# SNV-FEAST: microbial source tracking with single nucleotide variants

Leah Briscoe[1]\*, Eran Halperin[2,3,4,5,6], Nandita R. Garud[3,7]\*

1. Bioinformatics Interdepartmental Program, University of California Los Angeles, Los Angeles, Los Angeles, CA, United States of America
2. Department of Computer Science, University of California Los Angeles, Los Angeles, CA, United States of America
3. Department of Human Genetics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA, United States of America
4. Department of Computational Medicine, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA, United States of America
5. Department of Anesthesiology and Perioperative Medicine, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA, United States of America
6. Institute of Precision Health, University of California Los Angeles, CA, United States of America
7. Department of Ecology and Evolutionary Biology, University of California Los Angeles, CA, United States of America

\*Correspondence to leahpbriscoe@gmail.com and ngarud@ucla.edu

22    **ABSTRACT**

23    Elucidating the sources of a microbiome can provide insight into the ecological dynamics

24    responsible for the formation of these communities. "Source tracking" approaches to date

25    leverage species abundance information, however, single nucleotide variants (SNVs) may be

26    more informative because of their high specificity to certain sources. To overcome the

27    computational burden of utilizing all SNVs for a given sample, we introduce a novel method to

28    identify signature SNVs for source tracking. We show that signature SNVs used as input into a

29    previously designed source tracking algorithm, FEAST, can more accurately estimate

30    contributions than species and provide novel insights, demonstrated in three case studies.

31

32    **Keywords**: Source tracking, microbiome, single nucleotide variants, transmission, strains

33

34    **BACKGROUND**

35         Understanding the sources that could contribute to the formation of a given microbiome

36    is of great interest in elucidating the ecological processes that give rise to these complex

37    communities and the impact of these communities on human and environmental health. For

38    example, a hospital environment may introduce antibiotic resistance genes to an infant gut

39    microbiome, and local selective pressures may result in vastly different microbial compositions

40    in different parts of the ocean. Approaches for determining the proportion of a microbiome of

41    interest (the "sink") that is attributed to different microbiomes (the "sources") is known as

42    "source tracking" (Knights et al., 2011; Shenhav et al., 2019). Source tracking is useful for

43    forensics, categorization of samples, detecting contamination, and tracing transmissions between

44    different hosts or environments. While source tracking was developed as a way to quantitatively

45    characterize a sample based on a set of samples with known origin, in most studies, the true

46    source of samples may never be collected. In these cases, source tracking approaches are useful

47    in identifying similarities between microbiome samples even if they cannot be used to

48    definitively identify the true source of origin.

49         Current approaches for source tracking include the Bayesian approach, SourceTracker

50    (Knights et al., 2011) and more recently the expectation-maximization approach, FEAST

51    (Shenhav et al., 2019). These source tracking methods use species abundance profiles of the

52    sample of interest (the sink) and of potential sources and compute percentages of sinks that are

53    attributable to each potential source. However, species abundance profiles miss important sub-

54    species single nucleotide variants (SNVs), which may provide higher resolution information than

55    species about transmission patterns. For example, (Nayfach et al., 2016) found that the sharing of

56    microbiome SNVs private to mothers and their infants decreases over the first year of the

57    infant's life while species sharing increases. This suggests that while the infant microbiome

58    increasingly resembles the adult microbiome ecologically, sources other than the mother also

59    colonize the infant. Thus, species-level resolution may obscure true sources of microbes while

60    SNVs can reveal actual transmissions to the infant.

61         While tracking strain transmissions with SNVs has been highly successful in a number of

62    studies (Asnicar et al., 2017a; Ferretti et al., 2018; Korpela et al., 2018; Li et al., 2016; Nayfach

63    et al., 2016; Olm et al., 2021; Schmidt et al., 2019) current approaches to strain tracking are

64    limited. These methods provide binary information by inferring whether or not a strain

3

65    transmission has occurred per species but they do not shed light on the relative proportions of

66    microbiomes that are similar. A specific example of this is inStrain (Olm et al., 2021) which

67    computes a pairwise population-level average nucleotide identity (popANI) between two

68    samples. If an infant harbors several strains derived from the mother at low frequency, these

69    shared strains will have high popANI values, but they will represent a relatively small proportion

70    of the infant's microbiome. By contrast, source tracking allows us to simultaneously infer the

71    putative proportions for multiple sources contributing to a given sink, integrated over all

72    community members in the sink. As shown in **Figure 1**, one may be able to estimate that an

73    infant microbiome is  explained 25% by the mother, 10% by the dog, and 30% by unknown

74    sources (Knights et al., 2011; Shenhav et al., 2019). In other words, source tracking with SNVs

75    leverages not only the genetic variants within species, but also the relative abundances of the

76    species that carry the SNVs.

77          Here, we evaluate whether source contributions estimated with SNVs are more accurate

78    than with only species when provided as inputs to FEAST (Shenhav et al., 2019) (hereafter

79    referred to as SNV-FEAST and species-FEAST, respectively). FEAST (Shenhav et al., 2019) is

80    faster and more accurate than previous source tracking tools (Knights et al., 2011), and therefore,

81    is ideal for adaptation to SNV source tracking since it can accept larger numbers of features and

82    input sources. Despite this improved computational efficiency, the potentially millions of single

83    nucleotide variants across all microbiome species in a given host still can computationally

84    overwhelm FEAST. To address this, we introduce a novel approach to determine signature SNVs

85    that can be used as input to FEAST. This both reduces memory requirements and computation

86    time in the FEAST estimation, allowing us to optimally estimate the source contribution of a

87    sink. We find that SNV-FEAST and species-FEAST yield different outcomes when applied to

88    simulated data, with SNV-FEAST frequently out-performing species-FEAST. We apply SNV-

89    FEAST to three real-world case studies, including source tracking between infants and their

90    mothers in the first year of life, between infants and the neonatal intensive care unit (NICU), and

91    between oceans around the world. We confirm the ability of SNV-FEAST by recapitulating

92    several previously published findings in our case studies, as well as discover new source tracking

93    patterns across oceans. In sum, we show that SNVs can be used to estimate potential

94    transmissions across hosts and across environments.

95

96  **RESULTS**

97

98  **SNV-FEAST algorithm**

99  Here we adapt FEAST to accept SNV abundance instead of species abundance as input.

100  A computational challenge in using SNVs instead of species as input to FEAST is that SNVs

101  contribute a significantly larger feature space. The number of different species comprising a

102  microbiome can range from a few hundred to a few thousand, while the number of possible

103  SNVs for a given species alone can be in the thousands (Schloissnig et al., 2013). This difference

104  in number of input features can result in FEAST runtimes that last several hours instead of a few

105  minutes and memory intensive storage of read counts at all sites of variation.

106  We devised a likelihood-based approach for selecting a set of informative or "signature"

107  SNVs for a given source tracking analysis, allowing us to overcome the time and memory

108  intensive challenges of utilizing SNV-level data. We identify these informative SNVs by

109  computing a signature score (**Figure 1A**) (see **Methods)** that quantifies the extent to which

110  SNVs in the sink are most likely derived from one of the potential sources. This is analogous to

111  identifying SNVs private to sources and their sinks, but more generalized to include SNVs that

112  may be found in multiple sources, albeit at higher frequency in one of the potential sources (see

113  **Methods**).

114  To compute a signature score for a given SNV, two hypotheses are compared for each

115  potential source: (1) that one source solely explains the observed allele counts in the sink and (2)

116  all sources except that one source collectively explain the observed allele counts in the sink. For

117  each hypothesis, we calculate the binomial log-likelihood for the estimate of the allele frequency

118  in the sink, θ.

119  **Hypothesis 1:** Source $i$ with allele frequency $p_i$ explains the allele counts in the sink.

$$\hat{\theta} = p_i$$

120  **Hypothesis 2:** A combination of all other sources except $i$ (sources $j \neq i$) explain the observed

121  allele count distribution in the sink. The estimate of the sink allele frequency is computed using a

122  mixture of the allele frequencies $p_j$ from those sources. The mixing parameter $\alpha_j$ is learned using

123  Sequential Least Squares Programming with the constraint that $\sum_{j \neq i} \alpha_j = 1$.

124

125    The binomial log-likelihood is calculated as follows, where there are *n* reads with the reference

126    allele and *m* reads with the alternative allele in the sink.

127

128

129    A log likelihood ratio representing the support for hypothesis 1 relative to hypothesis 2 is

130    calculated per site per potential source. The maximum log likelihood ratio per site is the

131    signature score for that SNV, representing how favorably one of the sources explains the sink

132    over all other sources. Signature SNVs are those with scores greater than two standard deviations

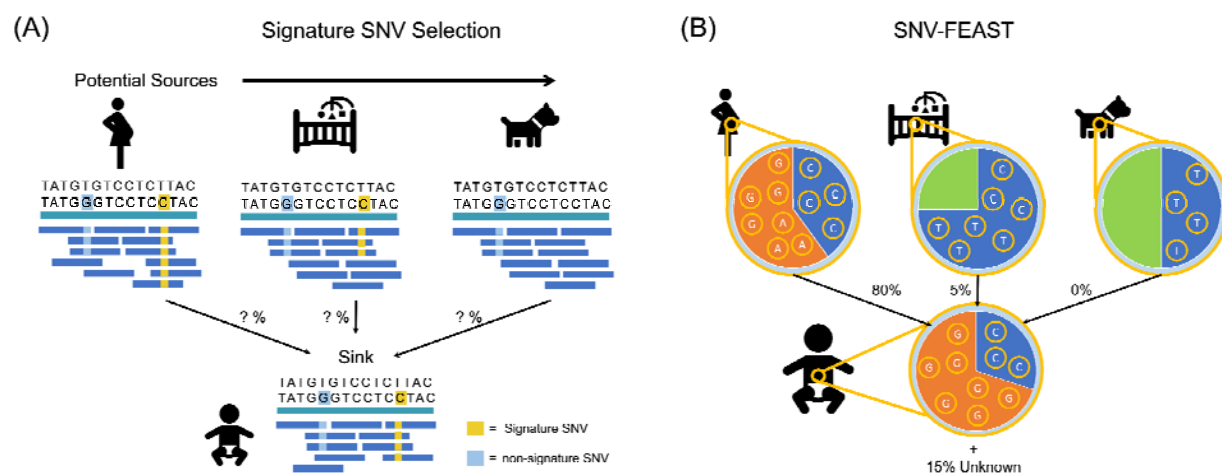133    over the mean signature score computed for all SNVs (**Methods**).

134

135

136



137    **Figure 1**: Signature SNV selection and SNV-FEAST. (A) A signature SNV is present in one or

138    few but not all sources. By contrast, a non-signature SNV is generically present in multiple

139    sources and thus provides little discriminating information.  (B) SNV-FEAST estimates the

140    proportion a given sink derived from various sources using the read counts for each allele in

141    sinks and sources.

142

143

144    **Evaluation of SNV-FEAST in simulations**

6

145   To compare the accuracy of species-FEAST and SNV-FEAST, we performed simulations

146 mimicking mother-infant transmissions with the goal of estimating contributions of different

147 sources to an infant sink. Our simulations tested the ability of SNVs and species to recapitulate

148 the true source composition in synthetic samples comprised of a mixture of reads drawn from

149 multiple real fecal adult samples. To construct these synthetic infant microbiomes, we mixed

150 metagenomic data from mothers sampled in a mother-infant dataset (Bäckhed et al., 2015) at

151 various proportions as described below (**Methods**).

152   The difficulty of source tracking increases with the number of contributing sources

153 (Shenhav et al., 2019). Thus, we simulate infants that have a small ($<=5$) versus large ($6 - 10$)

154 number of contributing sources (**Supplementary Table 1**), including an unknown source (e.g. a

155 randomly selected unrelated mother). Known source contributions to the simulated gut

156 microbiome sample of the infant were varied between 1 and 90% while the unknown

157 contribution varied between 10 and 90%. The unknown source was not presented to FEAST as a

158 potential known source.

159   Additionally, not all species in a mother are transmitted to the infant (Asnicar et al.,

160 2017b; Ferretti et al., 2018; Korpela et al., 2018; Sprockett et al., 2020; Yassour et al., 2018).

161 Thus, in our simulations, species transmission rates were determined using a beta distribution,

162 which is a natural model for values between $(0,1)$ and often proposed for microbial abundance

163 data (E. Z. Chen & Li, 2016; Martin et al., 2020; Sloan et al., 2006, 2007) (see **Methods**). We

164 therefore consider four simulated scenarios: a combination of low versus high number of sources

165 and low versus high transmission rates (see **Methods**).

166   **Figure 2** compares the performance of SNV-FEAST and species-FEAST in estimating

167 the true contribution of sources. FEAST using SNVs has equal if not better performance than

168 species in most scenarios, and performs especially well when transmission rates are low and

169 unknown source proportions are high. SNVs have a lower root mean squared error (RMSE)

170 compared to species in three of the four scenarios and higher Pearson correlation between true

171 and estimated contributions in all four scenarios. The difference in these correlations for SNVs

172 versus species is significant in all four cases when using a paired Wilcoxon signed rank test (high

173 transmission: p-value = 0.00560, 0.00251 for small and large number of sources, low

174 transmission: p-value = 0.00024, 0.002340 for small and large number of sources). These results

175 suggest that SNVs may offer useful signatures of transmission.

176

177



**Figure 2: Ability of SNV and species-FEAST to recapitulate true contributions in simulations.** Estimated known and unknown source proportions for infant microbiomes simulated with in silico mixtures of real maternal fecal microbiomes under different scenarios: either low number of contributing sources (<=5) or high number of sources (6-11), and high transmission rate of species or low transmission rate. Transmission rate is the probability of an infant being colonized by a given species, simulated using a beta distribution centered on the relative abundance of species in sources (**Methods**). 23 infants were simulated with five or fewer sources and 19 infants were simulated with a large number of sources (**Table S1**). The black line indicates the ground truth for proportions. For each simulated infant, there are 11 points plotted,

188    whereby 10 correspond to known sources (some of which have zero contribution), and one

189    corresponds to an unknown source which are indicated by a hollow circles in the plot.

190

191         To assess whether all species and all signatures SNVs in the sink are needed for accurate

192    source tracking, we varied the proportion of species (from 10%, 50% or 100%) and SNVs (from

193    10%, 50% or 100%) included as inputs to the algorithm (**Figure S1**). We used Pearson

194    correlation between the true and estimated proportions to represent accuracy of SNV-FEAST.

195    When decreasing the percentage of SNVs used, there is no statistically significant change in the

196    performance. However, when decreasing the percentage of species used, there are statistically

197    significant decreases in the performance (**Figure S1**).

198         To illustrate the advantage of SNV-FEAST over traditional strain tracking approaches

199    such as inStrain (Olm et al., 2021), we used the same synthetic communities produced in the

200    above simulation for inStrain profiling between each infant and each of their potential

201    contributing sources (**Figure S2**).  InStrain computes a popANI score, which represents the

202    average nucleotide identity between two different metagenomic samples for a given species. As

203    per the inStrain paper, popANI values > 99.999% represent the same strain for that species being

204    shared between samples (**Methods**). However, this approach provides a binarization as to

205    whether or not a strain was transmitted, and does not account for the relative abundance of the

206    strain in the sink. Thus, we computed the fraction of each infant's species that have popANI

207    ≥99.999%, with each potential source.

208         As expected, both SNV-FEAST and inStrain produce estimates of sharing that correlate

209    positively with the ground truth mixture proportions of the contributing source samples in each

210    infant (**Figure S2**). We found inStrain results yielded a 0.742 Pearson correlation ($p < 1 \times 10^{-12}$)

211    with the true mixture proportions, whereas SNV-FEAST has a 0.866 Pearson correlation ($p < 1 \times$

212    $10^{-12}$) with the true proportions. The higher correlation values for SNV-FEAST likely reflect that

213    relative abundances of strains and their genomic identities are simultaneously taken into account

214    for source tracking, whereas inStrain only accounts for genomic identities. Finally, several of the

215    shared species in the simulations had popANI values < 99.999%, reflecting the complex

216    mixtures from multiple sources.

217         We next compared SNV-FEAST with the strain tracking procedure in Nayfach et al.

218    2016. Again, we used the same synthetic communities produced in the simulation to determine

219    marker alleles as defined in Nayfach et al. 2016 (**Methods**). Here a marker allele is determined

220    to be a SNV that is private to mother, infant, or the mother-infant dyad, and absent from the

221    background population, which consisted of other samples in the dataset as well as samples from

222    United States adults in the Human Microbiome Project (**Methods**). Species with $\geq 5\%$ marker

223    allele sharing between mother and infant were deemed to share a strain (**Methods**). We found a

224    high correlation between the true mixture proportions (on x-axis) and the percentage of species

225    with transmission events (y-axis) (Pearson correlation 0.915 , $p < 1 \times 10^{-16}$). The higher

226    correlation for the Nayfach et al. 2016 approach compared to the inStrain approach possibly

227    reflects horizontal gene transfers between lineages residing in infants and mothers. By contrast,

228    there was a lower correlation between the true mixture proportions (x-axis) and the sharing for

229    all marker alleles across species present in the infant (y-axis) and (0.575 Pearson correlation, $p <$

230    $1 \times 10^{-16}$) (**Figure S3B**).

231

**Source tracking in infants over the first year of life**

233    Having assessed the abilities of SNV-FEAST in synthetic data, we next estimated the

234    contribution from the true mother over time to the true infant with SNV and species-FEAST in

235    the Backhed et al. 2015 dataset. This dataset is composed of metagenomic samples from infants

236    collected at four days, four months, and 12 months after birth, as well as their mothers at the time

237    of delivery.  Previous analyses on this data have shown that even while species similarity

238    increases, infants and their mothers share fewer proportions of strains over time as revealed by

239    sharing of SNVs private to mother-infant dyads (Nayfach et al., 2016). Thus, SNVs belonging to

240    strains shared only by the infant and their mother may be more informative of the true source

241    compared to species. Here we sought to test whether SNV and species-FEAST recapitulate these

242    results (**Methods**).

243    In applying FEAST to the Backhed et al. 2015 dataset, we estimated the proportion of

244    infant at birth attributable to mother. For 4 month infants, we estimated the proportion

245    attributable to the mother and itself at birth. For 12 month infants, we estimated the proportion

246    attributable to the mother and itself at birth and four months (Shenhav et al. 2019). This allowed

247    "unknown" to be more strictly defined as the component of the infant microbiome that could not

248    be explained by the mother. It also allowed us to better discern if completely new strains were

10

249   acquired at the 4[th] and 12[th] months of life (that were not already acquired during previous life

250   stages).

251      First, consistent with previous findings made with species and SNVs (Nayfach et al.,

252   2016), species-FEAST estimates an increasing contribution of the mother over time (t-test p-

253   value = 5.1 x 10[-4]), but SNV-FEAST estimates a decrease over time (p-value = 0.063) (**Figure**

254   **3**).

255      Second, we assessed the ability of species and SNV-FEAST to distinguish the true

256   mother from three randomly selected unrelated mothers. Species-FEAST estimates an increasing

257   contribution of unrelated mothers over time (t-test p-value = 0.014) while SNV-FEAST

258   estimates no significant change over time (t-test p-value = 0.59) (**Figure 3**). The increase in

259   contribution from unrelated mothers with species-FEAST does not suggest that these particular

260   unrelated mothers are seeding the infant. Rather, the opposing trend observed with SNVs

261   suggests that similarity at the species level is consistent with the maturation of the infant

262   microbiome over time.

263      Finally, we estimated contributions from unknown sources, i.e. the proportion of the

264   infant microbiome not explainable by the true mother, the three randomly selected unrelated

265   mothers, or any previous time point. Species-FEAST estimates a sharp decline in contribution of

266   unknown sources over the first year of life (t-test p-value =7.1 x 10[-12]) (**Figure 3**). This

267   significant decrease in unknown at the species level reflects the infant microbiome maturation

268   over the first year of life. By contrast, SNV-FEAST estimates little change in the contribution of

269   unknown sources (t-test p-value = 0.49) (**Figure 3**). Note that this unknown component reflects

270   what was gained since a previous time point. In other words, at 12 months, the infant on average

271   acquired the same fraction of unknown as it did at 4 months and birth. When source tracking is

272   run without including previous time points as sources, the unknown component increases over

273   the first year of life for SNVs only (**Figure S5**).

274      Next, we sought to understand the effect of swapping sink and source in the re-analysis of

275   Backhed et al. 2015 data. In **Figure 3G and H**, the infant at birth is the potential source and

276   mother is the sink. The estimated contribution from baby to mother is significantly smaller

277   (species-FEAST: 11.9 difference,  Wilcoxon rank sum test p-value = 0.013; SNV-FEAST: 16.0

278   difference, p-value = 2.2 x 10[-5]) compared to that of mother to baby. This trend may be

11

279    suggestive, but is not conclusive, of directionality, whereby a less diverse source is seeded by a
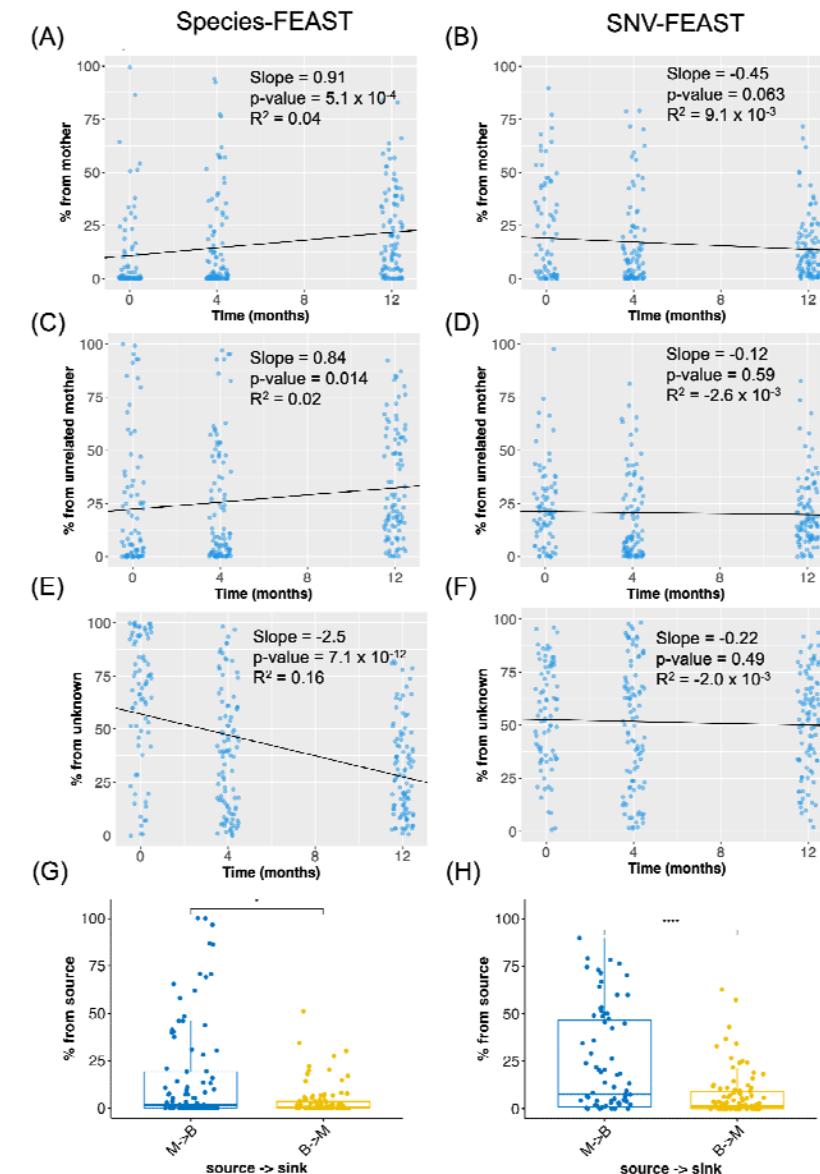
280    more diverse source.

281



282

283    **Figure 3. Source tracking in the infant gut microbiome over the first year of life.** Species-

284    and SNV-FEAST were applied to Backhed et al. 2019 data to estimate the contribution of (A, B)

285    mother, (C, D) unrelated mothers and (E, F) unknown sources to infants sampled at birth, four

286    months, and twelve months. The black line and inset statistics pertain to the linear regression fit

287    for the source estimates as a function of age of the infant. (G, H) are flipped source tracking

288    analyses with mother and infant swapped when using species-FEAST and SNV-FEAST,

12

289   respectively. **Figure S4** shows the species that were included in species-FEAST and species that

290   had SNVs included in SNV-FEAST. **Figure S5** shows the estimate of the unknown component

291   when previous time points of the infant are excluded from the sources.

292

293   **Contribution of the NICU built environment to infant microbiomes**

294   Next, we re-analyzed a metagenomic dataset studying the contribution of the hospital

295   environment to the infant gut microbiome in the neonatal intensive care unit (NICU) (Brooks et

296   al. 2017). This dataset is composed of microbiomes of infant stool, as well as the NICU rooms of

297   the same infants at frequently touched surfaces, sink basins, the floor, and isolette-top sampled

298   over an 11-month period (Brooks et al., 2017). We applied SNV and species-FEAST to assess

299   the contribution of the infant's own NICU room as well as a different NICU room in the vicinity

300   of the infant's gut microbiome (**see Methods**).

301   Concordant with the findings of Brooks et al., both SNV and species-FEAST detected

302   that the most common source contributing to the infant microbiome was the floor and isolette-top

303   from the infant's own room (**Figures 4A and B**). SNV-FEAST found Infant 18 also had large

304   contributions from their own room's touched surfaces at multiple time points (**Figure 4B**), which

305   is consistent with a finding by Brooks et al. that three strains found in Infant 18 perfectly

306   matched ($> 99.999\%$ average nucleotide identity) strains found in the touched surfaces samples

307   of Infant 18's own room. Lastly, both species-FEAST and SNV-FEAST found Infant 6's

308   microbiome was explained almost entirely by samples from a different room with SNV-FEAST

309   finding a sizeable contribution from both the floor and isolette top and the sink basin in this

310   different room. This is concordant with Brooks et al.'s finding of multiple cases of strain sharing

311   across rooms of Infant 6 and 12 for the different surfaces. FEAST with both data types is able to

312   quantify the extent to which Infant 6's microbiome was influenced by strains present in the built

313   environment.

314   Through application of SNV and species-FEAST, we are able to quantify any trends over

315   time the influence of the built environment on the infant microbiome (**Figures 4A and B**). SNV-

316   FEAST more consistently finds that contribution from the infant's own room exceeds

317   contributions from a different room over time (paired Wilcoxon signed rank test for same room >

318   different room: Infant 3: p-value = $1.95 \times 10^{-9}$, Infant 6: 1, Infant 12: $3.05 \times 10^{-5}$, Infant 18: 3.81

319   x $10^{-6}$) as compared to species-FEAST (Infant 3: p-value = 0.41, Infant 6: 1, Infant 12: $5.8 \times 10^{-4}$,

13

320    Infant 18: 3.81 x $10^{-6}$). Interestingly, species-FEAST assigns one dominant source primarily,

321    whereas SNV-FEAST more often finds a combination of sources for a given sample.

322         Additionally, both SNV and species-FEAST estimated a large unknown component for

323    all four infants, with Infant 18 showing the largest mean unknown component across the NICU

324    stay based on SNVs (**Figure S6**). This unknown component is important because it signifies the

325    extent to which other sources such as the mother and diet impact infant gut colonization.

326         We then asked the question: is the infant more explained by the built environment rather

327    than vice-versa, the built environment is more explained by the infant. We tested this by

328    swapping the infant and each of the three built environment sources (**Figure 4C and D**). The

329    estimated contribution of room to infant is significantly higher than the estimated contribution of

330    infant to room, but this asymmetry is more pronounced with SNV-FEAST. SNV-FEAST showed

331    significantly higher contribution of room to infant for two of the three surface types (floor and

332    isolette top: Wilcoxon rank sum test p-value = 7.00x $10^{-9}$, touched surface: p-value = 0.0058,

333    sink basin: p-value = 0.274) while species-FEAST found this to be true for one of the three

334    surface types (floor and isolette top: Wilcoxon rank sum test p-value = 7.1x $10^{-5}$, touched

335    surface: p-value = 0.968, sink basin: p-value =  0.998). Interestingly, the built environments of

336    different rooms highly resemble each other. This is especially apparent with species-FEAST,

337    suggestive of similar ecological forces operating in similar built environments. By contrast,

338    SNV-FEAST reveals a higher diversity of contributing sources of the built environment samples

339    to other NICU built environments, once again highlighting the utility of performing source

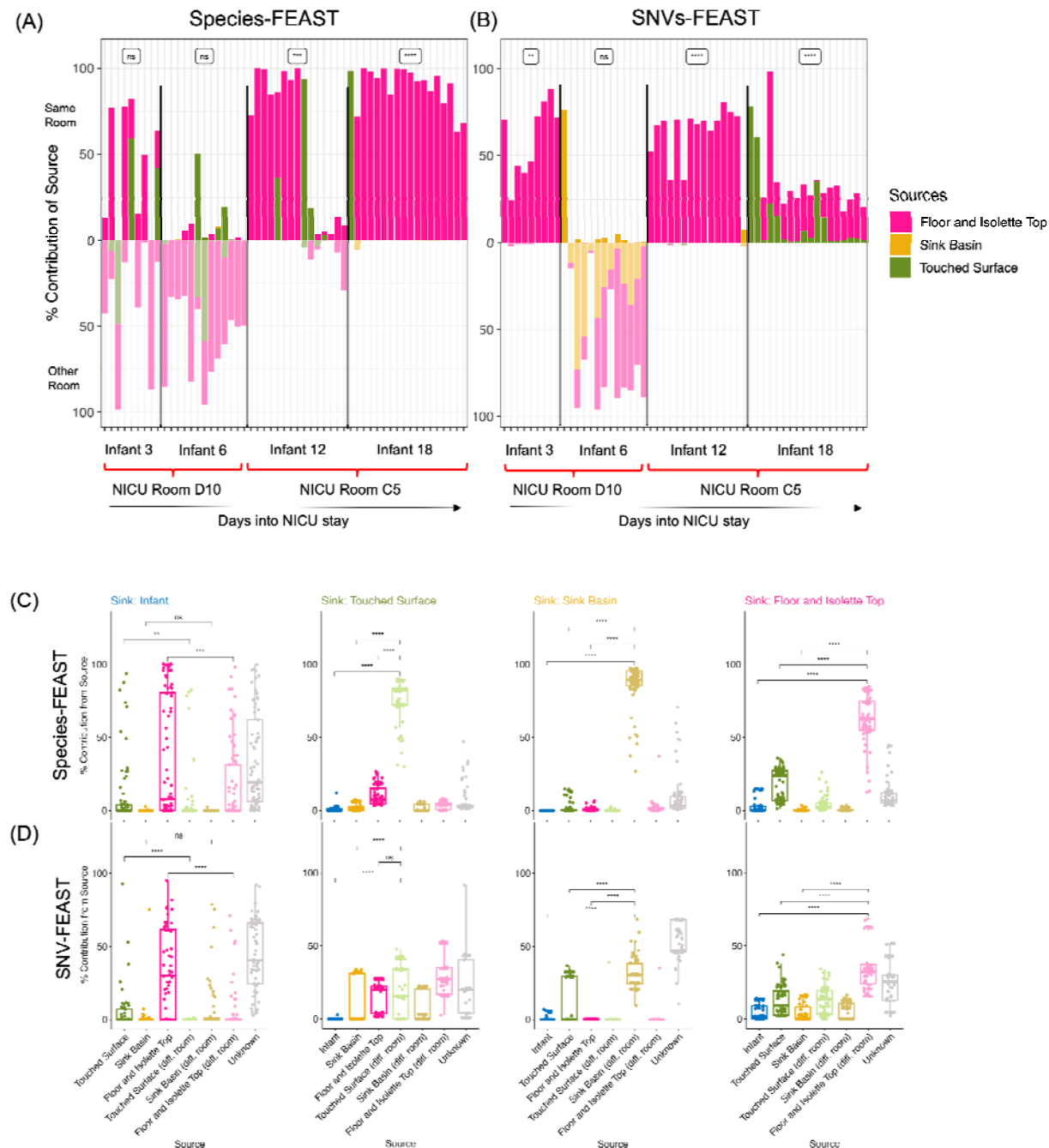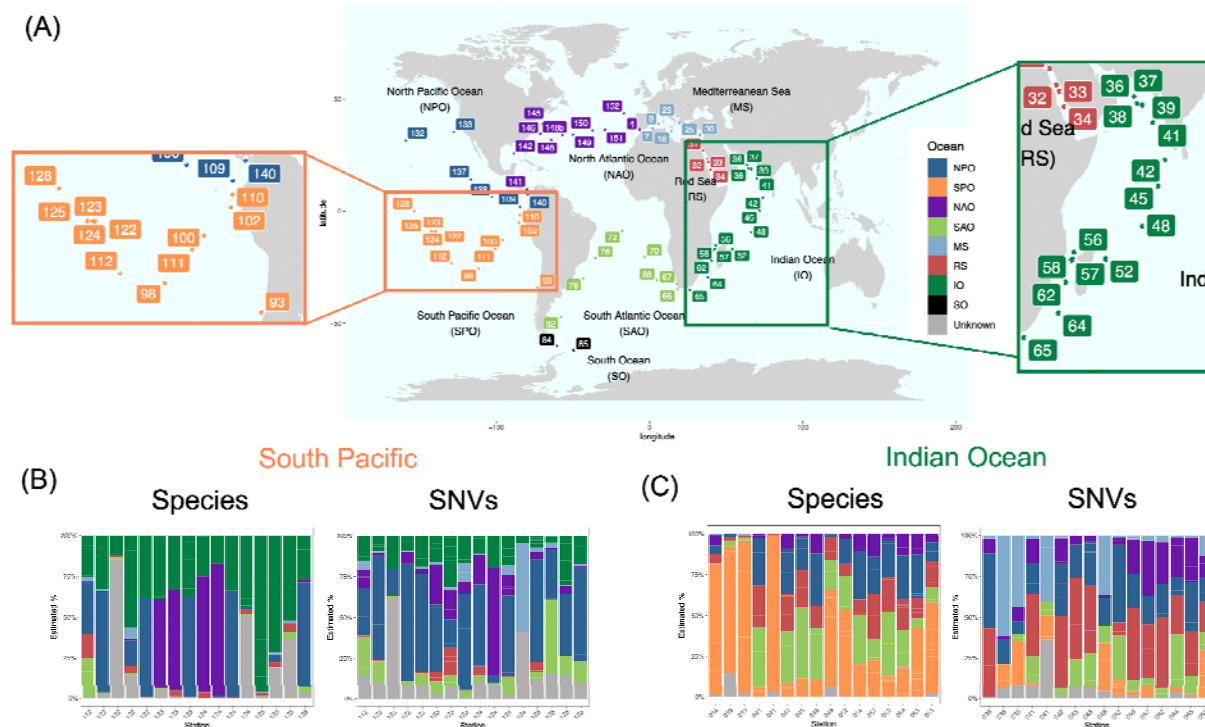340    tracking with SNVs.

341

342

**Figure 4: Source tracking of infant gut microbiome in the NICU.** (A) species-FEAST and (B) SNV-FEAST applied to infants in the NICU. Each bar represents one sampling day in the NICU stay of an infant. Infants 3 and 6 stayed in the same room, but at different times. The same applies to Infants 12 and 18. The contribution of a different room was determined by using samples from Infant 12's room for Infants 3 and 6, and samples from Infants 6's room for Infants 12 and 18 for each of the categories of surfaces per infant: touched surface, sink basin, or floor

15

349    and isolette top surface. The asterisks represent the result of a paired Wilcoxon signed rank test

350    indicating whether the total contribution of surfaces from the infant's own room were higher than

351    contributions from the other room: **** for p-value < 0.0001, *** for p-value < 0.001, ** for p-

352    value < 0.001, * for p-value < 0.05, and n.s. for p-value > 0.05. Iterative swapping of the infant

353    sink and each potential source for source tracking with (C) species-FEAST and (D) SNV-

354    FEAST. The first column shows source tracking results in which the infant was treated as the

355    sink. In each column after the first column, a different environmental source was swapped with

356    the infant and considered as a sink.

357

**Global source tracking of ocean microbiomes**

359         The ocean microbiome is a complex community that displays biogeography at the species

360    and functional levels (Nayfach et al., 2016; Sunagawa et al., 2015). To further understand global

361    patterns of ocean microbiomes, we applied SNV and species-FEAST to the Tara Oceans

362    microbiome dataset (Sunagawa et al., 2015). In the source tracking context, rather than defining

363    sharing as evidence of a transmission event (which is more likely in mother-infant data),

364    estimated source contributions at best explain the extent to which a given ocean sample

365    resembles other ocean samples. On one extreme, an ocean sample might be entirely explainable

366    by a single ocean's samples, and at the other extreme, an ocean sample might be explainable by

367    multiple oceans at the same time. Another alternative is for an ocean sample to not be

368    explainable by any of the provided sources, resulting in a high unknown component and

369    potentially suggesting high endemism. These source tracking estimates could be indicative of the

370    extent to which oceans mix or may be reflective of similar niches.

371         Tara Oceans is composed of 182 whole metagenomic sequencing samples derived from

372    64 stations at multiple depths. Previous research indicates that temperature is one of the highest

373    drivers of variability in microbial composition in the ocean (Ladau et al., 2013; Sunagawa et al.,

374    2015). For this reason, we restricted the source tracking analysis to sinks and sources from the

375    same temperature and depth range: above 20 degrees Celsius and within an average of 5 meters

376    below the surface.

377

**Figure 5. Microbial source tracking in the Tara Oceans dataset with SNV and species-FEAST.** World map indicating the location of sampling sites (A). Source tracking estimates for the contribution of different oceans to the South Pacific (n=16) (B) and Indian Oceans (n=16) (C) are depicted with vertical bars. In each experiment, all stations around the world excluding those from the "sink" ocean are considered potential sources. Light blue, for example, represents the total contribution of the four stations from the Mediterranean Sea that had samples in the surface layer that were also greater than 20°C in temperature.
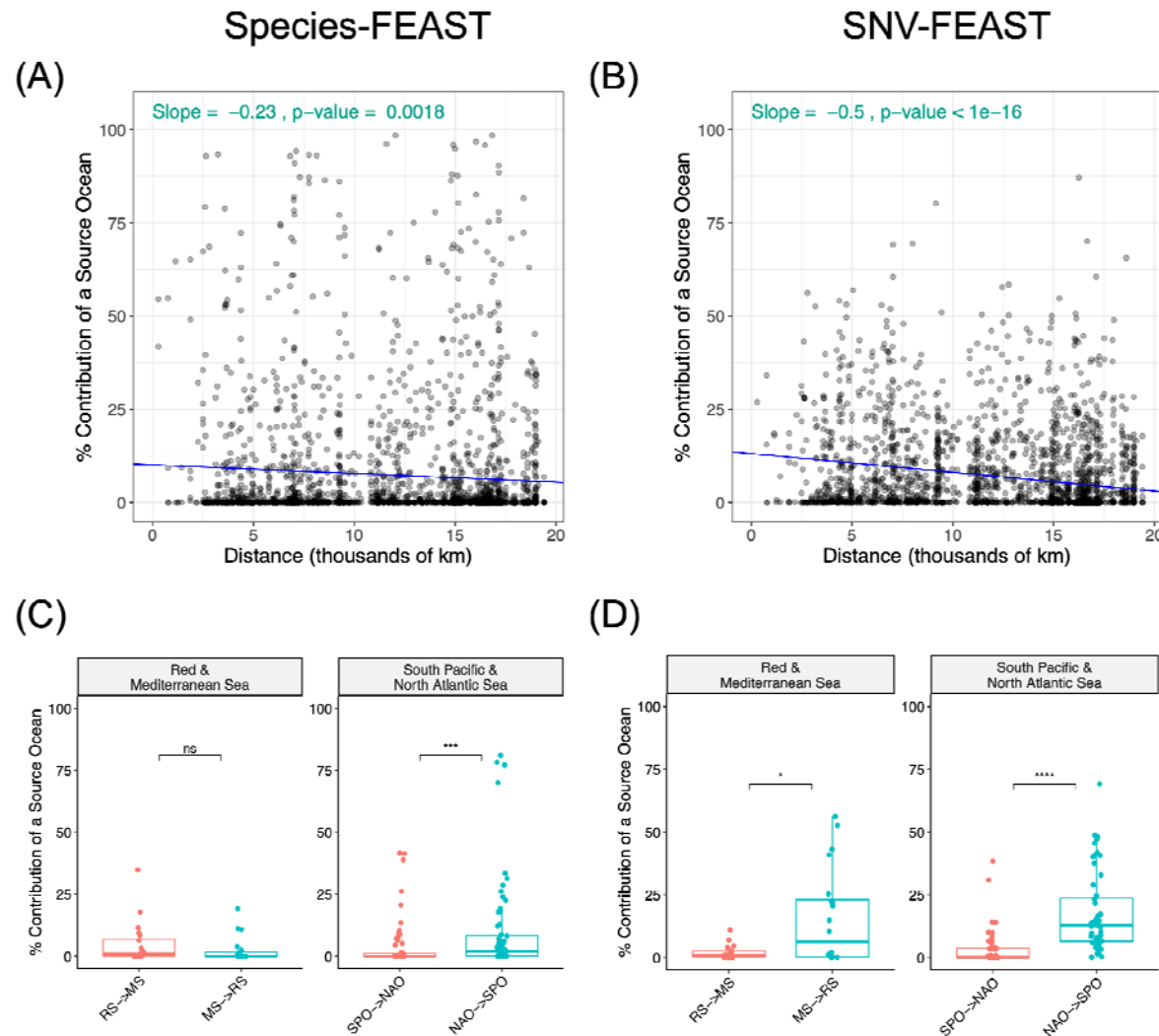
17

**Figure 6. Source tracking with ocean samples**.  Distance decay in contribution of a "source" ocean to a "sink" ocean when using (A) species-FEAST and (B) SNV-FEAST. In each experiment, only stations from one ocean were considered as sources for a given sink station. For example, when performing source tracking between the mediterranean and north atlantic, for each mediterranean station, the 10 available north atlantic stations were considered as potential sources. Thus, plotted are 10 points for a given mediterranean sink, where each point represents the contribution of a source station from the North Atlantic to the Mediterranean sink station in question. Shown in inset text are the slope and t-test p-value for the slope. (C) and (D) are flipped source tracking analysis with the Red Sea and Mediterranean, as well as the South Pacific Ocean and North Atlantic Ocean using species-FEAST and SNV-FEAST, respectively.

18

399     First, we performed source tracking between oceans using SNV and species-FEAST. We

400     treated each station around the world as a sink and estimated the contribution of different oceans

401     around the world to that sink (**Methods**). Unknown represents any portion of the microbiome in

402     these sink samples that cannot be explained by any of the provided source samples. We found

403     that species and SNV-FEAST estimate different amounts of sharing between oceans, where

404     SNVs estimate a higher unknown on average, potentially indicative of endemism. The finding

405     that SNV-FEAST estimates a higher unknown contribution on average is most evident in the

406     North Pacific, North Atlantic, South Atlantic, and Mediterranean oceans (**Figure S7**).

407     Additionally, in some oceans, SNVs identify contributions from oceans that species-FEAST does

408     not detect (**Figure 5, Figure S7**).  For example, in applying FEAST to Indian Ocean samples we

409     find that there is measurable sharing of microbes with the Mediterranean Sea, but this is not

410     detected with species (**Figure 5C**). Such differences were found in samples from other oceans as

411     well (**Figure S7**).

412     Next, we assessed whether source tracking estimates display a distance-decay

413     relationship. Previous studies found that genetic distance, such as that represented by fixation

414     index $F_{ST}$, increases with geographic distance between populations (Cavalli-Sforza & Feldman,

415     2003; DeGiorgio & Rosenberg, 2013). Based on these findings, our expectation was that samples

416     that are further away from a given station will have reduced resemblance to that station. To

417     assess this distance-decay relationship, we plotted pairwiseh source tracking results across all

418     possible pairs of ocean samples (**Figure 6A and B**). We found that indeed as the distance

419     increases, the % explainability of a given source ocean decreases -0.23 % per thousand km

420     according to species-FEAST (t-test p-value $< 1 \times 10^{-16}$), and -0.5% per thousand km according to

421     SNV-FEAST (t-test p-value = 0.0018 ). The steeper slope for SNV-FEAST suggests that SNVs

422     may be more sensitive to distance decay signals on a global level.

423     Finally, we investigated whether some oceans have higher estimated contributions to

424     other oceans than vice versa, potentially indicative of the directionality of transmissions (though

425     see Discussion).  Specifically, we investigated the relationship between the Red Sea to the

426     Mediterranean Sea (**Figure 6C and D**). Migration from the Red Sea to the Mediterranean,

427     known as Lessepsian migration, is well-documented for not only microorganisms but also

428     macroorganisms like fish (Bentur et al., 2008; Bianchi & Morri, 2003; Golani, 2009). Anti-

429     Lessepsian migration (Red Sea to Mediterranean), on the other hand, has been primarily thought

19

430  to be rare due to the Additionally, recent studies suggest that there is also evidence for anti-

431  Lessepsian migration of bacteria (Mediterranean to Red Sea) may be more common than

432  Lessepsian migration (Elsaeed et al., 2021). Research studies find that Mediterranean has brine

433  pools that produce similar a similar environment to the Red Sea's (Antunes et al., 2011),

434  allowing for bacteria from the MS to potentially thrive in the RS.

435      By swapping the Red Sea and Mediterranean as source and sink, we found that there was

436  indeed a significant difference in the estimated contribution from one direction to another with

437  SNVs but not species (**Figure 6C and D**). SNV-FEAST found the Mediterranean explained an

438  average of 15% of the Red Sea, while the Red Sea explained an average of 1.8% of the

439  Mediterranean (Wilcoxon rank sum test, p-value =0.02), consistent with anti-Lessepsian

440  migration.  Meanwhile, a similar analysis with species-FEAST found the Mediterranean

441  explained 2.5% of the Red Sea and the Red Sea explained 4.9% of the Mediterranean (Wilcoxon

442  rank sum test, p-value = 0.25). In a similar analysis between North Atlantic and South Pacific we

443  found that both species and SNVs supported significantly greater contributions from the North

444  Atlantic to the South Pacific, with SNV-FEAST estimating a greater contribution (17%,

445  Wilcoxon rank sum test p-value = $5.1 \times 10^{-11}$) than species-FEAST (10%, Wilcoxon rank sum

446  test p-value =$1.8 \times 10^{-4}$).

447      Together, these results suggest that on average, SNV and species FEAST generate similar

448  source tracking results in the Tara Oceans dataset, with SNVs displaying stronger signals of

449  endemism, distance-decay relationships, and potential directionality of transmission.

450

451

452  **DISCUSSION**

453      Source tracking provides insight into potential source contributions to a metagenomic

454  sample as well as similarities between metagenomic samples. While species abundances have

455  been informative in source tracking in several studies (Flores et al., 2011; Knights et al., 2011;

456  McGhee et al., 2020; Shenhav et al., 2019), they may be limited in their resolution. SNVs

457  provide a potential alternative because of their ability to distinguish sources of strain

458  transmissions. Here we compared the ability of a previously published source tracking algorithm

459  FEAST using species versus SNVs as input data. In application of species and SNV-FEAST to

460  simulations as well as three case studies, we demonstrate that the two input types can provide

461    distinct insights into microbial sharing and similarities across different environments. As a

462    hypothetical example, two unrelated samples may have very similar species composition due to

463    similar colonization processes and similar environmental influences without any actual microbial

464    sharing. It would be unlikely for these two unrelated samples to share rare SNVs, however. This

465    distinction suggests that SNVs indeed can provide insight into the ecological processes shaping

466    microbial communities that species information alone cannot, and our three case studies are able

467    to demonstrate this.

468         In the first case study, we confirmed previous findings that SNV sharing between

469    mothers and infants decreases over the first year of life while species sharing increases (Nayfach

470    et al., 2016), suggesting that while the infant microbiome matures to resemble adults at the

471    species level, sources other than the mother may seed the infant over time. In the second case

472    study, we confirmed source contributions from the NICU built environment to the infant

473    microbiome (Brooks et al., 2017), and found that SNVs detect a more consistent estimate in

474    source contributions overtime compared to species as well as detecting contribution from sources

475    not detectable by species-FEAST.

476         In the third case study, we perform source tracking in the Tara oceans dataset and found

477    SNVs display a stronger distance decay relationship. These distance-decay results parallel recent

478    findings made with gene content (Dlugosch et al., 2022). While previous studies have examined

479    the biogeography of the ocean using species profiles, genes (Dlugosch et al., 2022; Nayfach et

480    al., 2016) or amino acid variants from a single species (SAR11) (Delmont et al., 2019), for the

481    first time, we leverage the use of SNVs across all detected prevalent species in the ocean

482    microbiome to identify proportions of sharing across oceans. A benefit of using SNVs in the

483    ocean microbiome is that SNVs can track fragments of DNA that have moved due to horizontal

484    gene transfer in the distant past rather than relying on inference of whole genomes or presence of

485    private SNVs that may been transmitted in the recent past. This global-level source tracking is

486    analogous to admixture estimation with human genotypes (Alexander et al., 2009; Chiu et al.,

487    2022).

488         We note that source tracking provides insights into similarities between microbiomes and

489    potential transmissions, though the directionality is less conclusive. It is possible that increased

490    contributions in one direction but not the other is suggestive of directionality of transmission. For

491    example, in the case of the mother-infant data from Backhed et al. 2015, FEAST predicted

492   higher contribution from mother to baby than vice versa. This is consistent with work done on

493   crAss-like phage transmissions between mother and infant in the same dataset that showed

494   evidence of directionality by tracking the accumulation of mutations over time that are private to

495   the infant and absent from the mother (Siranosian et al., 2020).  But in the case of the ocean, it is

496   possible that over longer time periods, differences in relative contributions from one part of the

497   world to another (e.g. Mediterranean to Red Sea) are more reflective of local selection pressures

498   that permit certain species and genotypes (Delmont et al., 2019). Thus, source tracking in certain

499   instances, such as the ocean microbiome, at best reflects the extent of similarity between samples

500   and is less conclusive about directionality.

501       A popular approach used to track strain transmissions is by detecting high average

502   nucleotide identity (ANI) for species shared between source and sink. For example, inStrain

503   (Olm et al., 2021) identifies a match between samples for a given species when ANI exceeds

504   99.999%. However, it is to be noted that inStrain provides distinct and complementary

505   information from FEAST given its binarization of whether or not a strain is shared. For

506   illustration purposes, if an infant harbors 100 species, of which only 1 came from their mother,

507   but that species' strain's relative abundance is 50% of the infant's microbiome, SNV-FEAST

508   would infer that the mother's contribution is 50%, while inStrain would infer that only $1/100^{th}$ of

509   the species are derived from the mother.

510       Other studies rely on tracking transmissions of strains with private SNVs shared only

511   between the sink and putative source (Bäckhed et al., 2015; Korpela et al., 2018; Nayfach et al.,

512   2016; Schmidt et al., 2019). The private marker allele tracking approach in Nayfach et al. 2016

513   provides an improved estimate of true percentage of species that share some portion of their

514   genome with putative sources compared to inStrain (**Figure S2, S3**). It is possible that requiring

515   only 5% of marker alleles to be shared rather than a 99.999% ANI permits detection of

516   horizontal gene transfers between lineages residing in mothers and infants (D. W. Chen &

517   Garud, 2022; Vatanen et al., 2022). However, in FEAST, by using any SNV with an informative

518   distribution across sources as determined by our signature scoring method, we are able to

519   quantify the relative contribution of all the sampled environments and assign a proportion to

520   these putative sources. Another advantage to FEAST is that the contribution of unknown sources

521   can be quantified. For example, the significant fraction of marine biodiversity estimated to be

22

522    'unknown' may be endemic, as previously noted in the Mediterranean (Katsanevakis et al.,

523    2014).

524        A drawback, however, with using SNVs over species is deeper, whole genome

525    sequencing is required to accurately call SNVs. Moreover, even when there is sufficient

526    coverage, there is still the challenge of a large number of SNVs. We demonstrate one way to

527    subset SNVs that uses a scoring method for informativeness, but there may yet be other methods

528    for filtering SNVs to the most informative set. Another potential caveat of SNV filtering is that

529    not all species present will be represented in the final signature SNV set (**Figure S4**).  Species

530    with higher abundance are more likely to be represented in the signature SNV set.  However, we

531    show that not all species need to contribute signature SNVs in order to make accurate inferences,

532    and likewise, not all SNVs are needed to make accurate inferences (**Figure S1**).

533        Ascertainment of SNVs from metagenomic data in a high-throughput manner, especially

534    common SNVs with microbiome genotyping technology (Shi et al., 2021), is becoming an

535    increasing priority for the field as metagenomic datasets become more abundant. A genotyper for

536    prokaryotes has already been developed and tested on a catalog of over 100 million SNVs in

537    order to characterize population structure (Shi et al., 2021). Such a catalog of informative SNVs

538    could be invaluable for source tracking. With source tracking enabling us to characterize samples

539    by their relationship to known samples, we have a powerful tool to explore samples in new

540    contexts we have yet to discover.

541

542    **METHODS**

543    *Data*

544        For simulations and analyses of infant microbiomes in the first year of life (Bäckhed et

545    al., 2015), we downloaded the raw shotgun metagenomic sequencing reads from public read

546    archives under accession number PRJEB6456. We downloaded the raw sequence reads for the

547    NICU analysis (Brooks et al., 2017) from accession number PRJEB323631, and the equivalent

548    for the Tara Oceans analyses (Sunagawa et al., 2015) were downloaded from accession number

549    PRJEB402. Data from the HMP Consortium (Methé et al., 2012) and Lloyd-Price et al (Lloyd-

550    Price et al., 2017) was downloaded from the following

551    URL: https://aws.amazon.com/datasets/human-microbiome-project/.

552

23

*Estimation of species and SNV content of metagenomic samples*

We used MIDAS (Metagenomic Intra-Species Diversity Analysis System, version 1.2, downloaded on November 21, 2016 (Nayfach et al., 2016) to estimate species abundance and SNV content per species in each metagenomic shotgun sequencing sample. The database we used to apply MIDAS consisted of 31,007 bacterial genomes that are clustered into 5,952 species. The parameters we used to estimate species abundances and SNVs were described in (Garud et al., 2019). A species was considered present if there are at least 3 reads mapping to a set of single copy marker genes on average. To call SNVs, we used the default MIDAS settings in order to map reads to a single representative reference genome. The mapping was done with Bowtie 2 (Langmead & Salzberg, 2012): global alignment, MAPID$\geq$94.0%, READQ$\geq$20, ALN_COV$\geq$0.75, and MAPQ$\geq$20, where species with reads mapped to less than 40% of the genome were excluded from the SNV calls. We excluded samples with depth lower than 5 reads, and excluded genetic sites using the default site filters of MIDAS (e.g. ALLELE_FREQ$\geq$0.01, with the exception of SITE_DEPTH which was set to 3.

*Application of FEAST algorithm*

FEAST, originally introduced by Shenhav et al., is an R-based method that models the mixture proportions for various "source" microbial samples for a given "sink" (Shenhav et al., 2019). This method utilizes expectation maximization to estimate the proportions when given any sort of count-based feature matrix representing the potential sources and sinks. The intuition behind the estimation process is that a source with a similar species distribution to the sink would have a higher contribution estimate to the sink. A species with non-zero counts only in source *j* and the sink would increase the estimated contribution of source *j*. However, in many cases, the same species are found in multiple sources simultaneous. The algorithm does not uniquely assign a species to a source but rather simultaneously utilizes all species information to infer the source contributions. The method was originally tested and evaluated on species and not on more fine scale genetic data such as SNVs. The number of different species, on average, range in number from a few hundred to a few thousand, while the number of possible nucleotide sites that vary across different sources can number in millions. For this reason, a SNV-filtering process is necessary so that the algorithm can run within a reasonable time and with reasonable memory requirements.

*Application of FEAST to the Backhed et al. 2015 dataset:*

585    For both species and SNV-FEAST, the same set of sources and sinks were fed into the

586    FEAST algorithm. In the case study of infants in the first year of life (Bäckhed et al., 2015), the

587    sink consisted of the infant fecal sample at either four days, four months, or 12 months and the

588    potential sources consisted of fecal samples from the true mother, three randomly selected

589    mothers from the same dataset, and also any previous time points for the infant.

590    Species-FEAST utilized all species present in the infant whereas SNV-FEAST used

591    signature SNVs from the subset of species that had signature SNVs. Shown in **Figure S3** are the

592    distribution of species included in species and SNV-FEAST.

593

594    *Application of FEAST to the Brooks et al. 2017 dataset:*

595    For the case study of infants in the NICU (Brooks et al., 2017),  the sink consisted of the

596    fecal sample of the infant at a given time point and the potential sources consisted of pooled

597    reads from the touched surfaces, the sink basin and the floor and isolette top from both the

598    infant's own room as well as a different room. The different room was Infant 12's room for

599    Infants 3 and 6, Infants 6's room for Infants 12 and 18.

600

601    *Application of FEAST to the Sunagawa et al. 2015 dataset: Determining the signature SNV set*

602    Signature SNVs were identified as described in the main text. We provide specific steps for

603    determining signature SNVs:

604    (1) Filter sites: only sites of the genome with at least the required number of reads mapping

605        to the site are considered. In the case study of infants in the first year of life (Bäckhed et

606        al., 2015) and infants in the NICU (Brooks et al., 2017), the minimum coverage

607        requirement is 10 across the sink and $J$ sources. For the Tara Ocean (Sunagawa et al.,

608        2015) samples, the minimum coverage is five reads (Sunagawa et al., 2015).

609        Additionally, sites that are biallelic must have more than one read mapped to each allele

610        to be considered.

611    (2) Perform per site per source parameter estimates: for each potential source compute the

612        estimated allele frequency in the sink θ under two different hypotheses:

613        **Hypothesis 1:** Source $i$ with allele frequency $p_i$ explains the allele counts in the sink.

$$\hat{\theta} = p_i$$

614     **Hypothesis 2:** A combination of all other sources except $i$ (sources $j \neq i$) explain the

615     observed allele count distribution in the sink. The estimate of the sink allele frequency is

616     computed using a mixture of the allele frequencies $p_j$ from those sources. The mixing

617     parameter $\alpha_j$ is learned using Sequential Least Squares Programming (scipy.minimize() )

618     with the constraint of summing to 1 with bounds of 0 to 1 inclusive: $\sum_{j \neq i} \alpha_j = 1$.

619

$$\hat{\theta} = \sum_{j \neq i} \alpha_j p_i$$

620     **(3)** Compute per site per source log likelihoods: Compute the binomial log-likelihood under

621     hypotheses 1 and 2, given $n$ reads with the reference allele and $m$ reads with the

622     alternative allele in the sink:

623

$$l(\hat{\theta}) = n \, log \, \hat{\theta} + m \, log(1 - \hat{\theta})$$

624     (4) Compute per site per source log likelihood ratio:

$$l_1(\theta) - l_2(\theta)$$

625

626     (5) Compute per site summary signature score: The maximum log likelihood ratio per site is

627     the signature score for that SNV, representing how favorably one of the sources explains

628     the sink over all other sources

629     (6)  Filtering of SNVs using signature score: One signature score for that SNV represents

630     how favorably one source explains the sink better than all other sources. All the scores

631     are ranked across SNVs and SNVs with scores that are greater than two standard

632     deviations over the mean signature score within each 200 kbp window of the genome are

633     retained as signature SNVs. This window size was chosen for to optimize run time and

634     memory requirements.

635

636     Note, if only one source passes minimum coverage filtering, $l_2(\theta) = 0$  resulting in a

637     very high likelihood ratio as represented by $l_1(\theta)$ for the one source. These SNVs are

638     more likely to pass the signature score filtering. One exception for SNVs that are

639     included in the signature SNV set without passing signature score filtering are SNVs with

640     an allele that is completely unique to the infant, as these represent SNVs that are

26

641    potentially derived from an unknown source. Signature SNVs are obtained from the SNV

642    profile of every species for which there is MIDAS output.

643

644    *Simulating mother to infant transmission*

645    The mixture proportions for 28 simulated infants is shown in **Table S1**. Four possible

646    scenarios are simulated using a combination of either low or high number of sources and low or

647    high transmission probabilities of species. High transmission of species was simulated by

648    drawing separate transmission probabilities for each species in each contributing source based on

649    a beta distribution with a mean equal to the species relative abundance and variance equal to 0.1,

650    a value selected to emulate Backhed et al.'s mean relative abundance and variance. For the low

651    transmission scenario, transmission probabilities were drawn from a beta distribution with mean

652    0.1 times the relative abundance and variance at 0.1. To determine if a species from each source

653    was transmitted to a given infant, a binomial draw was performed $J$ times, where $J$ = number of

654    sources, and the probability of a mother transmitting the species is $p_j$ based on the beta-drawn

655    transmission probability. If any of the draws yields a one, that species is transmitted to the infant

656    from all sources. The same simulated data under these scenarios is used for both SNV and

657    species source tracking.

658    The source tracking estimates are compared to the true mixing proportions using

659    Spearman correlation. The significance of correlation is calculated using the stat_cor function in

660    the 'ggpubr' package (*CRAN - Package Ggpubr*, n.d.).

661

662    *Comparison to inStrain*

663    We ran inStrain (Olm et al., 2021) on the same synthetic samples as described above.

664    InStrain "profile" (Olm et al., 2021)  and inStrain "compare" (Olm et al., 2021) were run for

665    every possible infant-source pair. For example, for simulated infant 1 there were 10 putative

666    sources, therefore inStrain compare was run 10 times for each putative source. InStrain reports

667    popANI calculated per scaffold for a given species. To compute a single statistic per species, we

668    computed the average popANI across scaffolds for a given species. The percent infant

669    microbiome species that had strains shared with mother was computed as the number of species

670    in which popANI was >= 99.999% divided by the total number of species with coverage >= 5.

671    PopANI was only calculated in scaffolds that had >=5 coverage in both samples of the pair.

27

672

*Comparison with strain tracking approach in Nayfach et al. 2016*

We applied the strain tracking approach in Nayfach et al. 2016 on the same synthetic communities described above. In Nayfach et al. 2016, strain transmissions are tracked by identifying 'marker alleles' which are private to the infant, mother, or infant-mother dyad, and absent from the broader population. A strain is considered to be shared if at least 5% of all marker alleles for a mother-infant dyad are shared. Note that the approach for strain tracking proposed in Nayfach et al. 2016 utilizes SNV information outputted by MIDAS, but is not a part of MIDAS.

Each simulated infant had up to 10 sources that were real maternal samples from Backhed et al. 2015 For each possible pair of infants and maternal sources (10 pairings per infant, with 48 infants), we found the set of infant-only marker alleles, mother-only marker alleles, and mother-infant dyad marker aleles. As described in Nayfach et al, 2016, only sites with minimum 30 reads and only alleles that were supported by at least 10% of the total reads aligned to that site were considered. The infant marker allele and mother marker allele were defined as alleles that were present only in the focal sample and absent from the background samples (or below 3 reads = 10% * 30 reads). For the infant, the background consisted of all mothers (including mothers that were used to simulate the infant), real infant samples (excluding infants of mothers used to simulate the infant), and 337 samples of adults from the United States in the HMP (which includes 180 unique adults) that were obtained from the metagenomics repository of HMP under project ID SRP002163 and SRP056641 (Lloyd-Price et al., 2017; Methé et al., 2012). For the mother, the background consisted of all mother and infant samples in addition to the HMP samples. For computing shared marker alleles, an allele must be present in both the mother and infant but absent from the background, which consisted of all mothers and the HMP samples.

To compute sharing, two quantities were considered: "total sharing", defined as % shared marker alleles/ (infant marker alleles + mother marker alleles + shared marker alleles) and proportion of infant marker alleles that are shared: % shared marker alleles/ (infant marker alleles + shared marker alleles). The first quantity compared to FEAST estimates was the percentage of infant species in which the "total sharing" was at least 5%. The second quantity

28

702 compared to FEAST was the pooled proportion of infant marker alleles that are shared across all

703 species.

704

705 *Distance Decay Analysis*

706 To study the relationship between source tracking estimates and geographic distance, we

707 analyzed all oceans as either a sink or source against all other possible oceans. To compute

708 geographic distance between stations, we applied the Haversine distance to the longitude and

709 latitude of the sampling sites provided by (Sunagawa et al., 2015) using the package "geosphere"

710 (Hijmans et al., 2021). Source tracking estimates were computed as described above using either

711 SNV-FEAST or Species FEAST. The regression line for the distance decay analysis was

712 computed using a linear mixed model "contribution ~ distance + (1| sink_ocean)".

713

**Ethics declarations**

714

715 The authors declare that they have no competing interests.

716

724

**Availability of Data and Materials**

725

726 All metagenomic data was obtained from public repositories. The applicable accessions numbers

727 are PRJEB6456 for Backhed et al. 2015 (mother-infant), PRJEB323631 for Brooks et al. 2017

728 (NICU), PRJEB402 for Sunagawa et al. 2015 (Tara Oceans), and SRP002163 and SRP056641

729 for HMP.

730

731 Source code for SNV-FEAST signature SNV selection as well as analyses in this paper are

732 available at GitHub (https://github.com/garudlab/Signature-SNVs), Zenodo (DOI

733    10.5281/zenodo.7515044), and PyPi for pip installation (https://pypi.org/project/Signature-

734    SNVs/0.0.1/).

735

736

737    **References**

738    Alexander, D. H., Novembre, J., & Lange, K. (2009). Fast model-based estimation of ancestry in

739        unrelated individuals. *Genome Research*, *19*(9), 1655–1664.

740        https://doi.org/10.1101/GR.094052.109

741    Antunes, A., Ngugi, D. K., & Stingl, U. (2011). Microbiology of the Red Sea (and other) deep-

742        sea anoxic brine lakes. *Environmental Microbiology Reports*, *3*(4), 416–433.

743        https://doi.org/10.1111/J.1758-2229.2011.00264.X

744    Asnicar, F., Manara, S., Zolfo, M., Truong, D. T., Scholz, M., Armanini, F., Ferretti, P., Gorfer,

745        V., Pedrotti, A., Tett, A., & Segata, N. (2017a). Studying Vertical Microbiome

746        Transmission from Mothers to Infants by Strain-Level Metagenomic Profiling. *MSystems*,

747        *2*(1). https://doi.org/10.1128/msystems.00164-16

748    Asnicar, F., Manara, S., Zolfo, M., Truong, D. T., Scholz, M., Armanini, F., Ferretti, P., Gorfer,

749        V., Pedrotti, A., Tett, A., & Segata, N. (2017b). Studying Vertical Microbiome

750        Transmission from Mothers to Infants by Strain-Level Metagenomic Profiling. *MSystems*,

751        *2*(1). https://doi.org/10.1128/MSYSTEMS.00164-16/ASSET/54C4C531-C6DB-421B-

752        8C8A-10C0ECFE3BE9/ASSETS/GRAPHIC/SYS0011720800004.JPEG

753    Bäckhed, F., Roswall, J., Peng, Y., Feng, Q., Jia, H., Kovatcheva-Datchary, P., Li, Y., Xia, Y.,

754        Xie, H., Zhong, H., Khan, M. T., Zhang, J., Li, J., Xiao, L., Al-Aama, J., Zhang, D., Lee, Y.

755        S., Kotowska, D., Colding, C., … Jun, W. (2015). Dynamics and stabilization of the human

756        gut microbiome during the first year of life. *Cell Host and Microbe*, *17*(5), 690–703.

757        https://doi.org/10.1016/j.chom.2015.04.004

758    Bentur, Y., Ashkar, J., Lurie, Y., Levy, Y., Azzam, Z. S., Litmanovich, M., Golik, M., Gurevych,

759        B., Golani, D., & Eisenman, A. (2008). Lessepsian migration and tetrodotoxin poisoning

760        due to Lagocephalus sceleratus in the eastern Mediterranean. *Toxicon*, *52*(8), 964–968.

761        https://doi.org/10.1016/J.TOXICON.2008.10.001

762    Bianchi, C. N., & Morri, C. (2003). Global sea warming and "tropicalization" of the

763        Mediterranean Sea: biogeographic and ecological aspects. *Biogeographia – The Journal of*

764       *Integrative Biogeography*, *24*(1). https://doi.org/10.21426/B6110129

765   Brooks, B., Olm, M. R., Firek, B. A., Baker, R., Thomas, B. C., Morowitz, M. J., & Banfield, J.

766       F. (2017). Strain-resolved analysis of hospital rooms and infants reveals overlap between

767       the human and room microbiome. *Nature Communications*, *8*(1), 1–7.

768       https://doi.org/10.1038/s41467-017-02018-w

769   Cavalli-Sforza, L. L., & Feldman, M. W. (2003). The application of molecular genetic

770       approaches to the study of human evolution. *Nature Genetics 2003 33:3*, *33*(3), 266–275.

771       https://doi.org/10.1038/ng1113

772   Chen, D. W., & Garud, N. R. (2022). Rapid evolution and strain turnover in the infant gut

773       microbiome. *Genome Research*, *32*(6), 1124–1136.

774       https://doi.org/10.1101/GR.276306.121/-/DC1

775   Chen, E. Z., & Li, H. (2016). A two-part mixed-effects model for analyzing longitudinal

776       microbiome compositional data. *Bioinformatics*, *32*(17), 2611–2617.

777       https://doi.org/10.1093/BIOINFORMATICS/BTW308

778   Chiu, A. M., Molloy, E. K., Tan, Z., Talwalkar, A., & Sankararaman, S. (2022). Inferring

779       population structure in biobank-scale genomic data. *The American Journal of Human*

780       *Genetics*. https://doi.org/10.1016/J.AJHG.2022.02.015

781   *CRAN - Package ggpubr*. (n.d.). Retrieved March 6, 2022, from https://cran.r-

782       project.org/web/packages/ggpubr/index.html

783   DeGiorgio, M., & Rosenberg, N. A. (2013). Geographic Sampling Scheme as a Determinant of

784       the Major Axis of Genetic Variation in Principal Components Analysis. *Molecular Biology*

785       *and Evolution*, *30*(2), 480–488. https://doi.org/10.1093/MOLBEV/MSS233

786   Delmont, T. O., Kiefl, E., Kilinc, O., Esen, O. C., Uysal, I., Rappé, M. S., Giovannoni, S., &

787       Eren, A. M. (2019). Single-amino acid variants reveal evolutionary processes that shape the

788       biogeography of a global SAR11 subclade. *ELife*, *8*. https://doi.org/10.7554/ELIFE.46497

789   Dlugosch, L., Poehlein, A., Wemheuer, B., Pfeiffer, B., Badewien, T. H., Daniel, R., & Simon,

790       M. (2022). Significance of gene variants for the functional biogeography of the near-surface

791       Atlantic Ocean microbiome. *Nature Communications 2022 13:1*, *13*(1), 1–11.

792       https://doi.org/10.1038/s41467-022-28128-8

793   Elsaeed, E., Fahmy, N., Hanora, A., & Enany, S. (2021). Bacterial Taxa Migrating from the

794       Mediterranean Sea into the Red Sea Revealed a Higher Prevalence of Anti-Lessepsian

795       Migrations. *OMICS A Journal of Integrative Biology*, *25*(1), 60–71.

796       https://doi.org/10.1089/omi.2020.0140

797 Ferretti, P., Pasolli, E., Tett, A., Asnicar, F., Gorfer, V., Fedi, S., Armanini, F., Truong, D. T.,

798       Manara, S., Zolfo, M., Beghini, F., Bertorelli, R., De Sanctis, V., Bariletti, I., Canto, R.,

799       Clementi, R., Cologna, M., Crifò, T., Cusumano, G., … Segata, N. (2018). Mother-to-Infant

800       Microbial Transmission from Different Body Sites Shapes the Developing Infant Gut

801       Microbiome. *Cell Host and Microbe*, *24*(1), 133-145.e5.

802       https://doi.org/10.1016/j.chom.2018.06.005

803 Flores, G. E., Bates, S. T., Knights, D., Lauber, C. L., Stombaugh, J., Knight, R., & Fierer, N.

804       (2011). Microbial Biogeography of Public Restroom Surfaces. *PLoS ONE*, *6*(11), e28132.

805       https://doi.org/10.1371/journal.pone.0028132

806 Garud, N. R., Good, B. H., Hallatschek, O., & Pollard, K. S. (2019). Evolutionary dynamics of

807       bacteria in the gut microbiome within and across hosts. *PLoS Biology*, *17*(1), e3000102.

808       https://doi.org/10.1371/JOURNAL.PBIO.3000102

809 Golani, D. (2009). Distribution of Lessepsian migrant fish in the Mediterranean.

810       *Http://Dx.Doi.Org/10.1080/11250009809386801*, *65*(S1), 95–99.

811       https://doi.org/10.1080/11250009809386801

812 Hijmans, R. J., Karney, C., Geographiclib, ] (, Williams, E., Vennes, C., & Maintainer, ]. (2021).

813       *Package "geosphere."* https://doi.org/10.1007/s00190012

814 Katsanevakis, S., Coll, M., Piroddi, C., Steenbeek, J., Lasram, F. B. R., Zenetos, A., & Cardoso,

815       A. C. (2014). Invading the Mediterranean Sea: Biodiversity patterns shaped by human

816       activities. *Frontiers in Marine Science*, *1*(SEP), 32.

817       https://doi.org/10.3389/FMARS.2014.00032/ABSTRACT

818 Knights, D., Kuczynski, J., Charlson, E. S., Zaneveld, J., Mozer, M. C., Collman, R. G.,

819       Bushman, F. D., Knight, R., & Kelley, S. T. (2011). Bayesian community-wide culture-

820       independent microbial source tracking. *Nature Methods*, *8*(9), 761–765.

821       https://doi.org/10.1038/nmeth.1650

822 Korpela, K., Costea, P., Coelho, L. P., Kandels-Lewis, S., Willemsen, G., Boomsma, D. I.,

823       Segata, N., & Bork, P. (2018). Selective maternal seeding and environment shape the

824       human gut microbiome. *Genome Research*, *28*(4), 561–568.

825       https://doi.org/10.1101/GR.233940.117/-/DC1

Ladau, J., Sharpton, T. J., Finucane, M. M., Jospin, G., Kembel, S. W., O'Dwyer, J., Koeppel, A. F., Green, J. L., & Pollard, K. S. (2013). Global marine bacterial diversity peaks at high latitudes in winter. *The ISME Journal 2013 7:9*, *7*(9), 1669–1677. https://doi.org/10.1038/ismej.2013.37

Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods 2012 9:4*, *9*(4), 357–359. https://doi.org/10.1038/nmeth.1923

Li, S. S., Zhu, A., Benes, V., Costea, P. I., Hercog, R., Hildebrand, F., Huerta-Cepas, J., Nieuwdorp, M., Salojärvi, J., Voigt, A. Y., Zeller, G., Sunagawa, S., De Vos, W. M., & Bork, P. (2016). Durable coexistence of donor and recipient strains after fecal microbiota transplantation. *Science*, *352*(6285), 586–589. https://doi.org/10.1126/SCIENCE.AAD8852/SUPPL_FILE/LI-SM.PDF

Lloyd-Price, J., Mahurkar, A., Rahnavard, G., Crabtree, J., Orvis, J., Hall, A. B., Brady, A., Creasy, H. H., McCracken, C., Giglio, M. G., McDonald, D., Franzosa, E. A., Knight, R., White, O., & Huttenhower, C. (2017). Strains, functions and dynamics in the expanded Human Microbiome Project. *Nature 2017 550:7674*, *550*(7674), 61–66. https://doi.org/10.1038/nature23889

Martin, B. D., Witten, D., & Willis, A. D. (2020). MODELING MICROBIAL ABUNDANCES AND DYSBIOSIS WITH BETA-BINOMIAL REGRESSION. *The Annals of Applied Statistics*, *14*(1), 94. https://doi.org/10.1214/19-AOAS1283

McGhee, J. J., Rawson, N., Bailey, B. A., Fernandez-Guerra, A., Sisk-Hackworth, L., & Kelley, S. T. (2020). Meta-SourceTracker: application of Bayesian source tracking to shotgun metagenomics. *PeerJ*, *8*, e8783. https://doi.org/10.7717/peerj.8783

Methé, B. A., Nelson, K. E., Pop, M., Creasy, H. H., Giglio, M. G., Huttenhower, C., Gevers, D., Petrosino, J. F., Abubucker, S., Badger, J. H., Chinwalla, A. T., Earl, A. M., Fitzgerald, M. G., Fulton, R. S., Hallsworth-Pepin, K., Lobos, E. A., Madupu, R., Magrini, V., Martin, J. C., … White, O. (2012). A framework for human microbiome research. *Nature 2012 486:7402*, *486*(7402), 215–221. https://doi.org/10.1038/nature11209

Nayfach, S., Rodriguez-Mueller, B., Garud, N., & Pollard, K. S. (2016). An integrated metagenomics pipeline for strain profiling reveals novel patterns of bacterial transmission and biogeography. *Genome Research*, *26*(11), 1612–1625. https://doi.org/10.1101/gr.201863.115

857   Olm, M. R., Crits-Christoph, A., Bouma-Gregson, K., Firek, B. A., Morowitz, M. J., & Banfield,

858        J. F. (2021). inStrain profiles population microdiversity from metagenomic data and

859        sensitively detects shared microbial strains. *Nature Biotechnology 2021 39:6*, *39*(6), 727–

860        736. https://doi.org/10.1038/s41587-020-00797-0

861   Schloissnig, S., Arumugam, M., Sunagawa, S., Mitreva, M., Tap, J., Zhu, A., Waller, A., Mende,

862        D. R., Kultima, J. R., Martin, J., Kota, K., Sunyaev, S. R., & Weinstock, G. M. (2013).

863        Genomic variation landscape of the human gut microbiome. *Nature*.

864        https://doi.org/10.1038/nature11711

865   Schmidt, T. S. B., Hayward, M. R., Coelho, L. P., Li, S. S., Costea, P. I., Voigt, A. Y., Wirbel, J.,

866        Maistrenko, O. M., Alves, R. J. C., Bergsten, E., de Beaufort, C., Sobhani, I., Heintz-

867        Buschart, A., Sunagawa, S., Zeller, G., Wilmes, P., & Bork, P. (2019). Extensive

868        transmission of microbes along the gastrointestinal tract. *ELife*, *8*.

869        https://doi.org/10.7554/ELIFE.42693

870   Shenhav, L., Thompson, M., Joseph, T. A., Briscoe, L., Furman, O., Bogumil, D., Mizrahi, I.,

871        Pe'er, I., & Halperin, E. (2019). FEAST: fast expectation-maximization for microbial

872        source tracking. *Nature Methods*, *16*(7). https://doi.org/10.1038/s41592-019-0431-x

873   Shi, Z. J., Dimitrov, B., Zhao, C., Nayfach, S., & Pollard, K. S. (2021). Fast and accurate

874        metagenotyping of the human gut microbiome with GT-Pro. *Nature Biotechnology 2021*

875        *40:4*, *40*(4), 507–516. https://doi.org/10.1038/s41587-021-01102-3

876   Siranosian, B. A., Tamburini, F. B., Sherlock, G., & Bhatt, A. S. (2020). Acquisition,

877        transmission and strain diversity of human gut-colonizing crAss-like phages. *Nature*

878        *Communications 2020 11:1*, *11*(1), 1–11. https://doi.org/10.1038/s41467-019-14103-3

879   Sloan, W. T., Lunn, M., Woodcock, S., Head, I. M., Nee, S., & Curtis, T. P. (2006). Quantifying

880        the roles of immigration and chance in shaping prokaryote community structure.

881        *Environmental Microbiology*, *8*(4), 732–740. https://doi.org/10.1111/J.1462-

882        2920.2005.00956.X

883   Sloan, W. T., Woodcock, S., Lunn, M., Head, I. M., & Curtis, T. P. (2007). Modeling taxa-

884        abundance distributions in microbial communities using environmental sequence data.

885        *Microbial Ecology*, *53*(3), 443–455. https://doi.org/10.1007/S00248-006-9141-

886        X/FIGURES/4

887   Sprockett, D. D., Martin, M., Costello, E. K., Burns, A. R., Holmes, S. P., Gurven, M. D., &

888    Relman, D. A. (2020). Microbiota assembly, structure, and dynamics among Tsimane

889    horticulturalists of the Bolivian Amazon. *Nature Communications 2020 11:1*, *11*(1), 1–14.

890    https://doi.org/10.1038/s41467-020-17541-6

891    Sunagawa, S., Coelho, L. P., Chaffron, S., Kultima, J. R., Labadie, K., Salazar, G., Djahanschiri,

892    B., Zeller, G., Mende, D. R., Alberti, A., Cornejo-Castillo, F. M., Costea, P. I., Cruaud, C.,

893    d'Ovidio, F., Engelen, S., Ferrera, I., Gasol, J. M., Guidi, L., Hildebrand, F., … Bork, P.

894    (2015). Structure and function of the global ocean microbiome. *Science*, *348*(6237).

895    https://doi.org/10.1126/SCIENCE.1261359

896    Vatanen, T., Jabbar, K. S., Vlamakis, H., Knip, M., & Correspondence, R. J. X. (2022). Mobile

897    genetic elements from the maternal microbiome shape infant gut microbial assembly and

898    metabolism. *Cell*, *185*, 4921-4936.e15. https://doi.org/10.1016/j.cell.2022.11.023

899    Yassour, M., Jason, E., Hogstrom, L. J., Arthur, T. D., Tripathi, S., Siljander, H., Selvenius, J.,

900    Oikarinen, S., Hyöty, H., Virtanen, S. M., Ilonen, J., Ferretti, P., Pasolli, E., Tett, A.,

901    Asnicar, F., Segata, N., Vlamakis, H., Lander, E. S., Huttenhower, C., … Xavier, R. J.

902    (2018). Strain-Level Analysis of Mother-to-Child Bacterial Transmission during the First

903    Few Months of Life. *Cell Host & Microbe*, *24*(1), 146-154.e4.

904    https://doi.org/10.1016/J.CHOM.2018.06.007

905

906