tested using
 93 Friedman's test, – a non-parametric version of two-way ANOVA which takes into account
 94 possible correlation between the measurements (MATLAB® documentation, 2010). For those
 95 samples, where we did not expect the correlation, Kruskal-Wallis test was used. The null
 96 hypothesis was rejected at the level of 5 %.

98 Information on the most dominant bacteria was subsequently resolved using Multivariate Curve
 99 Resolution analysis (MCR-ALS). This analysis allows recovery of the common information
 100 contained between the samples of interest into so-called components, as well as simultaneous
 101 relative quantification of this information in all the samples (Zimonja, *et al.*, 2008). Taxonomic
 102 level of components' resolution for non-defined bacterial assemblages directly depends on the
 103 diversity represented within a dataset (Rudi, *et al.*, 2012, Sekelja, *et al.*, 2012). If a given
 104 phylum is represented by one clearly dominant genus, then the signature sequence for this genus
 105 will be resolved as a component. Whilst if there were several equally distributed genera within
 106 the same family, then the signature sequence for this family would have been recovered. Prior
 107 to MCR-ALS analysis, one needs to specify the number of components to be resolved. In case
 108 the set number is too high, the 'real' component would be split and thus at least two of the
 109 resolved components would contain the same information. This can be detected by biological
 110 reasoning since these components will then represent the same taxonomic group. To define the
 111 initial number of components (initial estimates i), we used both Principal Component Analysis

112 (PCA) and Evolving Factor Analysis (EFA) as recommended (Tauler, *et al.*, 1995). The detailed
113 description of use of MCR-ALS analysis for mixed sequence resolution can be found in
114 Avershina *et al.* (Avershina, *et al.*, 2013). Resolved components spectra were manually base-
115 called and classified by Ribosomal Database Project (RDP) hierarchical classifier (Wang, *et*
116 *al.*, 2007).

117 To address the longitudinal structure of the MCR-ALS score data, i.e. relative abundance of
118 resolved components, Parallel Factor Analysis (PARAFAC) method was used. PARAFAC is a
119 multi-way generalization of the two-way PCA. However, unlike PCA the rotation problem is
120 omitted so that pure components can be resolved (Bro, 1997). The core consistency index was
121 used as a criterion for determining the number of components.

122 **Real-time quantitative PCR**

123 We have previously qPCR-amplified the 16S rRNA gene of commonly identified gut bacteria,
124 as well as some pathogenic bacterial species (Storro, *et al.*, 2011) for the same study cohort.
125 Among tested species were *Bacteroides fragilis*, *Bifidobacterium longum*, *Bifidobacterium*
126 *breve*, *Bifidobacterium animalis* subsp. *lactis*, genus *Bifidobacterium*, *Clostridium difficile*,
127 *Clostridium perfringens*, *Lactobacillus rhamnosus*, *Lactobacillus reuteri* and *Helicobacter*
128 *pylori*. For this work, we binarized these data based on whether the given bacterium was or
129 wasn't detected in a sample. For every age unweighted Cohen's kappa indices (Sim & Wright,
130 2005) were calculated to evaluate whether there was an agreement between detection of a given
131 bacteria in mothers and children. Interpretation of the index was performed using guidelines
132 provided in the MATLAB[®] script for Cohen's kappa index calculation (Cardillo, 2007). The
133 relative amount of the detected vs non-detected populations of bacteria is represented in
134 Supplementary Figure 1. "Non-detected" populations were defined as populations that did not
135 show amplification after 40 cycles. Some bacteria (*L. rhamnosus* and *C. difficile*) were not

136 detected in any of the mothers, whereas others (e.g. *H. pylori*) were detected only in two
137 mothers (Supplementary Table 1). Therefore, to ensure sufficient amount of information, only
138 bacterial groups that were detected in more than 11 mothers were included in the analysis. The
139 bacterial groups that satisfied this criterion were: *B. longum*, genus *Bifidobacterium*, *B. fragilis*
140 and *E. coli*. We also addressed the persistence patterns of these four bacteria in a population by
141 calculating the fraction of individuals, in which the species was detected at a time point ‘*x*’
142 given it was detected at a time point ‘*x-1*’.

143 **Pyrosequencing analysis**

144 A subset of seven random mother and child pairs were selected for deep 454-sequencing from
145 the pairs with the most complete temporal series in the main study cohort. DNA isolation,
146 amplicon and PCR conditions were the same as for direct sequencing approach. The only
147 difference was the modification of PCR primers targeting V3 – V4 region of 16S rRNA, to be
148 adapted to the GS-FLX Titanium Chemistry (454 Life Sciences, USA). Sequencing was
149 performed according to the manufacturer’s recommendations at the Norwegian High-
150 Throughput Sequencing Centre (Oslo, Norway). Pyrosequencing data were analyzed using
151 QIIME pipeline (Caporaso, *et al.*, 2010). Error-correction, chimera removal and operational
152 taxonomic unit (OTUs) clustering was performed using USEARCH quality filtering with
153 QIIME, which incorporates UCHIME (Edgar, *et al.*, 2011) and a 97 % sequence identity
154 threshold. The RDP classifier (Wang, *et al.*, 2007) was used to assign taxonomic identity to the
155 resulting OTUs. For a phylogeny-based diversity assessment, we used weighted UniFrac
156 hierarchical clustering (Lozupone & Knight, 2005) based on 10 rarefactions with 1600
157 randomly selected sequences per sample for each rarefaction.

158 In order to investigate what shapes gut microbiota both in infancy and adulthood, we fitted
159 observed species distributions to common used distributions using the Species Diversity and

160 Richness v. 4.1.2 (PISCES Conservation Ltd., UK) software. Hubbell's model of neutrality,
161 often used as a null model of community structure (Magurran, 2004), assumes that when an
162 individual dies in a saturated community, the probability of its replacement by an offspring of
163 rare species is the same as by an offspring of a more abundant species. Jabot and Chave (2011)
164 have developed a generalization of this model introducing a parameter δ . This parameter
165 estimates the non-neutrality of the system based on the deviation of observed species evenness
166 as opposed to the system being best described by neutral model. When δ is positive, dominant
167 species have higher chance of taking the place of the dead individual, whereas negative values
168 indicate that rare species' chances increase. Based on 1000 randomly selected sequences per
169 sample from the chimera- and noise-free pyrosequencing dataset, we calculated non-neutrality
170 parameter δ using Parthy v. 1.0 software (Jabot & Chave, 2011).

171 **RESULTS**

172 **Mixed sequence analysis**

173 Nucleotide alpha-diversity (Simpson's diversity index) of mixed spectra ranged from 1.77 ± 0.10
174 [mean \pm standard deviation] at 4 month old to 1.91 ± 0.09 at 2 year old infants (Figure 1A).
175 Generally, diversity of adult' stool samples was higher than that of newborns ($p = 0.0001$) and
176 4 month old infants ($p = 2.26 * 10^{-9}$). At 1 year of age, the diversity increased compared to 4-
177 month-olds ($p = 0.0028$) and then further increased by 2 years of age ($p = 0.0054$).

178 Newborns exhibited highest beta-diversity between individuals (modified Bray-Curtis index
179 $BC = 0.20 \pm 0.02$ and 0.18 ± 0.03 for 3- and 10-days-old infants respectively; Figure 1B). By the
180 age of 4 months, the variation within the population had significantly decreased ($p = 7.51 * 10^{-13}$)
181 and remained the same up to 1 year. Though the beta-diversity between stool samples from
182 2-year-olds was significantly lower than that of 1-year-olds ($p = 1.54 * 10^{-5}$), it was still
183 significantly higher than the beta-diversity between adult stool samples ($p = 4.38 * 10^{-6}$). In

184 addition to inter-individual comparisons, beta-diversity estimations were used to analyze intra-
185 individual variation that developed within an individual from one time point to another (Figure
186 1C). The highest variation (highest beta-diversity) was observed between the spectra of mothers
187 at their late pregnancy stage and 3 days old infants ($BC = 0.21 \pm 0.04$), as well as between 4
188 months old and 1 year old children ($BC = 0.20 \pm 0.04$), whereas the least variation (lowest beta-
189 diversity) was observed between stool samples collected from mothers at two pregnancy
190 trimesters ($BC = 0.08 \pm 0.03$) and also between 1- and 2-year-olds ($BC = 0.12 \pm 0.02$).

191 Both PCA and EFA suggested six components to be resolved by MCR-ALS. When six
192 components were used, the information on *Bacteroidetes* group was entirely absent. Therefore
193 MCR-ALS analysis was repeated by gradually increasing the number of components to be
194 resolved until the duplication event. In total, eight components accounting for 70 % of the
195 variation in the system was resolved by MCR-ALS and classified by RDP classifier
196 (Supplementary Table 2).

197 Taxonomically, stool samples analyzed from mothers were rich in *Lachnospiraceae*- and
198 *Faecalibacterium*-affiliated components (Figure 2). At 3 days, all eight components seemed to
199 be evenly represented, but by the age of 10 days there was a significant decrease in the level of
200 *Lactobacillales* ($p = 0.0191$). By the age of four months, bifidobacteria constituted 57.6 % of
201 total gut microbiota, whereas *Lactobacillales*- and *Streptococcus*-affiliated components were
202 diminished ($p = 0.0135$ and $p = 0.0001$ respectively). At 1 and 2 years of age, average
203 composition resembled that of pregnant women, though there were several pronounced
204 differences. For example, the *Bifidobacterium*-affiliated ($p = 0.0042$ and $p = 0.0021$ for 1 and
205 2 years respectively), and other Actinobacteria- ($p = 0.0016$ and $p = 2.3 \cdot 10^{-5}$ for 1 and 2 years
206 respectively) components were higher in children than in their mothers, whereas
207 *Faecalibacterium*- ($p = 4.3 \cdot 10^{-6}$ and $p = 5.9 \cdot 10^{-7}$ for 1 and 2 years respectively) and

208 *Bacteroides*-affiliated ($p = 1.4 \times 10^{-5}$ and $p = 5.6 \times 10^{-8}$ for 1 and 2 years respectively) components
209 were lower.

210 Due to the fact that the majority of infants were born vaginally, at term and were breast-fed
211 during the first days of life, we could not investigate the effect of birth mode and diet. However,
212 we could test whether implementation of solid food (wheat, rice, corn) at four months would
213 affect fecal microbial composition. These analyses showed no significant difference in relative
214 composition of gut microbiota.

215 In order to investigate longitudinal structure in the data (i.e. individual sharing of bacteria for
216 more than one time point), 3 components PARAFAC model was deduced based on a core
217 consistency index of more than 99 %. The loadings for the MCR-ALS components dimension
218 indicate that *Escherichia*-, *Bifidobacterium*- and *Lachnospiraceae*-affiliated components
219 influenced the longitudinal structure of the data (Figure 3A). In particular, the *Escherichia*-
220 affiliated component was associated with 3 and 10 days, *Bifidobacterium*- with 3 days, 10 days
221 and 4 months, while *Lachnospiraceae*-affiliated component was associated with early and late
222 pregnancy, in addition to 1 and 2 years (Figure 3B).

223 **Real-time quantitative PCR analysis of prevalence**

224 Figure 4 illustrates qPCR prevalence data calculated for selected bacterial groups both for the
225 whole study cohort, as well as for a subpopulation of children whose mothers tested positive
226 for the target bacterium (mother-child positive subpopulation). At 10 days, *E. coli* was more
227 frequently detected in those children whose mothers also tested positive for this bacterium ($p =$
228 0.002). Interestingly, the difference between detection frequencies of this bacterium in mother-
229 child positive subpopulation and total children population was higher in 10 days as compared
230 to 3 days. This may indicate either postnatal or very low at-birth transmission of this bacterial
231 species. *B. longum* was deemed to be one of the most persistent colonizers among the four

232 bacterial groups tested. Already by the age of 10 days, it was detected in nearly all infants who
233 tested positive at 3 days after birth (Figure 4). Even by the age of 2 years, this species persisted
234 in the majority of infants who previously tested positive. In contrast, *E. coli* detection was
235 observed to be stable during the first year (80 % – 85 % of population). However, by 2 years a
236 detection limit had decreased to 45 % of children who previously tested positive.

237 Cohen's kappa index was used to indicate the magnitude of agreement between the detection
238 of a given bacteria in an individual mother and her child (in the whole cohort). In our dataset
239 the index ranged from -0.05 (poor agreement) to 0.30 (fair agreement) and was observed to
240 decrease with age, indicating that the detection of a given bacterium in 1-2 year old children
241 was less dependent on their mother testing positive (Table 1). In concurrence with qPCR
242 prevalence data (Figure 4), Cohen's kappa indices indicated slight to fair agreement both for *E.*
243 *coli* and *B. fragilis*. The ranking is based on the guidelines to the MATLAB[®] script for the index
244 calculation (Cardillo, 2007). Bifidobacteria were observed to be negative at 4 months,
245 indicating poor agreement in mother-child detection patterns. High p-values ($p > 0.05$) also
246 support low correspondence between detection of a given bacteria in mothers and children.

247 **Pyrosequencing data analysis**

248 Eight samples, mostly belonging to one mother-child pair, were removed from the analysis due
249 to a low number of recovered sequences (less than 2000 sequences per sample). Therefore the
250 analysis was performed on a total of 39 samples from 6 children and 5 mothers. After quality
251 filtering, chimera-removal and normalization, 370207 sequences were used for subsequent
252 analysis with a mean of 9492 sequences per sample (ranging from 2146 to 21317 sequences per
253 sample). Apart from one sample, stool samples from mothers' and 1- and 2-years-old infants
254 clustered separately from stool samples of newborns and 4-month-olds based on weighted
255 UniFrac distances (1600 sequences per sample, bootstrap values are based on 10 rarefactions;

256 Supplementary Figure 3A). To examine how similar the fecal microbiota from different age
257 groups was, we used Jaccard distance index calculated for detected OTUs (Supplementary
258 Figure 3B). Overall, there was higher variation in microbiota from children when compared to
259 mothers ($p = 0.0011$ and $p = 0.0001$ at 3 days and 2 years of age respectively), although the
260 microbiota of newly-born children were more similar to each other than to their related ($p =$
261 0.0010 , $p = 0.0011$ and $p = 0.0034$ for 3 days, 10 days and 4 months respectively) and unrelated
262 mothers ($p = 0.0011$, $p = 0.0006$ and $p = 0.0024$ for 3 days, 10 days and 4 months respectively).
263 By the age of 1 year, their microbiota was as similar to adults as it was to other children from
264 the same age group.

265 We compared how many OTUs were shared between five children at various time points and
266 their mothers (both related and unrelated). In total, 30 samples were used for these comparisons.
267 From birth to 4 months of age, only one child had more OTUs shared with his own mother than
268 with any other unrelated mother. However, by the age of 2 years the number of children who
269 shared more OTUs with their mothers than with other unrelated mothers increased to 3 out of
270 5 (Supplementary Table 3). We also examined which OTUs were underrepresented in children
271 at various ages compared to their mothers (Supplementary Tables 4 – 8). In the immediate
272 period after birth (days 1-3), 1230 OTUs were absent in all infant samples, of which 44 % were
273 affiliated to the family of *Lachnospiraceae*. At ages 1-2 years, 500 OTUs were absent,
274 composed of approximately 30 % that were affiliated to the *Lachnospiraceae*. Overall
275 *Lachnospiraceae*-affiliated OTUs which had representatives in all children at a given age were
276 first detected at 1 year, although in one child OTUs affiliated to this clostridial family were
277 detected right after birth. In contrast, within the first days after birth only OTUs affiliated to the
278 *Bifidobacteriaceae*, *Streptococcaceae* and *Staphylococcaceae* were shared among all infants
279 and by four months only *Bifidobacteriaceae*-affiliated OTUs were shared. By the age of 1 year

280 the majority of OTUs were affiliated to the *Clostridiales*, whereas at 2 years shared
281 *Bacteroidales*-affiliated OTUs also appeared.

282 Depending on ecological forces that structure communities, species within these communities
283 may follow different distributions that can be described mathematically (Magurran, 2004). We
284 therefore fitted OTU distributions to these common distribution curves (Supplementary Table
285 9). The majority of samples fitted well to truncated log normal distribution, two samples,
286 belonging to one child at 3 and 10 days of age, fitted log series distribution. The geometric and
287 broken stick distributions didn't fit the data. We also tested whether distributions fitted a neutral
288 model and how much they deviate from it. All the samples showed higher dominance than it
289 would be expected in case of neutrality (Supplementary Figure 2), though there was a
290 significant difference in deviation between mothers and 3-days-olds ($p = 0.0091$). Moreover,
291 when combined, in infancy as well as at 4 months, the dominance was significantly higher than
292 in adults and 1- and 2-year-olds ($p = 0.0001$).

293 **Data consistency**

294 To address whether MCR-ALS and pyrosequencing predictions of fecal microbiota correspond
295 to each other, we selected all OTUs belonging to taxonomical groups predicted by MCR-ALS
296 from a pyrosequencing dataset. We then grouped those OTUs in correspondence with MCR-
297 ALS components and calculated their relative amounts based on the total number of OTUs.
298 Pearson's correlation analysis revealed high correlation between MCR-ALS predictions and
299 pyrosequencing results (correlation coefficient = 0.7463, $p = 4.47 \cdot 10^{-51}$).

300 **DISCUSSION**

301 Interestingly there was a significant drop in inter-individual beta-diversity in a short period of
302 time after birth (3 to 10 days), as assessed by mixed sequencing. Due to practical reasons, many
303 temporal research studies of fecal microbiota face a trade-off between sampling frequency and

304 number of individuals included in the study. To our knowledge, all temporal fecal microbiota
305 studies to date that have extensive sampling during first weeks of life (Favier, *et al.*, 2003,
306 Palmer, *et al.*, 2007, Koenig, *et al.*, 2011) have few individuals analyzed; whereas studies with
307 high sample numbers often have fewer or more infrequent time-points (Yatsunenko, *et al.*,
308 2012). However, our results illustrate that significant differences in average bacterial
309 composition and beta-diversity occurs between 3 and 10 days. These data therefore suggest that
310 to better understand the development of gut microbiota, gaps between sampling periods should
311 be reduced, particularly for those studies that compare different populations (Yatsunenko, *et*
312 *al.*, 2012).

313 Pyrosequencing and mixed sequence analysis both demonstrated individualized clustering of
314 the fecal microbiota during early and late pregnancy in our cohort, with little or no evidence for
315 population-based changes during pregnancy. We were therefore not able to reproduce the
316 results of a major change in the fecal microbiota between early and late pregnancy, as recently
317 reported by Koren *et al.* (Koren, *et al.*, 2012). Since our sampling times matches that of Koren
318 *et al.* with ± 3 weeks we believe that sampling time cannot explain the differences in microbiota
319 detected between the two studies. The most likely explanation would therefore be that there are
320 true differences in the gut microbiota composition among pregnant women in the two cohorts.

321 QPCR analysis suggested a relatively low direct transmission of gut bacteria from mother to
322 child; at 10 days of age there was better overall agreement between detection of bacteria in
323 mother-child pairs than at 3 days (Table 1). Even early colonizers such as *E. coli* were not likely
324 to be directly transmitted at birth, but rather during first days of life (Figure 4). The difference
325 in detection of this species in mother-child positive subpopulation and the total population was
326 higher at 10 days than at 3 days. Based on differences between weighted UniFrac (takes into
327 account relative amounts) and Jaccard (takes into account only presence/absence data)
328 distances, it may be suggested that by 1-2 years of age adult-characteristic OTUs already

329 appeared in the gut, though they were still rare. Interestingly, many OTUs affiliated to
330 *Lachnospiraceae* were shared between mothers and 1-2 year old children, suggesting that these
331 species possibly originate from the mother. PARAFAC data based on mixed sequencing also
332 supported sharing of this component between mothers and infants. Even though detection of
333 bifidobacteria seemed to be independent of the mother, frequency of *B. longum* was higher in
334 a mother-child positive sub-population, which is in line with a recent model suggesting
335 transmittance of *B. longum* subsp. *longum* from mother to child (Makino, *et al.*, 2011).

336 At 3 days of age, there was relatively high abundance of *Lactobacillales* in stool samples
337 (Figure 2). Lactobacilli are often isolated from human breast milk (Martin, *et al.*, 2003, Martin,
338 *et al.*, 2007), and it was noted that the majority of infants (98 %) in our cohort were exclusively
339 breast-fed during the first six weeks of life. Interestingly, by the age of 10 days the level of this
340 bacterial group was observed to decline despite no changes in diet with respect to breast milk
341 intake. As such, we hypothesize that lactobacilli detected in this study were possibly acquired
342 via the vaginal microbiota of the mother during the infant's passage through the birth channel.
343 If we assume that neutral processes (i.e. random replacement of a dead individual in a
344 community by an offspring of other species regardless of relative abundance of this species)
345 are not involved in shaping gut microbiota, one would expect low individual alpha-diversity
346 coinciding with high inter-individual beta-diversity. In contrast, we observed steady decreases
347 in beta-diversity over time (lowest among adult women) suggesting that overall microbiota
348 development is ultimately directed towards a more stable community. Furthermore, delta
349 values, which characterize a deviation from neutrality, were significantly lower in adulthood
350 than in infancy.

351 In contrast to our findings, it has recently been argued that niche selection is also the main force
352 shaping the distal gut community. This conclusion was based on the fact that microbial OTUs
353 in the gut were more closely related to each other than what would be expected in case of

354 neutrally shaped community (Jeraldo, *et al.*, 2012). The discrepancy, however, could be
355 explained by the fact that niche selection will always limit the phylotypes allowed in a given
356 environment (Magurran, 2004), and that the distal gut represents a highly selective environment
357 (Marchesi, 2011), whereas among the allowed phylotypes neutral processes could be important.
358 Probably, since we did not take phylogenetic distances into account we also discovered the
359 neutral processes as a potential contributor. This explanation is coherent with our recently
360 proposed interface model for bacterial-host interactions, suggesting host selection independent
361 of the actual services provided (Avershina & Rudi, 2013).

362 In conclusion, our analyses of a large longitudinal cohort of mothers and their children have
363 revealed new knowledge about the ecology of human gut bacteria, suggesting that there are still
364 important mechanisms that remain unknown.

365 **ACKNOWLEDGEMENTS**

366 Funding for the IMPACT study was obtained from GlaxoSmithKline AS, Norway. The PACT
367 study was funded by the Norwegian Department of Health and Social affairs from 1997–2003.
368 A university scholarship from NTNU funded the research fellows. The mixed sequencing
369 analyses were funded by a research levy on certain agricultural products from the Norwegian
370 Government. PBP is funded by Norwegian Research Council project 214042. Authors have no
371 conflict of interest to declare.

372 **REFERENCES**

- 373 Avershina E & Rudi K (2013) Is it who you are or what you do that is important in the human
374 gut? *Beneficial Microbes* in press.
375 Avershina E & Rudi K (2013) Is it who you are or what you do that is important in the human
376 gut? *Benef Microbes* **4**: 219-222.
377 Avershina E, Storro O, Oien T, Johnsen R, Wilson R, Egeland T & Rudi K (2013) Succession
378 and correlation-networks of bifidobacteria in a large unselected cohort of mothers and their
379 children. *Appl Environ Microbiol* **79**: 497-507.
380 Bjerke GA, Wilson R, Storro O, Oyen T, Johnsen R & Rudi K (2011) Mother-to-child
381 transmission of and multiple-strain colonization by *Bacteroides fragilis* in a cohort of mothers
382 and their children. *Appl Environ Microbiol* **77**: 8318-8324.

- 383 Bro R (1997) PARAFAC. Tutorial and applications. *Chemometrics and Intelligent Laboratory*
 384 *Systems* **38**: 149-171.
- 385 Caporaso JG, Kuczynski J, Stombaugh J, *et al.* (2010) QIIME allows analysis of high-
 386 throughput community sequencing data. *Nat Methods* **7**: 335-336.
- 387 Cardillo G (2007) Cohen's kappa: compute the Cohen's kappa ratio on a 2x2 matrix. Vol. 2012
 388 ed.^eds.), p.^pp. MathWorks, MATLAB Central File Exchange.
- 389 Claesson MJ, Cusack S, O'Sullivan O, *et al.* (2011) Composition, variability, and temporal
 390 stability of the intestinal microbiota of the elderly. *Proc Natl Acad Sci U S A* **108 Suppl 1**:
 391 4586-4591.
- 392 De Filippo C, Cavalieri D, Di Paola M, *et al.* (2010) Impact of diet in shaping gut microbiota
 393 revealed by a comparative study in children from Europe and rural Africa. *Proc Natl Acad Sci*
 394 *U S A* **107**: 14691-14696.
- 395 de Muinck EJ, Øien T, Storrø O, Johnsen R, Stenseth NC, Rønningen KS & Rudi K (2011)
 396 Diversity, transmission and persistence of *Escherichia coli* in a cohort of mothers and their
 397 infants. *Environmental Microbiology Reports* **3**: 352-359.
- 398 Edgar RC, Haas BJ, Clemente JC, Quince C & Knight R (2011) UCHIME improves sensitivity
 399 and speed of chimera detection. *Bioinformatics* **27**: 2194-2200.
- 400 Favier CF, de Vos WM & Akkermans AD (2003) Development of bacterial and bifidobacterial
 401 communities in feces of newborn babies. *Anaerobe* **9**: 219-229.
- 402 Ginzinger DG (2002) Gene quantification using real-time quantitative PCR: an emerging
 403 technology hits the mainstream. *Exp Hematol* **30**: 503-512.
- 404 Jabot F & Chave J (2011) Analyzing tropical forest tree species abundance distributions using
 405 a nonneutral model and through approximate Bayesian inference. *Am Nat* **178**: E37-47.
- 406 Jeraldo P, Sipos M, Chia N, *et al.* (2012) Quantification of the relative roles of niche and neutral
 407 processes in structuring gastrointestinal microbiomes. *Proc Natl Acad Sci U S A* **109**: 9692-
 408 9698.
- 409 Jimenez E, Marin ML, Martin R, *et al.* (2008) Is meconium from healthy newborns actually
 410 sterile? *Res Microbiol* **159**: 187-193.
- 411 Koenig JE, Spor A, Scalfone N, *et al.* (2011) Succession of microbial consortia in the
 412 developing infant gut microbiome. *Proc Natl Acad Sci U S A* **108 Suppl 1**: 4578-4585.
- 413 Koren O, Goodrich JK, Cullender TC, *et al.* (2012) Host remodeling of the gut microbiome and
 414 metabolic changes during pregnancy. *Cell* **150**: 470-480.
- 415 Ley RE, Peterson DA & Gordon JI (2006) Ecological and evolutionary forces shaping microbial
 416 diversity in the human intestine. *Cell* **124**: 837-848.
- 417 Lozupone C & Knight R (2005) UniFrac: a new phylogenetic method for comparing microbial
 418 communities. *Appl Environ Microbiol* **71**: 8228-8235.
- 419 Magurran AE (2004) *Measuring biological diversity*. Blackwell Science Ltd.
- 420 Makino H, Kushiro A, Ishikawa E, *et al.* (2011) Transmission of intestinal *Bifidobacterium*
 421 *longum* subsp. *longum* strains from mother to infant, determined by multilocus sequencing
 422 typing and amplified fragment length polymorphism. *Appl Environ Microbiol* **77**: 6788-6793.
- 423 Marchesi JR (2011) Human distal gut microbiome. *Environ Microbiol* **13**: 3088-3102.
- 424 Marques TM, Wall R, Ross RP, Fitzgerald GF, Ryan CA & Stanton C (2010) Programming
 425 infant gut microbiota: influence of dietary and environmental factors. *Curr Opin Biotechnol* **21**:
 426 149-156.
- 427 Martin R, Heilig HG, Zoetendal EG, Jimenez E, Fernandez L, Smidt H & Rodriguez JM (2007)
 428 Cultivation-independent assessment of the bacterial diversity of breast milk among healthy
 429 women. *Res Microbiol* **158**: 31-37.
- 430 Martin R, Langa S, Reviriego C, *et al.* (2003) Human milk is a source of lactic acid bacteria for
 431 the infant gut. *J Pediatr* **143**: 754-758.

- 432 Nadkarni MA, Martin FE, Jacques NA & Hunter N (2002) Determination of bacterial load by
 433 real-time PCR using a broad-range (universal) probe and primers set. *Microbiology* **148**: 257-
 434 266.
- 435 Oien T, Storro O & Johnsen R (2006) Intestinal microbiota and its effect on the immune system-
 436 -a nested case-cohort study on prevention of atopy among small children in Trondheim: the
 437 IMPACT study. *Contemp Clin Trials* **27**: 389-395.
- 438 Palmer C, Bik EM, DiGiulio DB, Relman DA & Brown PO (2007) Development of the human
 439 infant intestinal microbiota. *PLoS Biol* **5**: e177.
- 440 Ronaghi M (2001) Pyrosequencing sheds light on DNA sequencing. *Genome Res* **11**: 3-11.
- 441 Rudi K, Moen B, Sekelja M, Frisli T & Lee MR (2012) An eight-year investigation of bovine
 442 livestock fecal microbiota. *Vet Microbiol* **160**: 369-377.
- 443 Satokari R, Gronroos T, Laitinen K, Salminen S & Isolauri E (2009) Bifidobacterium and
 444 Lactobacillus DNA in the human placenta. *Lett Appl Microbiol* **48**: 8-12.
- 445 Sekelja M, Rud I, Knutsen SH, Denstadli V, Westereng B, Naes T & Rudi K (2012) Abrupt
 446 temporal fluctuations in the chicken fecal microbiota are explained by its gastrointestinal origin.
 447 *Appl Environ Microbiol* **78**: 2941-2948.
- 448 Sim J & Wright CC (2005) The kappa statistic in reliability studies: Use, interpretation, and
 449 sample size requirements. *Physical Therapy* **85**: 257-268.
- 450 Skanseng B, Kaldhusdal M & Rudi K (2006) Comparison of chicken gut colonisation by the
 451 pathogens *Campylobacter jejuni* and *Clostridium perfringens* by real-time quantitative PCR.
 452 *Mol Cell Probes* **20**: 269-279.
- 453 Storro O, Oien T, Dotterud CK, Jenssen JA & Johnsen R (2010) A primary health-care
 454 intervention on pre- and postnatal risk factor behavior to prevent childhood allergy. The
 455 Prevention of Allergy among Children in Trondheim (PACT) study. *BMC Public Health* **10**:
 456 443.
- 457 Storro O, Oien T, Langsrud O, Rudi K, Dotterud C & Johnsen R (2011) Temporal variations in
 458 early gut microbial colonization are associated with allergen-specific immunoglobulin E but
 459 not atopic eczema at 2 years of age. *Clin Exp Allergy* **41**: 1545-1554.
- 460 Tauler R, Smilde A & Kowalski B (1995) Selectivity, Local Rank, 3-Way Data-Analysis and
 461 Ambiguity in Multivariate Curve Resolution. *Journal of Chemometrics* **9**: 31-58.
- 462 Wang Q, Garrity GM, Tiedje JM & Cole JR (2007) Naive Bayesian classifier for rapid
 463 assignment of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol* **73**:
 464 5261-5267.
- 465 Whitman WB, Coleman DC & Wiebe WJ (1998) Prokaryotes: the unseen majority. *Proc Natl*
 466 *Acad Sci U S A* **95**: 6578-6583.
- 467 Yatsunenکو T, Rey FE, Manary MJ, *et al.* (2012) Human gut microbiome viewed across age
 468 and geography. *Nature* **486**: 222-227.
- 469 Yatsunenکو T, Rey FE, Manary MJ, *et al.* (2012) Human gut microbiome viewed across age
 470 and geography. *Nature* **486**: 222-227.
- 471 Zimonja M, Rudi K, Trosvik P & Næs T (2008) Multivariate curve resolution of mixed bacterial
 472 DNA sequence spectra: identification and quantification of bacteria in undefined mixture
 473 samples. *Journal of Chemometrics* **22**: 309-322.

474

475 **Tables**

476 Table 1

477 Cohan's kappa index – estimate of agreement in detection of a given bacteria in mothers and
 478 their infants. Calculations are based on detection of a given bacteria by RT-PCR.

Age	<i>B. fragilis</i>	<i>B. longum</i>	<i>Bifidobacterium</i>	<i>E. coli</i>
3 days	0.18	0.07	0.04	0.17
10 days	0.24	0	0.04	0.3
4 months	0.27	-0.03	-0.05	0.02
1 year	0.1	-0.02	-0.05	0.01
2 years	0.1	0	-0.04	-0.07

479

480

481 **Figures**

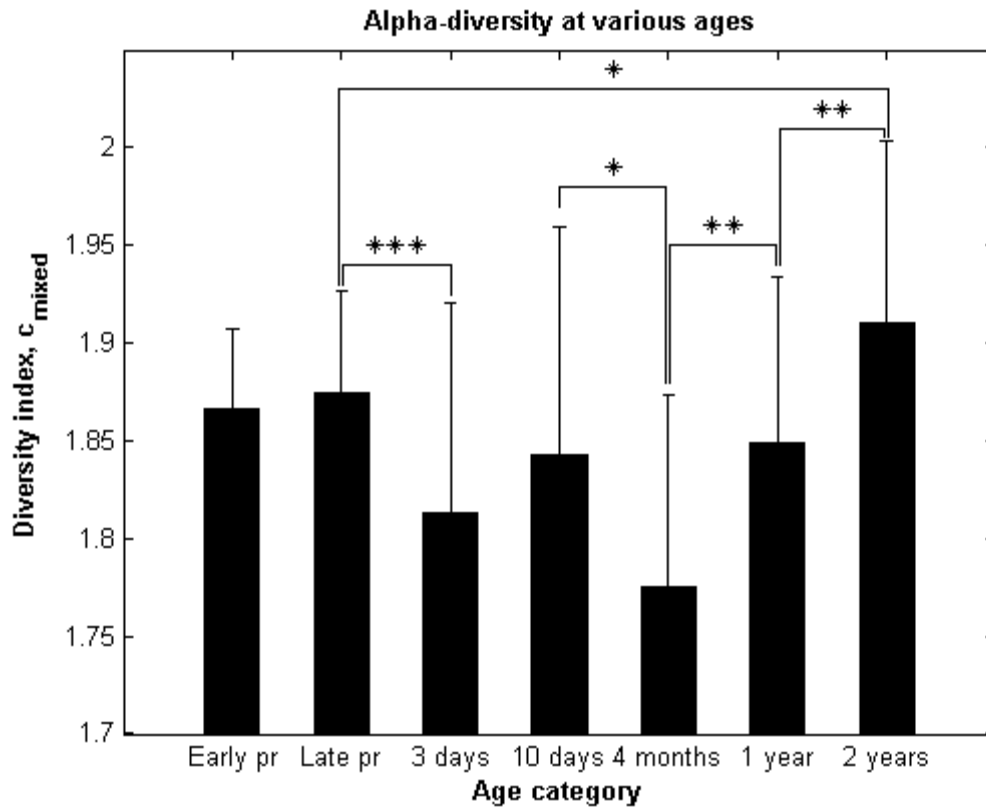
482 **Figure 1** Nucleotide diversity measurements. The significance in difference between diversity
 483 indices at two subsequent time points was calculated with the Friedman's (A and B) and Kruskal-
 484 Wallis (C) tests. * $p < 0.05$; ** $p < 0.01$ and *** $p < 0.001$. Early pr and Late pr: Early (8-20 weeks)
 485 and late (30-40 weeks) pregnancy periods, respectively. **A.** The modified Simpson's index of
 486 nucleotide spectra diversity c_{mixed} at various ages. **B.** The modified Bray-Curtis index of nucleotide
 487 dissimilarity (BC) between individuals at various ages. Early pr and Late pr: early (8-20 weeks) and
 488 late (30-40 weeks) pregnancy periods, respectively. **C.** The modified Bray-Curtis index of
 489 nucleotide dissimilarity (BC) between the subsequent time points. E-L pr: the period between early
 490 (8-20 weeks) and late (30-40 weeks) pregnancy periods; L pr – 3 d: comparison between 3 day-old
 491 newborns and their mothers during the late pregnancy stage; 3 d – 10 d: between 3 and 10 days of
 492 age; 10 d – 4 m: between 10 days and 4 months of age; 4 m – 1 y: between 4 months and 1 year of
 493 age; 1 y – 2 y: between 1 and 2 years of age. The error bars represent standard error of the mean.

494
 495 **Figure 2** Bacterial species composition in stool samples of infants (from 3 days to 2 years of age)
 496 and their mothers during pregnancy as revealed by MCR-ALS. Early pr and Late pr: early (8-20
 497 weeks) and late (30-40 weeks) pregnancy periods, respectively.

498
 499 **Figure 3** Summary of PARAFAC analysis on relative abundances of MCR-ALS resolved bacterial
 500 groups. C1, C2, C3 – PARAFAC components. Early pr and Late pr: early (8-20 weeks) and late
 501 (30-40 weeks) pregnancy periods, respectively. **A.** PARAFAC-suggested components C1, C2 and
 502 C3 represent *Bifidobacterium*, *Lachnospiraceae* and *Escherichia* components respectively. **B.**
 503 At early days of life, C1 and C3 determined the variation in the system, whereas at pregnancy,
 504 1 and 2 years of life, C2 became more important.

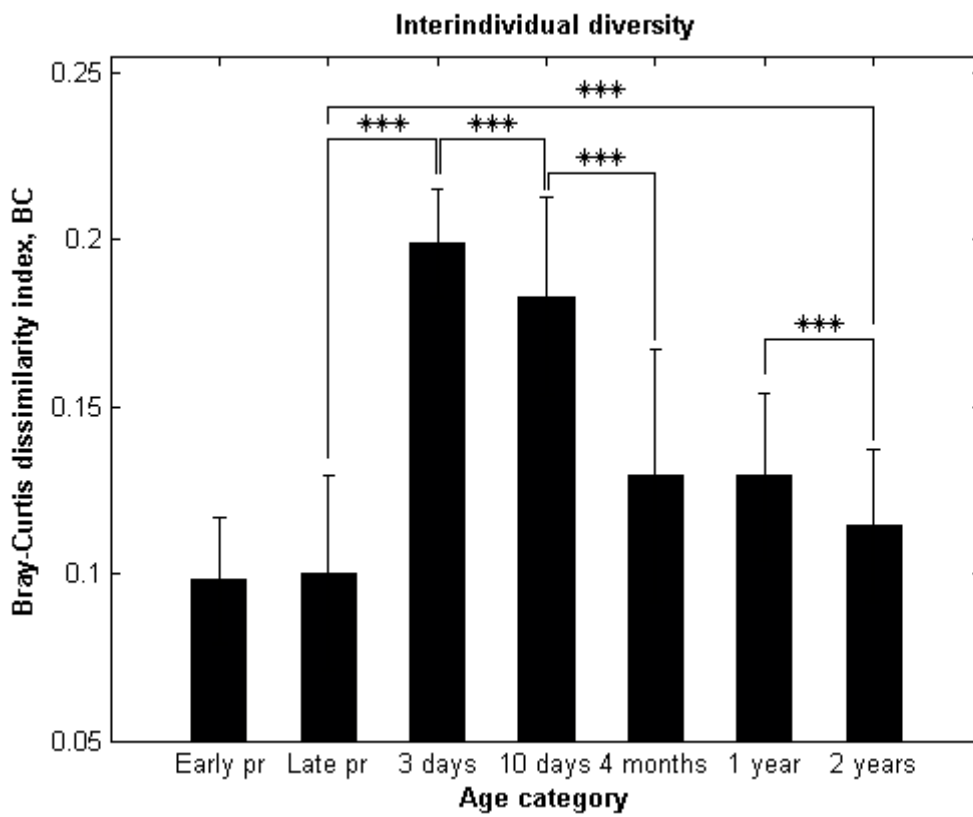
505
 506 **Figure 4** Prevalence of bacterial species in a population of children at various ages. Blue line
 507 indicates prevalence of bacteria in a subpopulation of children in whose mothers it was also
 508 detected; red line – in a total population of children of a given age. Black line depicts the percentage
 509 of individuals in who bacteria was detected both in a given and a previous time point compared to
 510 a total number of individuals where it was detected in a previous time point. Late pr: late (30-40
 511 weeks) pregnancy period. **one-sided binomial test p -value < 0.01 .

512
 513



514

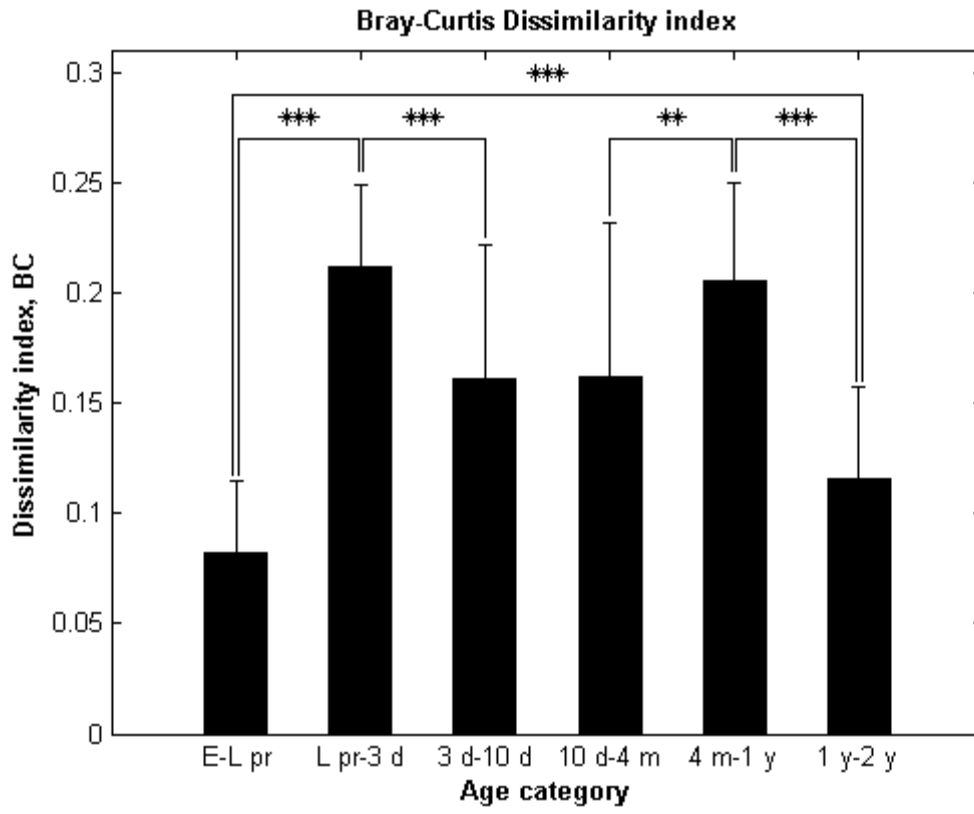
515 Figure 1A



516

517 Figure 1B

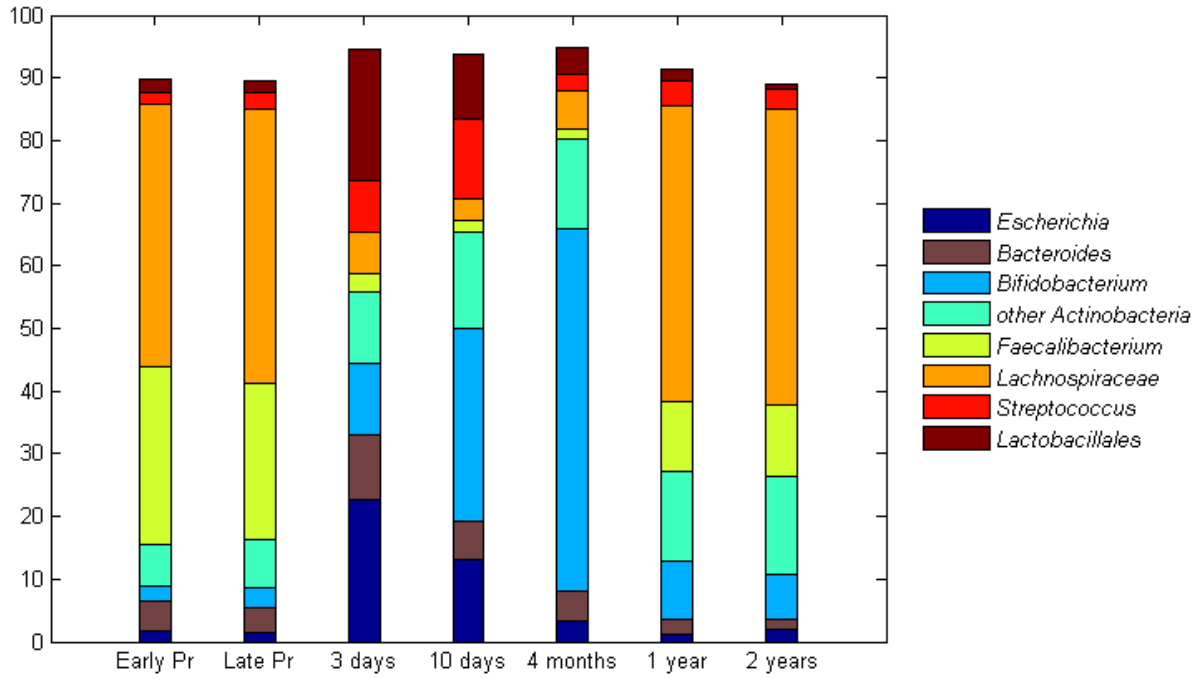
518



519

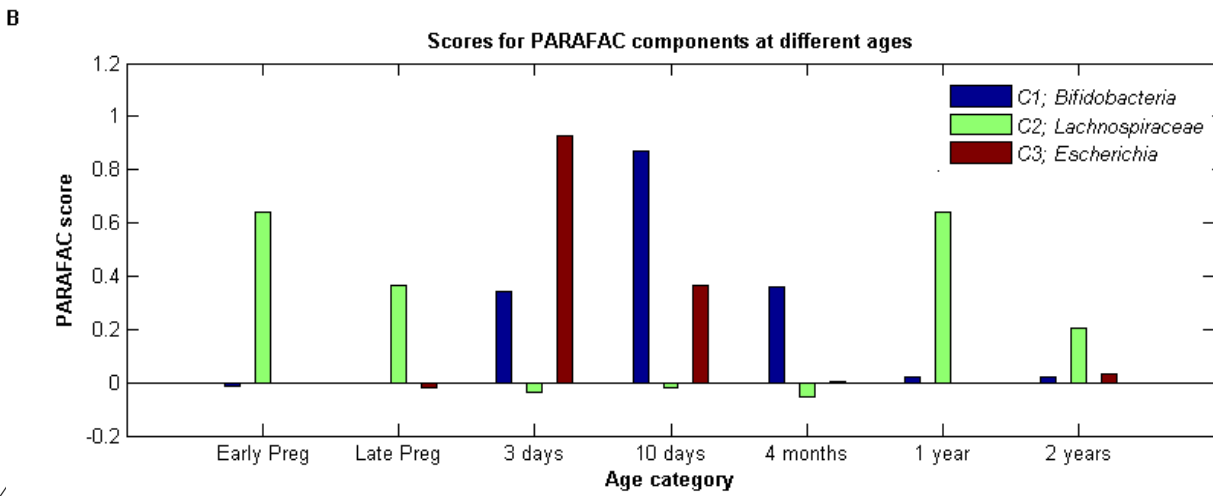
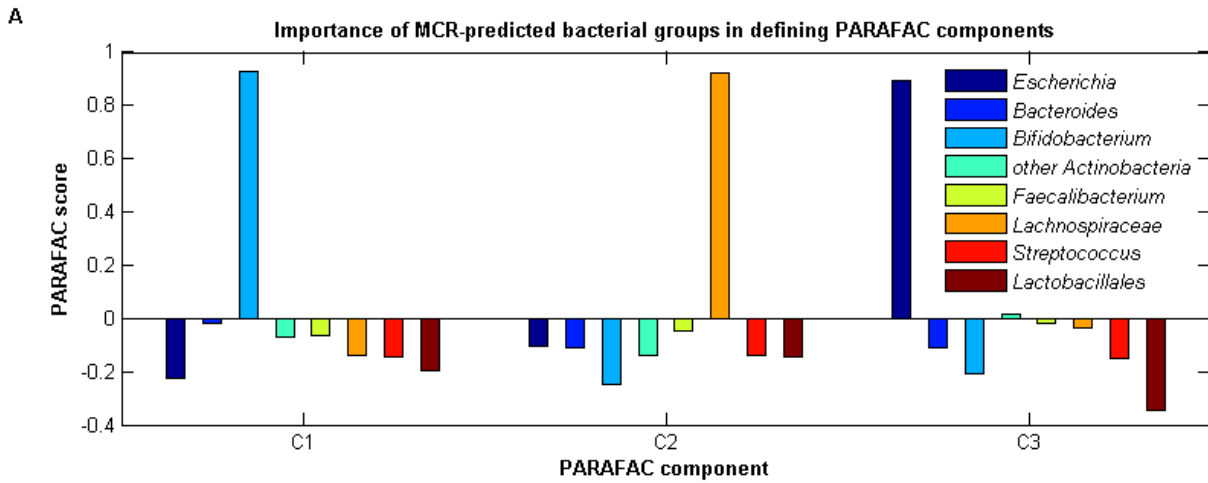
520 Figure 1C

521



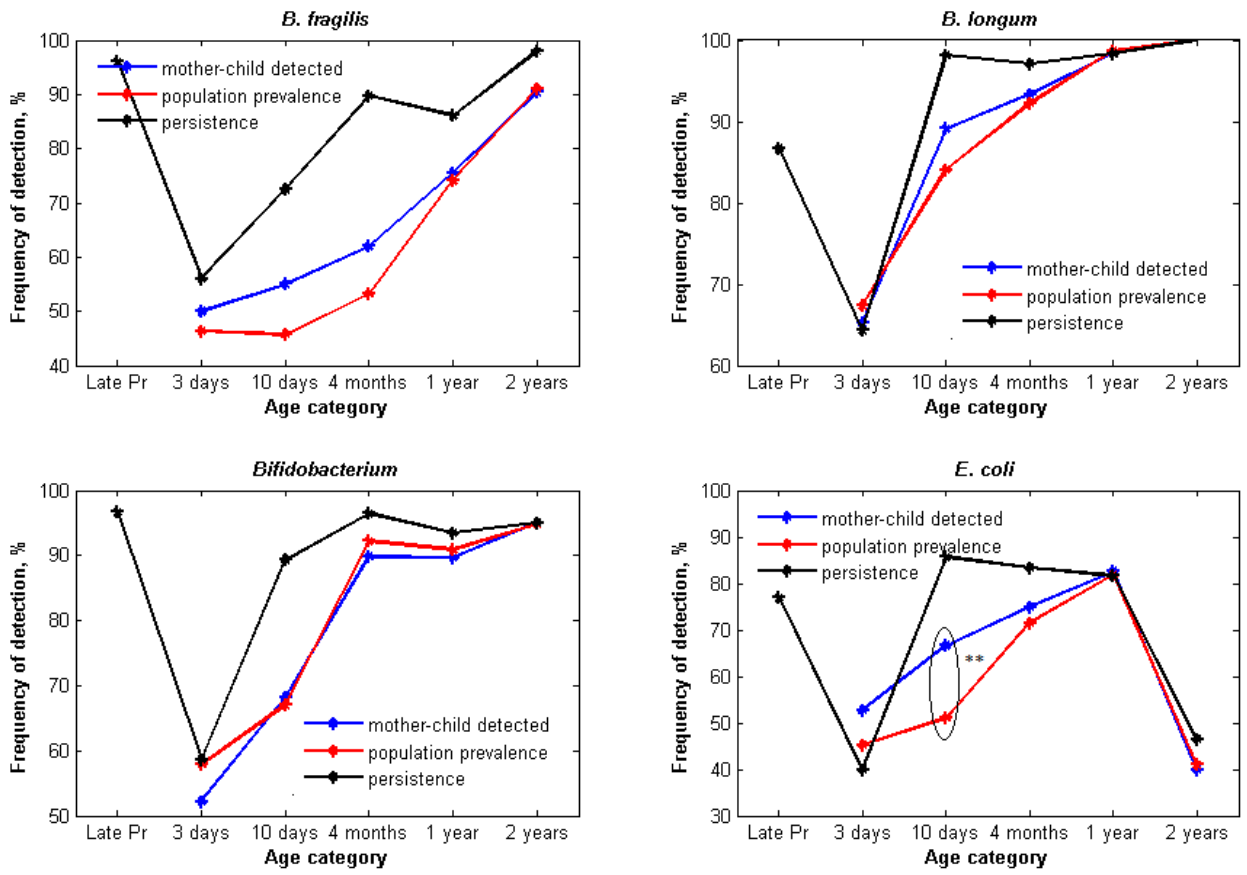
522

523 Figure 2



524

525 Figure 3



526

527 Figure 4