

1 **Longitudinal analysis of infant stool bacteria communities before and after**
2 **acute febrile malaria and artemether/lumefantrine treatment**

3
4 Rabindra K. Mandal¹, Rosie J. Crane^{2,3*}, James A. Berkley^{2,3}, Wilson Gumbi², Juliana
5 Wambua², Joyce Mwongeli Ngoi², Francis M. Ndungu², and Nathan W. Schmidt^{1*}

6
7 **SUPPLEMENTARY METHODS**

8 **Study site and sample collection**

9 Between August 2015 and January 2017, 100 infants resident within one contiguous 26x13km
10 zone comprising 12 villages in Kilifi County, Kenya, were observed from within 14 days of birth
11 until nine months of age. The primary purpose of this cohort study was to describe risk factors
12 for Environmental Enteric Dysfunction. Ethical approval for the cohort study was obtained from
13 KEMRI Scientific & Ethics Review Unit, Kenya (references 2983) and Oxford Tropical
14 Research Ethics Committee, UK (37-15).

15 Malaria transmission intensity in the area is between moderate and high. Malaria
16 episodes were identified through active and passive case detection. For active detection,
17 participants were visited at home once per week and axillary temperature measured. If fever was
18 present ($\geq 37.5^\circ\text{C}$), a capillary blood sample was obtained for both rapid diagnostic test (RDT,
19 CareStart kit; AccessBio) immediately at the participant's home, and parasite microscopy later
20 that day at the KEMRI Wellcome Trust Research Programme Clinical Trials Laboratory
21 (KEMRI-WTRP CTL). *Plasmodium* species and density were recorded. Participants with
22 positive RDT and/or slide microscopy were immediately started on a 3-day oral course of AL. If
23 the caretaker reported that the participant had been feverish, but temperature upon home visit

24 was $<37.5^{\circ}\text{C}$, repeat temperatures were measured at least twice during the following 24 hours. If
25 at any point the temperature exceeded 37.5°C , blood sampling, and treatment if malaria was
26 detected, was provided as detailed above. Passive case detection occurred in-between weekly
27 home visits when caretakers were encouraged to bring participants to the local primary
28 healthcare facility (Junju dispensary) in the event of fever during working hours, where they
29 would be assessed by the study clinician. Outside of working hours, participants were brought to
30 the home of their local study fieldworker. Antibiotic courses were also prospectively recorded on
31 the study database during weekly home visits by fieldworkers, by the study clinician at Junju
32 dispensary, and for any inpatient admissions by the cohort Principal Investigator (Dr. Rosie
33 Crane) at Kilifi County Hospital.

34 Stool samples were collected at home every 1-3 weeks by caretakers transferring feces
35 from disposable nappy to a sterile pot using a sterile spatula. Samples were transferred to cold
36 storage ($2-8^{\circ}\text{C}$) mostly within 1 hour then brought to Junju dispensary where four 0.5-2ml
37 aliquots were created then stored in a dry shipper at -198°C mostly within a further 1 hour.
38 Samples were transported in the dry shipper on a weekly basis to KEMRI-WTRP CTL where
39 they were transferred without thawing to -70°C storage.

40

41 **Selection of participants and stool samples**

42 Participants were selected for inclusion in this analysis if they had stool samples collected before
43 and after a confirmed clinical malaria episode(s) without antibiotics having been administered
44 between collection time points. Whereas multiple samples were collected for each of these
45 selected participants, only those stool samples meeting the criteria above were selected for
46 analysis, as shown in Figure 1A.

47

48 **DNA extraction, amplification and sequencing**

49 Stool samples were retrieved from -70°C storage and thawed on ice. From each sample, 200mg
50 was aseptically weighed into a container pre-loaded with ≈370mg of acid-washed 212-300 μM
51 glass beads (Sigma G1277-500G). DNA extraction was then carried out using Qiamp Fast DNA
52 Stool Mini Kit (Qiagen 51604). The manufacturer's protocol was modified in two ways to
53 enhance cell lysis: first by mechanical disruption through bead beating, and second through an
54 additional incubation of the lysate (95°C, 5 minutes) after addition of 1ml inhibitex buffer
55 followed by centrifugation (15000g, 1 minute 30 seconds) to pellet the stool particles. The
56 supernatant was then taken through further steps of inhibitor removal, purification and elution of
57 DNA using spin columns as per kit protocol. DNA yield was determined using Qubit 2.0
58 Fluorimeter/Qubit dsDNA HS Assay Kit as per protocol.

59 The V3- V4 hypervariable region of the 16S rRNA gene was targeted for sequencing. Primers
60 targeting this region were constructed with Illumina adapter overhang sequences added to the
61 gene-specific primer sequences. A region of 467 bp was targeted and amplified using the
62 following primers; 16S Amplicon PCR Forward Primer -
63 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCCTACGGGNGGCWGCAG and 16S
64 Amplicon PCR Reverse Primer -
65 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGACTACHVGGGTATCTAATCC.
66 Amplifications were done in 25 μl reactions with 12.5 μl Q5® Hot Start High-Fidelity 2X
67 Master Mix (NEB), 1 μl of 1μM forward and reverse 16S Amplicon PCR primer and 2.5 μl of
68 template. The reactions were performed on ABI Veriti thermocyclers (Applied Biosystems) under
69 the following conditions 95°C for 3 minutes, 25 cycles of; 95°C for 30 seconds, 55°C for 30

70 seconds, 72°C for 30 seconds, followed by 72°C for 5 minutes and a final hold at 4°C. The
71 amplified products were then verified on 1% agarose gel before purifying using Agencourt
72 AMPure XP beads (BeckmanCoulter). Libraries were prepared by ligating Illumina dual indices
73 and Illumina sequencing adapters to the prepared amplicons using the NexteraXT index kit
74 according to the Illumina 16S metagenomic sequencing library preparation protocol (Illumina).
75 Each library was then purified using Agencourt AMPure XP beads (BeckmanCoulter) and
76 thereafter the size and quantity was assessed using the Agilent 2100 Bioanalyzer (Agilent) and
77 the Qubit 2.0 fluorimeter (Life Technologies) respectively. The barcoded libraries were pooled at
78 equimolar concentration and 8 pM of the pooled library spiked with 5% Phix (v3) and then
79 sequenced on 300 PE Miseq run using the MiSeq® Reagent Kit v3 (600 cycle). 16S rRNA gene
80 sequences have been deposited in the NCBI short read archive under bioproject XXXXXX (note
81 to editor - to be populated upon acceptance).

82

83 **Sequence Data analysis**

84 The paired end sequenced 16S rRNA gene (V3-V4) amplicons were analyzed with QIIME 2
85 (qiime2-2017.8, <https://qiime2.org/>) using the command line interface. Data were analyzed with
86 three different approaches (Runs). In Run 1, full length forward and reverse reads were used with
87 140 bp overlapping region. Run 2 had trimmed forward (16 - 290 bp) and reverse reads (1 – 220
88 bp) based on the Phred quality score with 50 bp overlapping reads. And, Run 3 had high quality
89 forward reads (16 – 212 bp) spanning the complete V3 region of the 16S rRNA gene.

90 All plugins used for the data analysis were implemented in QIIME2. Raw sequencing
91 data were imported using import plugin, demultiplexed with demux plugin, and sequence quality
92 control and feature table were constructed with DADA2 plugin [1]. DADA2 corrects Illumina

93 amplicon sequence data where possible, remove phiX reads, and filters chimeric sequences. The
94 resulting feature table were visualized with feature-table plugin and tree for phylogenetic
95 diversity analyses were generated with alignment and phylogeny plugin. Alpha (within group)
96 and beta (between group) diversity analyses were performed using diversity plugin at sampling
97 depth of 30,000 sequences for Run 3 and 5000 sequence depth for Run 2. Principal coordinates
98 analysis (PCoA) plots were made with emperor plugin. Alpha diversity metrics used were OTUs
99 (richness), Shannon index (evenness and abundance), and pielou_e (evenness). Beta diversity
100 metrics were Bray Curtis, weighted UniFrac and unweighted UniFrac distance. Naive Bayes
101 classifier trained the classifier on the Greengenes 13_8 99% OTUs database on the trimmed 16S
102 sequence (V3-V4 region) that was amplified with forward primer (341F: 5'-
103 CCTACGGGNGGCWGCAG-3') and reverse primer (805R: 5'-
104 GACTACHVGGGTATCTAATCC-3') and taxonomic analysis was performed using feature-
105 classifier plugin. OTU interaction map was made using Cytoscape [2] and graphical
106 phylogenetic analysis plot was produced using GraPhlAn [3]. Metagenomic capacity were
107 predicted based on 16S rRNA gene using online tool Piphillin
108 (<http://secondgenome.com/solutions/resources/data-analysis-tools/piphillin/>) with default settings
109 [4]. Briefly, raw feature table (OTU Abundance Table) and representative sequence (req-seqs)
110 were uploaded and metagenomic capacity were predicted against Kyoto Encyclopedia of Genes
111 and Genomes (KEGG) database (KEGG; May 2017 version). In some cases heatmaps and PCoA
112 plots were drawn with ClustVis (<https://biit.cs.ut.ee/clustvis/>) [5]. GraphPad software (v 7.0b)
113 was used for statistical analysis and some data display.

114 Statistical significance for alpha and beta diversity between groups were performed with
115 linear mixed model effects model (LMEM) using QIIME2 longitudinal plugin for longitudinal

116 analysis. At first, differentially abundant bacterial taxa were screened with linear discriminant
117 analysis (LDA) effect size (LEfSe; <http://huttenhower.sph.harvard.edu/galaxy>) bioinformatics
118 pipeline [6]. Secondly, the significance level was verified using repeated measures (LMEM).
119 Analysis of differential metabolic pathways and genes before and after malaria episode/AL
120 treatment were performed using RNA-seq 2G online tool (<http://52.90.192.24:3838/rnaseq2g/>)
121 with default settings using two differential expression (DE) methods (Students T and DESeq2)
122 [7]. Log2FC (fold change) cutoff for differentially expressed KEGG orthologs was ≥ 1 and P
123 value ≤ 0.05 for both metabolic pathways and KEGG genes. Analysis includes all stool samples,
124 unless otherwise indicated when paired before and after malaria episode stool samples are
125 compared.

126

127 **SUPPLEMENTARY TEXT**

128 **Forward reads outperformed the joined reads**

129 Demultiplexing the raw sequencing reads from an amplicon spanning V3-V4 of the 16S rRNA
130 gene in all 44 stool samples (Figure 1A), plus an additional 4 samples that were not included in
131 the analysis owing to the identification of an intervening antibiotic treatment, produced an
132 average of 250,169.94 (SE \pm 6971.34) and median of 252,590.5 reads per sample (Minimum =
133 50,800 and Maximum = 335,554) totaling 12,008,157 reads (Supplementary Table 1). The 16S
134 rRNA gene sequencing data were analyzed by three different approaches with varying lengths of
135 forward and reverse reads based on the Phred quality score. In the first approach (Run 1) full
136 length forward and reverse read were used, in the second approach (Run 2) trimmed forward and
137 reverse reads were used, and in the third approach (Run 3) only trimmed high quality forward
138 reads were used (Supplementary Figure 1A-B, see Materials and Methods for more detail). After

139 the quality control steps using DADA2, Run 1 had very few reads (11.81 ± 2.54 per sample and
140 were not considered for downstream analysis), Run 2 had intermediate reads ($28,203 \pm 1,225$ per
141 sample), while Run 3 had the highest reads ($163,811 \pm 4,622$ per sample) that were assigned
142 taxonomically (Supplementary Figure 1C). Run 3 had significantly higher confidence ($0.91 \pm$
143 0.002) for taxonomically classified reads than Run 2 (0.88 ± 0.003) ($P < 0.0001$, Mann-Whitney
144 test, Supplementary Figure 1D). Consistent with the differential confidence in classification, Run
145 2 had a different profile than Run 3 at the phylum level (Supplementary Figure 1E-F), the ratio
146 of Firmicutes to Bacteroidetes was two times higher in Run 2 (9.52) than Run3 (4.72), and most
147 notably, Run 2 had a significantly higher portion of reads (16.97 ± 1.9 % per sample) classified
148 only to kingdom level (bacteria) than Run 3 (0.59 ± 0.1 %) ($P < 0.0001$, Mann-Whitney test,
149 Supplementary Figure 1E-F). Furthermore, there was a significantly higher number of
150 Observed_OTUs in Run 3 (104.1 ± 3.93) than Run 2 (62.81 ± 3.2) ($P < 0.0001$, Mann-Whitney
151 test, Supplementary Figure 1G). Consistent with the increased diversity of taxonomically
152 assigned reads in Run 3 (Supplementary Figure 1F) compared to Run 2 (Supplementary Figure
153 1E), the Shannon index and *pielou_e* were significantly lower in Run 3 compared to Run 2 ($P <$
154 0.0001 , Mann-Whitney test, Supplementary Figure 1H-I). Overall, Run 3 outperformed Run 2
155 with regards to Phred quality score, confidence in taxonomically classified reads and alpha
156 diversity. Error prone sequence correction and comparatively lower quality sequence reads at 3'
157 end of reads, a limitation of most DNA sequencers, might have contributed to the poor
158 performance of stitched reads (Run 2). Consequently, Run 3 was used for all the analysis
159 performed in this study.

160

161 **SUPPLEMENTARY TABLES**

162 Supplementary Table 1. Metadata file with read numbers and alpha diversity indices.

163 Supplementary Table 2. Antibiotic courses.

164 Supplementary Table 3. Malaria episode characteristics.

165 Supplementary Table 4. Alpha diversity analysis.

166

167 **SUPPLEMENTARY FIGURES**

168 **Supplementary Figure 1** Quality analysis of 16S rRNA gene sequencing. V3 and V4 region

169 were subjected to MiSeq sequencing. Sequence quality control and feature table were

170 constructed with DADA2 implemented inside QIIME2. Three runs were performed with varying

171 read length depending on the Phred quality score. Run 1: Full length forward and reverse reads

172 with 140 bp overlapping region. Run2: Forward reads (16 – 290 bp) and reverse reads (0 – 220)

173 were trimmed for low quality reads conserving 50 bp overlapping reads. Run 3: High quality

174 forward reads (16 – 212 bp). A and B) Overall Phred score of forward and reverse read

175 respectively. C) Reads taxonomically classified by DADA2 with three different runs. D)

176 Confidence of reads classified by DADA2. Box denotes 25th and 75th percentile with median in

177 between and whisker denotes lowest and highest value. Data were analyzed by the Mann-

178 Whitney test. E and F) Taxonomic classification at phylum level. Alpha diversity of Run 2 and

179 Run 3 indicated by Observed_OTUs (G), Shannon Index (abundance + evenness) (H) and

180 pielou_e (evenness) (I). G-I) Individual samples from each run are shown along with the mean ±

181 S.E. Data were analyzed by the Mann-Whitney test.

182

183 **Supplementary Figure 2** OTU interaction maps show no distinct pattern between paired before

184 and after malaria episode/AL treatment samples. A) Interaction between 24 paired samples

185 (Figure 1) and 997 sequence variants (SV). Nodes (SVs) shared by the most samples are placed
186 at the core of map as indicated by edge length from the samples. B) Overall top 300 SVs shared
187 between the before and after malaria episode/AL treatment stool samples in Kenyan infants
188 which represents the core of interaction map. C) Relative abundance of sequence variants shown
189 in B and the remaining 697 SVs.

190

191 **Supplementary Figure 3** No effect of infant age on alpha diversity. Pearson correlation of age
192 versus Observed_OTUS (A,C) and age versus Shannon Index (B,D) in overall 44 stool samples
193 (A,B) and infants that had not and had antibiotics course (C,D).

194

195 **Supplementary Figure 4** Phylogenetic graph shows no obvious difference in the microbiota
196 taxonomic composition between before and after malaria episode/AL treatment. Taxonomic
197 cladogram represents the phylogenetic analysis before malaria/AL treatment (A), after
198 malaria/AL treatment (B), and combined before and after malaria episodes/AL treatment (C).
199 Size of nodes correlates with their relative abundance and different colors indicate different
200 clades. The graph was produced using Graphical Phylogenetic Analysis (GraPhlAn) tool.

201

202 **Supplementary Figure 5** Minimal taxonomic difference between before and after Malaria/AL
203 treatment. A) Cladogram showing discriminant features at kingdom, phylum, class, order, family
204 and genus level. Rings are arranged according to the taxonomic level. Outermost- Genus and
205 inner most –Kingdom. B) Bar graph shows the fold change of differentially abundant features.
206 Alpha value for the pairwise Wilcoxon test during LEfSe analysis was set to 0.05. Threshold on

207 the logarithmic LDS score for discriminative features were set to 2 and 0.2. C and D).

208 Differentially abundant features re-evaluated using LMEM.

209

210 **Supplementary Figure 6** Taxonomic assignment at genus level shows no distinct clustering
211 between the before and after malaria episode/AL treatment stool samples. A) Heat map of
212 overall top 15 genera. B) PCoA plot based on top 15 genera. MS: Malaria status; B- Before
213 malaria; A- After malaria; PID: Participant ID.

214

215 **Supplementary Figure 7** Linear mixed effect model analysis of sequence variants identified by
216 LEfSe analysis (see Figure 5). A) LMEM analysis of sequence variants (SV38 and SV43) that
217 were identified as being overly abundant in before malaria/AL stool samples. B) LMEM analysis
218 of sequence variants (SV29 and SV72) that were identified as being overly abundant in after
219 malaria/AL stool samples.

220

221 **Supplementary Figure 8** Predicted metagenomic capacity at KEGG ortholog (KO) level
222 identifies minimal differences between the before and after malaria episode/AL treatment stool
223 samples. Volcano plot of KO with DeSeq2 (A) and student's T-test (B). C) Overlapping KO
224 identified by DeSeq2 and T-test. D) PCoA plot and E) Heat map and F) Pathway enrichment of
225 the 47 overlapping KOs identified by DeSeq2 and T-test.

226

227 **Supplementary Figure 9** KEGG orthologs (KOs) of the A) N-Glycan biosynthesis pathway and
228 B) Histidine metabolism pathway having differential abundance in the before and after malaria
229 episode/AL treatment stool samples (see Figure 6 and Supplementary Figure 8). The predicted

230 KOs are visualized on KEGG pathway constructed using Pathview. The annotation for the graph
231 are same as the online source (https://pathview.uncc.edu/overview#kegg_view).

232 **SUPPLEMENTARY REFERENCES**

233

234

- 235 1. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: high-
236 resolution sample inference from Illumina amplicon data. *Nature methods* **2016**; 13(7):
237 581-3.
- 238 2. Shannon P, Markiel A, Ozier O, et al. Cytoscape: a software environment for integrated
239 models of biomolecular interaction networks. *Genome research* **2003**; 13(11): 2498-504.
- 240 3. Asnicar F, Weingart G, Tickle TL, Huttenhower C, Segata N. Compact graphical
241 representation of phylogenetic data and metadata with GraPhlAn. *PeerJ* **2015**; 3: e1029.
- 242 4. Iwai S, Weinmaier T, Schmidt BL, et al. Piphillin: Improved Prediction of Metagenomic
243 Content by Direct Inference from Human Microbiomes. *PloS one* **2016**; 11(11):
244 e0166104.
- 245 5. Metsalu T, Vilo J. ClustVis: a web tool for visualizing clustering of multivariate data using
246 Principal Component Analysis and heatmap. *Nucleic acids research* **2015**; 43(W1):
247 W566-W70.
- 248 6. Segata N, Izard J, Waldron L, et al. Metagenomic biomarker discovery and explanation.
249 *Genome biology* **2011**; 12(6): R60.
- 250 7. Zhang Z, Zhang Y, Evans P, Chinwalla A, Taylor D. RNA-Seq 2G: Online Analysis Of
251 Differential Gene Expression With Comprehensive Options Of Statistical Methods.
252 *bioRxiv* **2017**: 122747.
- 253