

Supplemental Materials

Contents:

Computational Methods and Analysis

1	Polymorphism as a function of age across datasets	Page 8
2	Quasi-phasing	Page 17
3	False positive rate for SNV changes	Page 20

List of Figures:

S1	Sampling time points and sample sizes for infants and mothers analyzed in this paper	Page 3
S2	Shannon alpha diversity displayed by dataset.	Page 4
S3	Shannon alpha diversity in mothers' gut microbiomes at time of delivery versus healthy non-pregnant women.	Page 5
S4	Within-sample polymorphism over different life stages for the most prevalent species in our dataset	Page 6
S5	Within-sample polymorphism in the Shao et al. data set.	Page 7
S6	Between-dataset comparisons of <i>E. coli</i> polymorphism levels for matched life stages.	Page 8
S7	Number of QP versus non-QP samples	Page 9
S8	Proportions of high coverage species that are quasi-phaseable (QP) per host sample, categorized by life stage.	Page 10
S9	Decay in rates of SNV change, gene gain, gene loss and replacement rates over life stage.	Page 11
S10	Evolution and replacement rates in infants versus adults matched for duration of sampling.	Page 12
S11	Difference in rates of evolution and strain replacement for C-section versus vaginally born and breast versus formula fed babies.	Page 13
S12	Prevalence of genes that are gained or lost in putative SNV modification events with respect to infant and HMP adult cohorts.	Page 14
S13	Prevalence of sweeping SNVs in the infant, HMP adult, and mother cohorts	Page 15

S14	Comparison of HMP (adult) and infant prevalences of sweeping alleles involved in putative modification events.	Page 16
S15	Allele frequency distributions for three infant hosts.	Page 18
S16	Quasi-phasing of two hypothetical samples with difference allele frequencies.	Page 19
S17	Schematic of allele frequency changes detected.	Page 20

List of Tables:

Table S1: Statistical significance comparing polymorphism rates between datasets for matched age categories. These rates were computed with Mann-Whitney U tests and are compared against a Bonferroni-corrected significance level $\alpha=0.05/18$.

Table S2: Statistical significance comparing evolutionary rates and replacement rates between life stages, plotted in **Figure 2**. These significance values were computed with permutation tests and are corrected with the Benjamini-Hochberg method.

Table S3: Parallelism of SNV changes. Reported are numbers of SNV changes in evolutionary modification events, numbers of unique hosts experiencing such SNV changes, and expected number of SNV changes under the null, grouped by PATRIC gene ID.

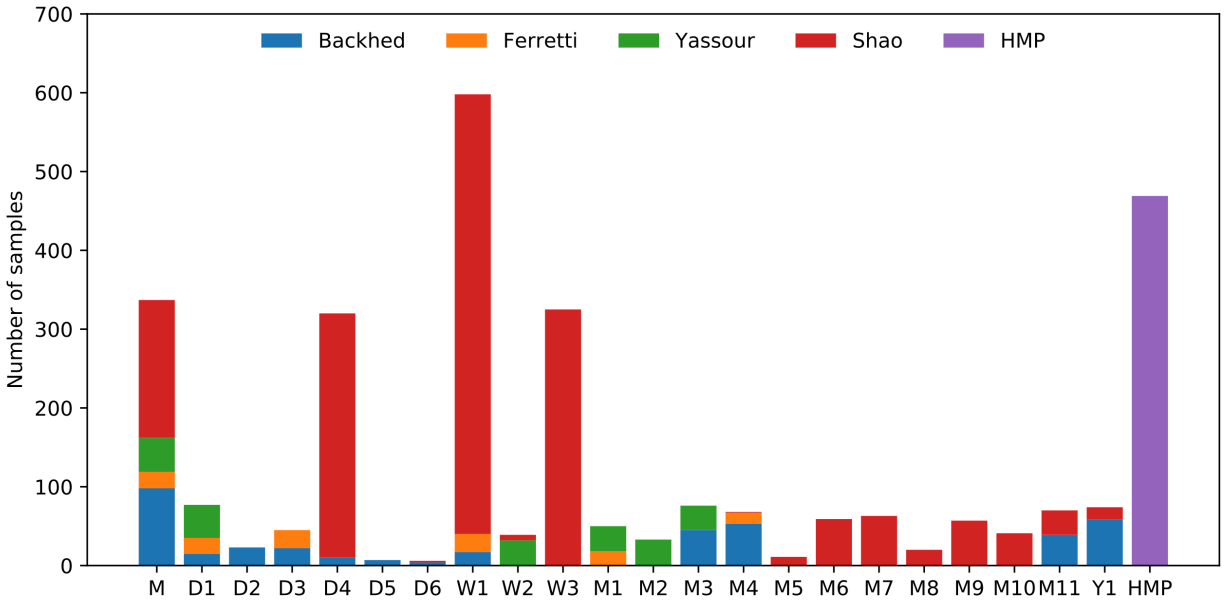


Figure S1: Sampling timepoints and sample sizes for infants and mothers analyzed in this paper. M indicates mother. Only those mothers sampled around the time of delivery were included in this analysis. “D1” indicates day 1, “W1” indicates week 1, “M1” indicates month 1, “Y1” indicates year 1, and so on. “M” indicates mother samples at delivery and “HMP” includes all HMP1-2 adult samples. Some infants from Shao et al. 2019 were sampled multiple times between day 7 and 14, and thus were grouped in the Week 1 category.

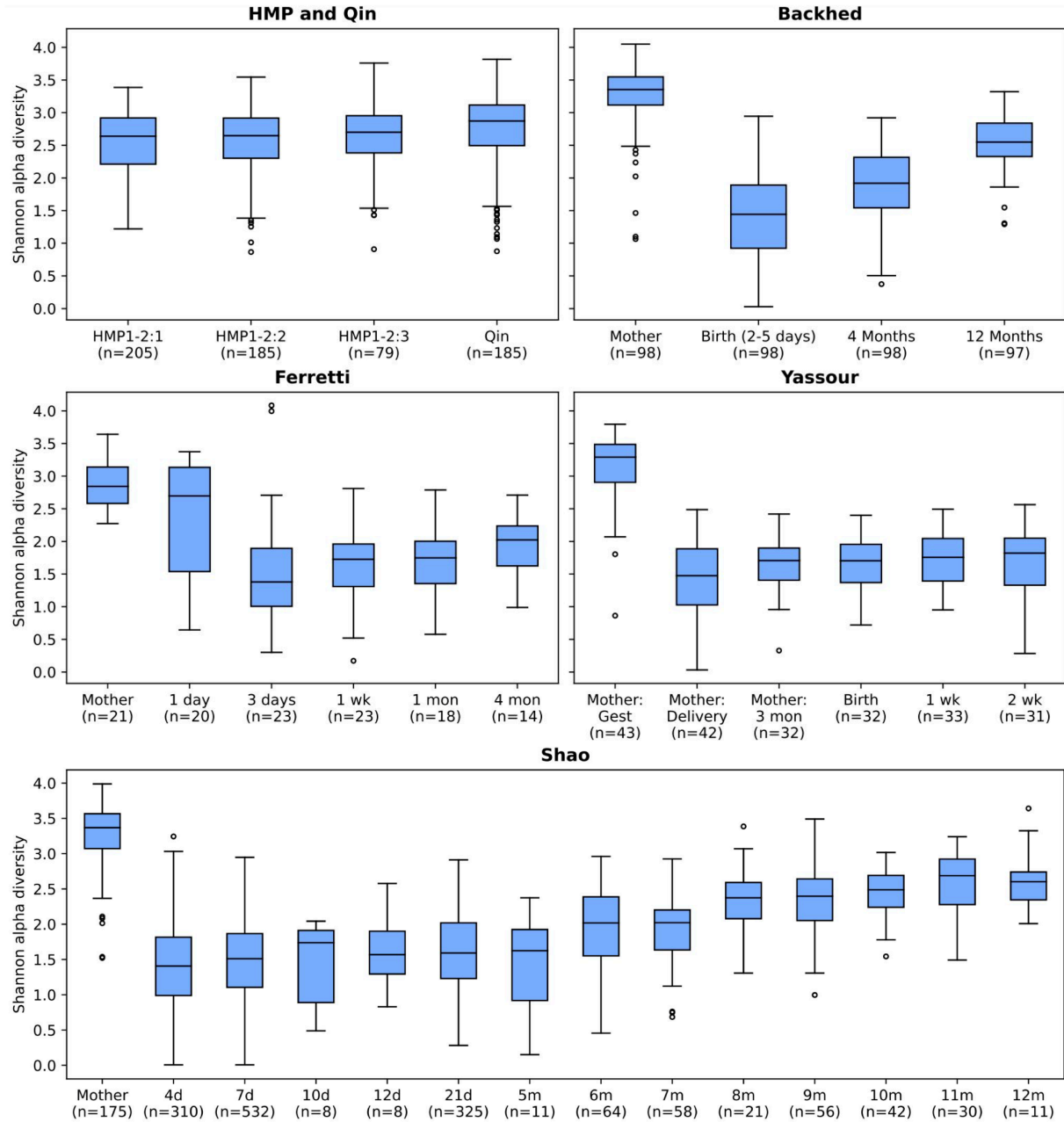


Figure S2: Shannon alpha diversity displayed by dataset. Samples sizes (n) are shown in parentheses. Overall, alpha diversity increases as a function of age of the host (p value = 2×10^{-16} , $\text{glmmTMB}(\text{alpha_div} \sim \text{day} + (1|\text{sample}) + (1|\text{subject}) + (1|\text{dataset}))$), using all samples from infants).

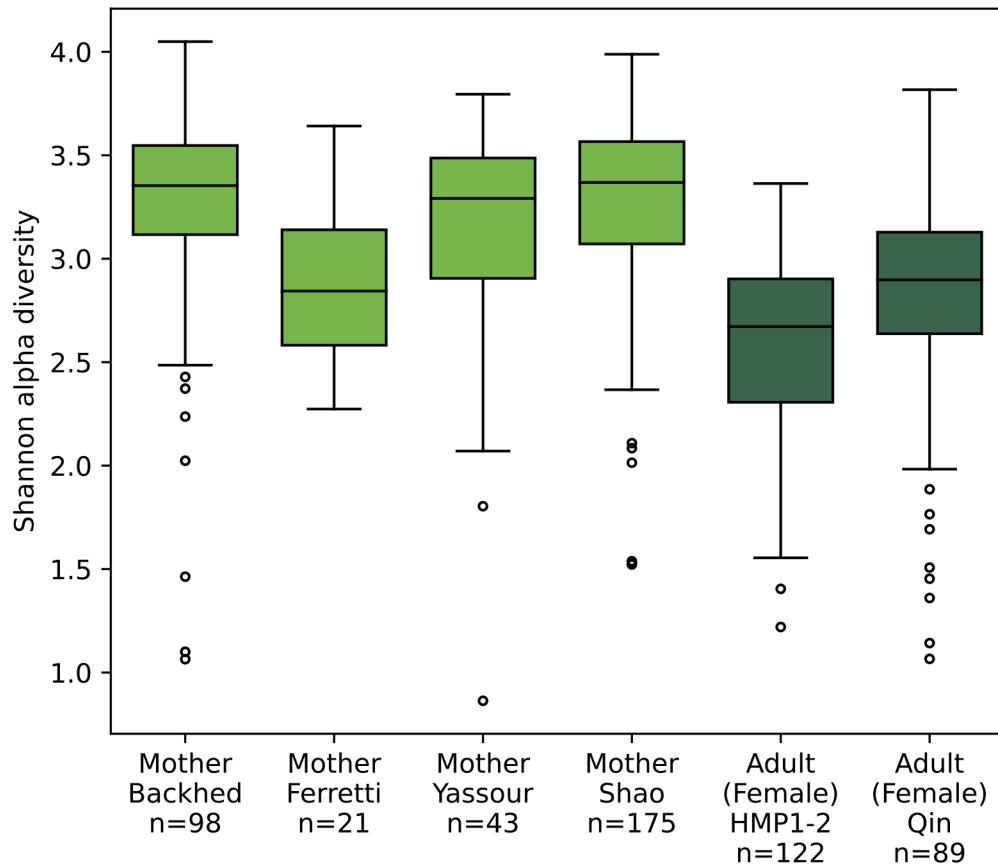


Figure S3: Shannon alpha diversity in mothers' gut microbiomes at time of delivery versus healthy non-pregnant women. Plotted are alpha diversity values for mothers at time of delivery in the Backhed, Ferretti, Yassour, and Shao datasets, and non-pregnant women in the HMP and Qin et al. datasets. To assess whether pregnancy results in a significant difference in alpha diversity, we fit a GLMM with pregnancy status as the predictor and study and host as random effects ($\text{alpha diversity} \sim \text{pregnancy status} + (1|\text{study}) + (1|\text{host})$) and obtained a coefficient of 0.44 for pregnancy status with a p value of 0.035.

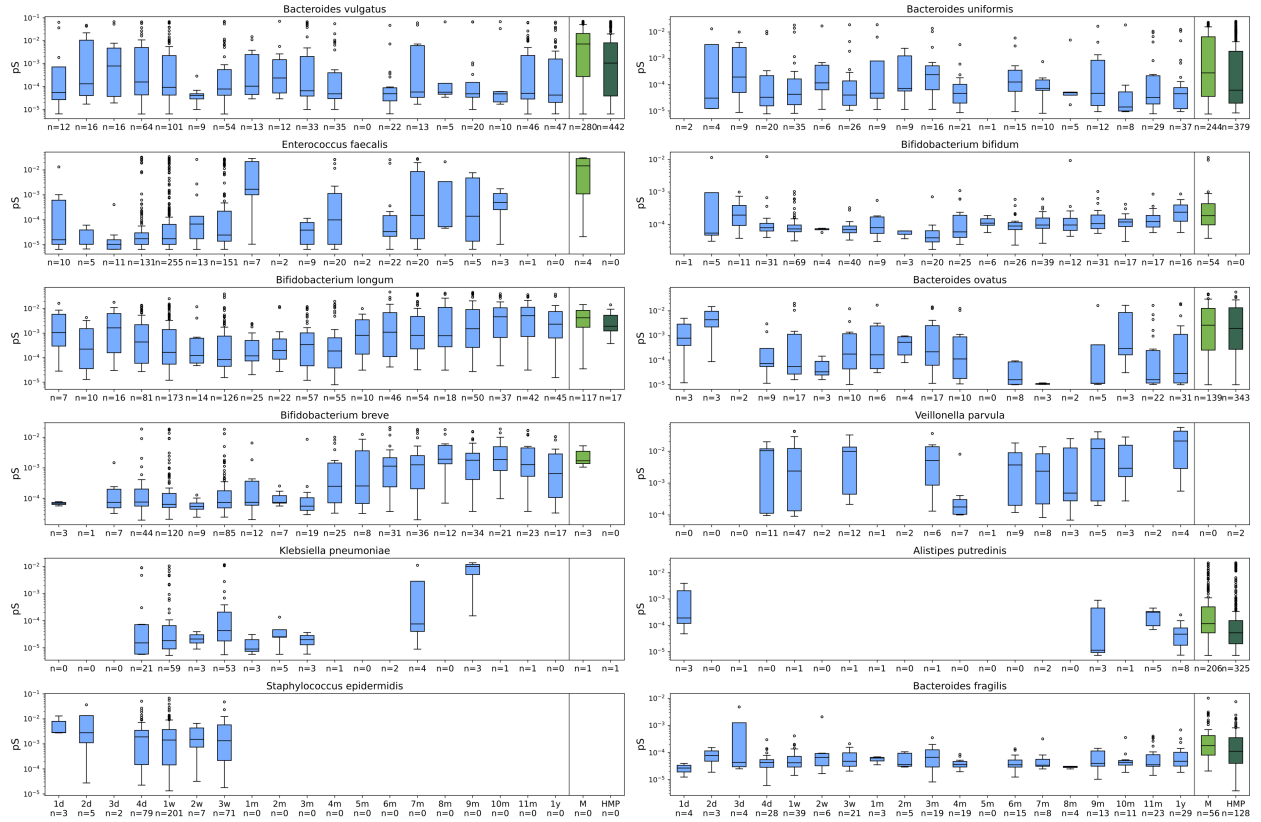


Figure S4: Within-sample polymorphism over different life stages for the most prevalent species in our dataset. The next 12 most prevalent species, after *E. coli*, which is shown in **Figure 1B**, are displayed here. The same life stages as in **Figure 1** are shown, where “B” indicates birth, “D1” indicates day 1, “W1” indicates week 1, “M1” indicates month 1, “Y1” indicates year 1, and so on. “M” indicates mother samples at delivery and “HMP” includes all HMP1-2 adult samples. pS indicates synonymous polymorphism rate, where a polymorphism was defined as a site with allele frequency between 0.2 and 0.8.

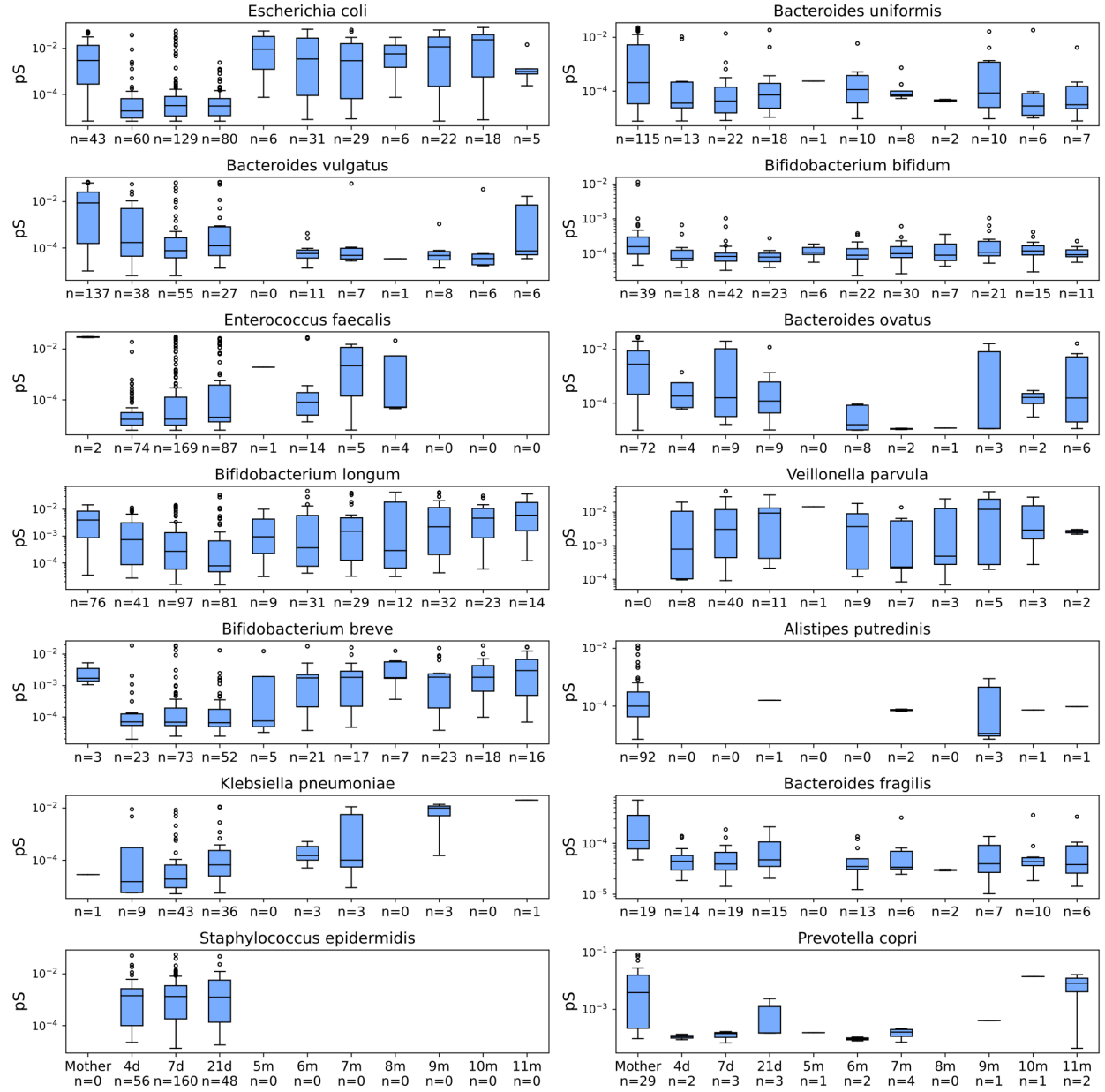


Figure S5: Within-sample polymorphism in the Shao et al. data set. To assess if polymorphism rates generally follow the same trends within a dataset versus across datasets, we examine polymorphism rates for the 14 most prevalent species in our dataset using data from Shao et al only.

Polymorphism as a function of age across datasets

To assess the effect of dataset on polymorphism, we compared distributions of polymorphism per sample-species pair for matched age categories across different datasets for the 13 most prevalent species (in infants). Up to 12 pairwise comparisons for each species were considered for the following age categories: 1 day (Backhed and Ferretti), 3 days (Backhed and Ferretti), 4 days (Backhed and Yassour), 1 week (Ferretti, Yassour and Shao), 2 weeks (Yassour and Shao), 1 month (Ferretti, Yassour and Shao), 4 months (Backhed and Ferretti), and 12 months (Backhed and Shao). Due to non-normal polymorphism distributions, we used the Mann Whitney U statistic to test for significant differences in polymorphism between datasets. Of 18 tests across the 13 species satisfying a minimum sample size requirement of 10, one test was significant (Backhed vs. Ferretti at 12 months, $P=0.00018$) using a Bonferroni-corrected significance level $\alpha=0.05/18$. Of 45 tests satisfying a laxer sample size threshold of 5, only the aforementioned comparison remained significant. All test results are reported in **Table S1**.

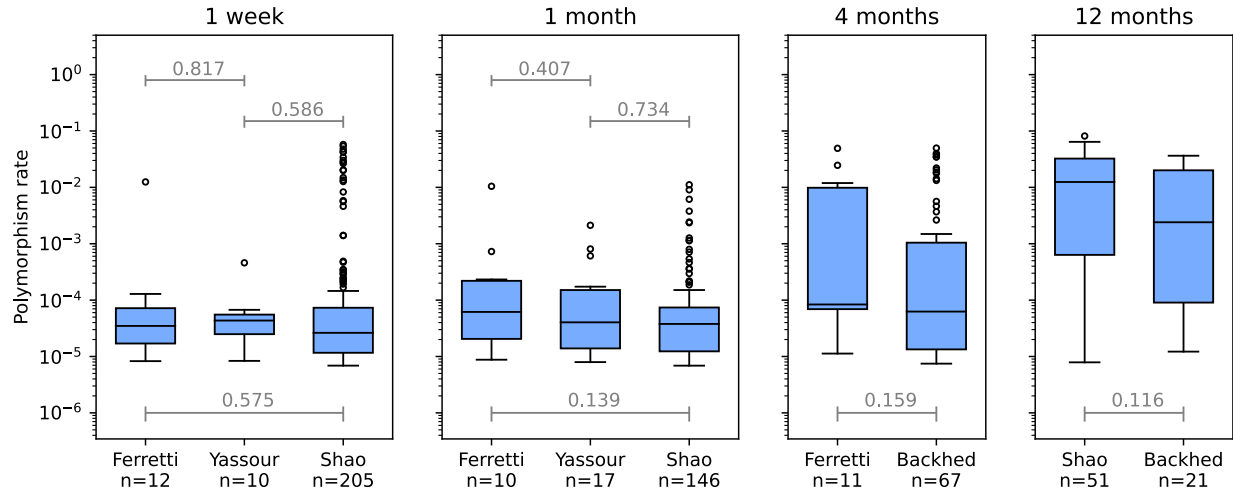


Figure S6: Between-dataset comparisons of *E. coli* polymorphism levels for matched life stages. *P* values from Mann-Whitney *U* tests are reported in grey.

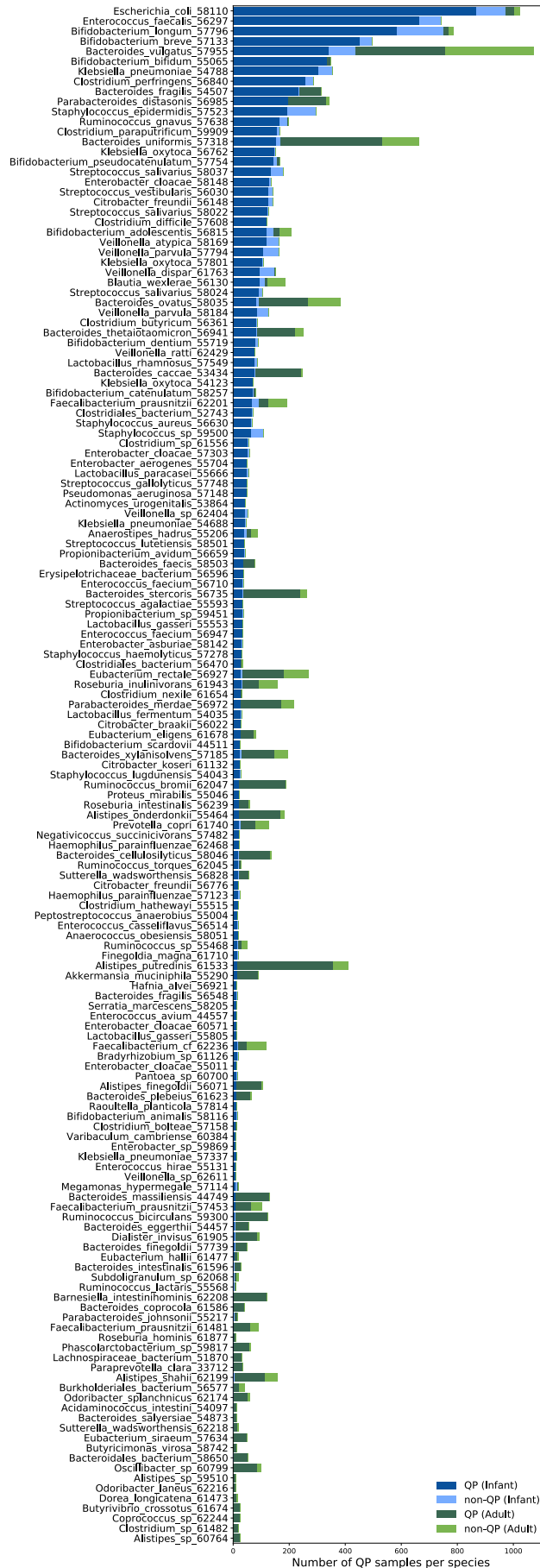


Figure S7: Number of QP versus non-QP samples. The 153 species that have greater than 10 QP samples across all hosts are shown; species are ordered by the number of QP samples in infants

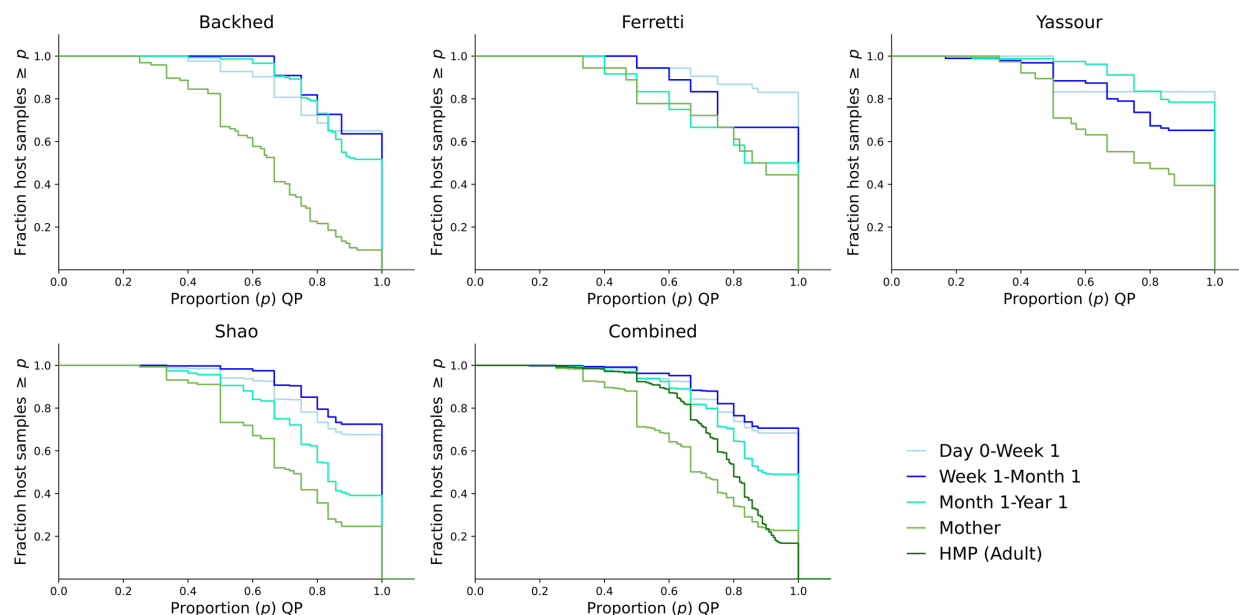


Figure S8: Proportions of high coverage species that are quasi-phaseable (QP) per host sample, categorized by life stage. Other than ‘Mother’ and ‘HMP’, category label refers to life stage within infants. Here, ‘host sample’ is defined as a single metagenomic sample. Distributions are plotted per infant dataset, as well as in the combined dataset.

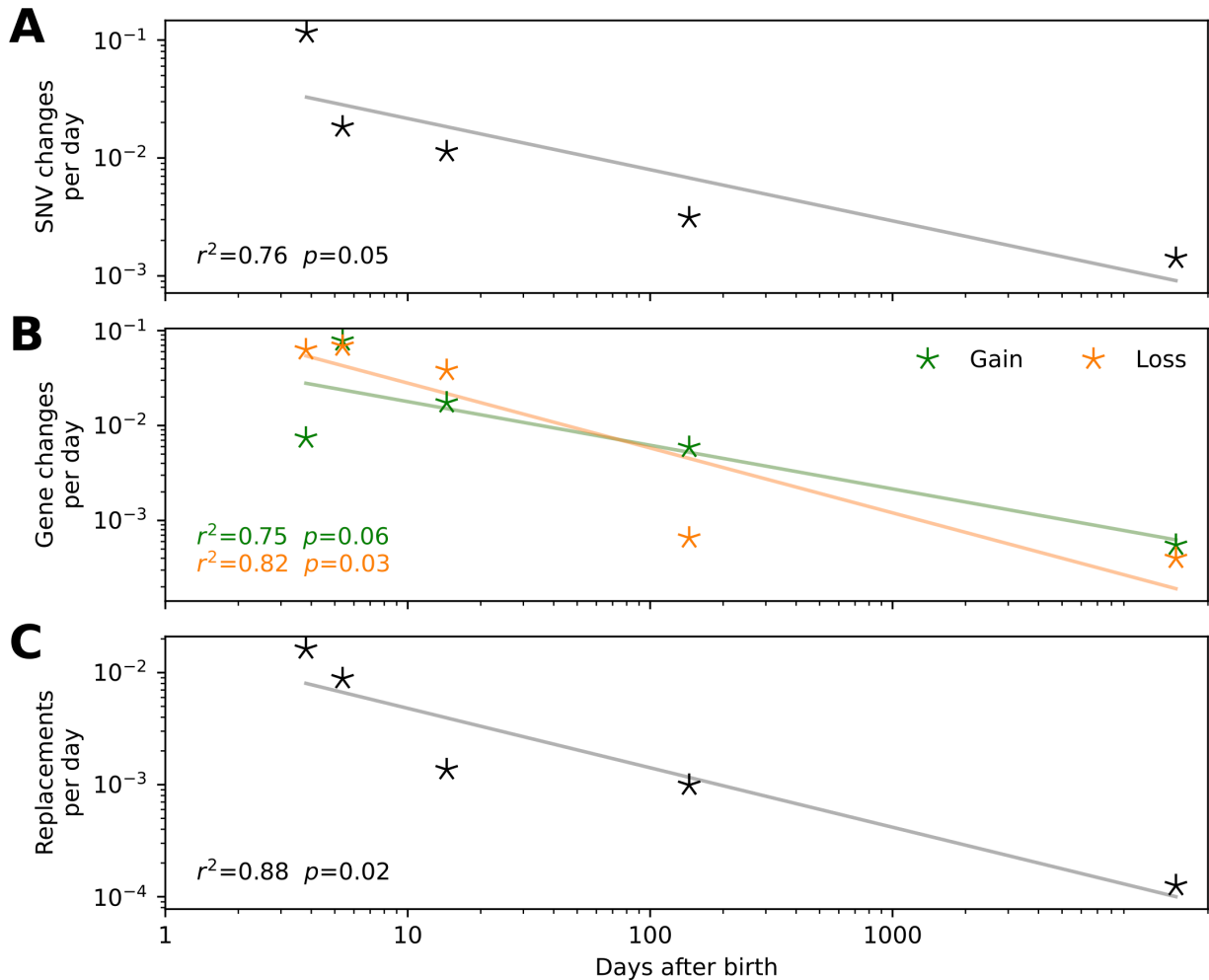


Figure S9: Decay in rates of SNV change, gene gain, gene loss and replacement rates over life stage. The same timepoint pair categories as in **Figure 2** were considered: mother-infant within the first week, infant-infant within the first week, infant-infant week 1-month 1, infant-infant month 1-year 1, and adult-adult. Days after birth were assigned to be the median day in a given life stage interval. Mothers were assigned day 0, infant meconium samples were assigned day 1, and adult timepoints were arbitrarily assigned to be approximately 40 years after birth. A linear model was fit with either rate of evolutionary change or replacement as the response variable and days after birth as the predictor. P values are reported, as are the correlation coefficients, r^2 .

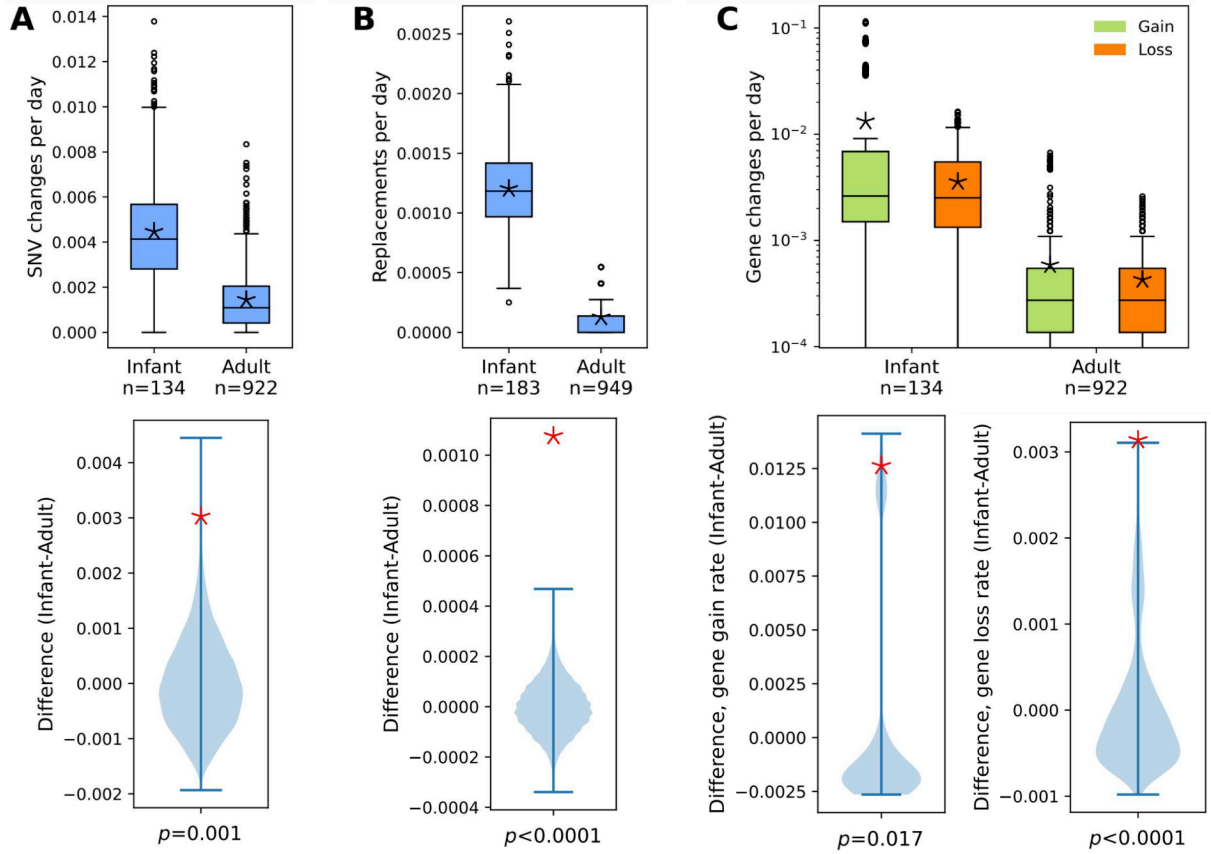


Figure S10: Evolution and replacement rates in infants versus adults matched for duration of sampling. Shown are 1000 bootstrapped (A) SNV change, (B) replacement, and (C) gene change rates for all infant-infant vs. HMP adult-adult QP pairs matched for duration between timepoints of 4 to 8 months. To assess if there are significant differences between infant and adult distributions for SNV change, replacement, gene gain, and gene loss rates, we performed a permutation test consisting of 10,000 permutations. The distributions of the difference in rates between infants and adults are plotted, with the observed value indicated with an asterisk. P values are reported below each permutation distribution plot. Adjusted p values with the Benjamini-Hochberg method are 0.0013, 0.0002, 0.017, and 0.0002, respectively.

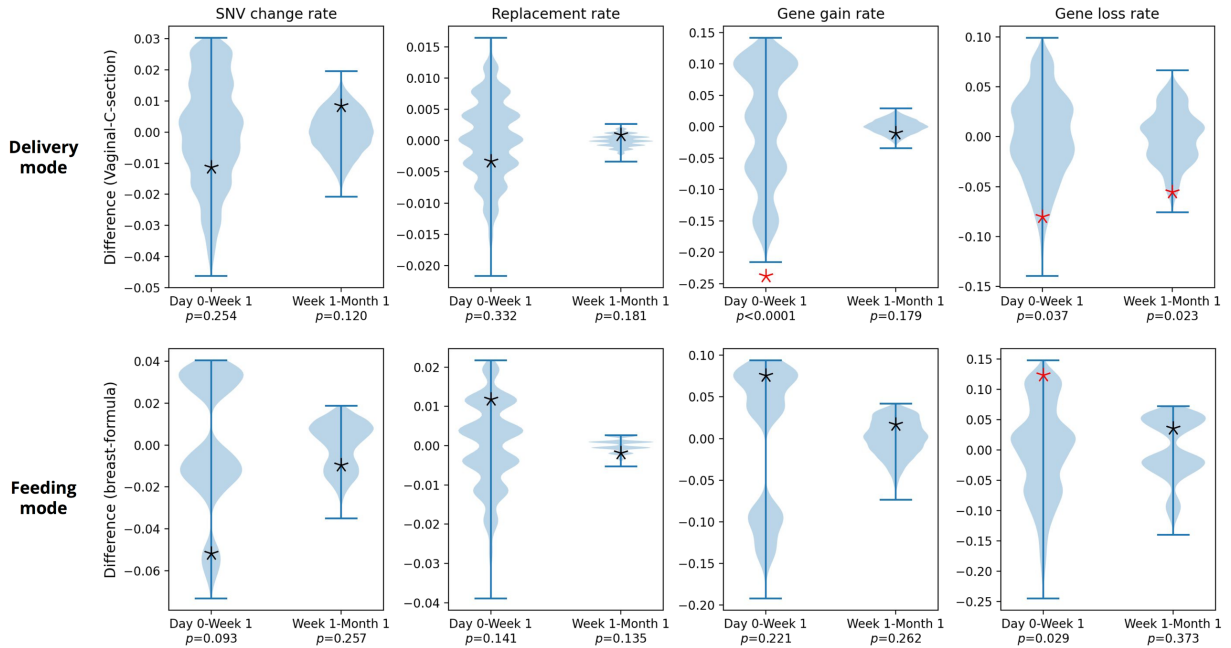


Figure S11: Difference in rates of evolution and strain replacement for C-section versus vaginally born and breast versus formula fed babies. We assessed whether rates of SNV changes associated with evolutionary modifications, strain replacement, gene gains, and gene loss are significantly different between C-section and vaginally born babies as well as between breast versus formula-fed babies. In the breast versus formula comparison, we did not include babies on mixed diets. We performed 10,000 permutations, and the resulting distributions are shown above. The asterisk shows the observed difference, and reported are associated p values. In red are those p values that pass the 0.05 significance threshold. Note that only the difference in gene gain rate for delivery mode for day 0 – week 1 survives multiple hypothesis correction with a Bonferroni adjusted p value. We analyzed infants in the day 0 to week 1 and week 1 to month 1 categories because these had the maximal number of samples available. The sample sizes were as follows: Day 0-Week 1: Vaginal: 127; C-section: 52; Breast: 112; Formula: 67. Week 1-Month 1: Vaginal: 149; C-section: 59; Breast: 136; Formula: 72.

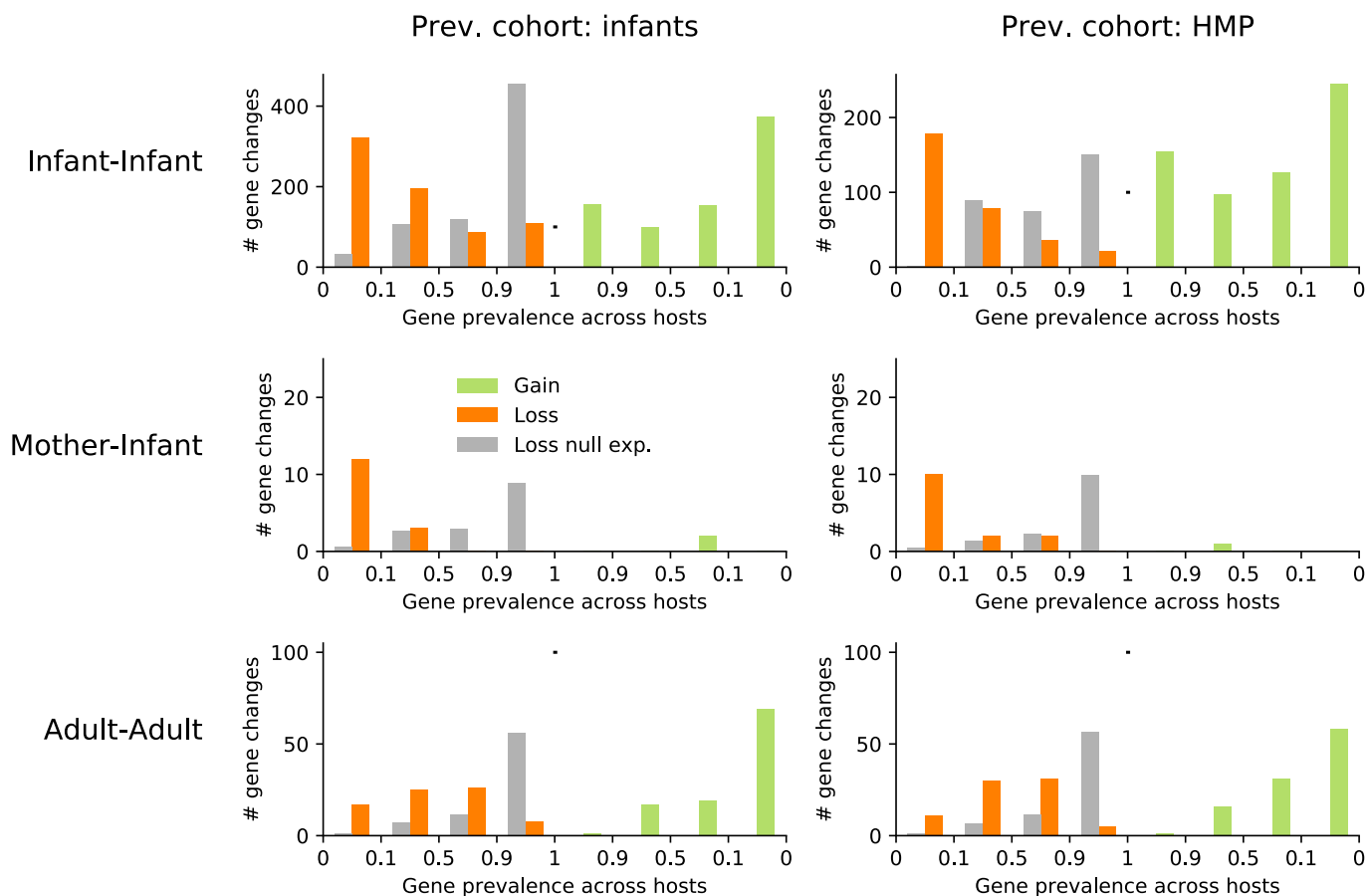


Figure S12: Prevalence of genes that are gained or lost in modification events with respect to infant and HMP adult cohorts.

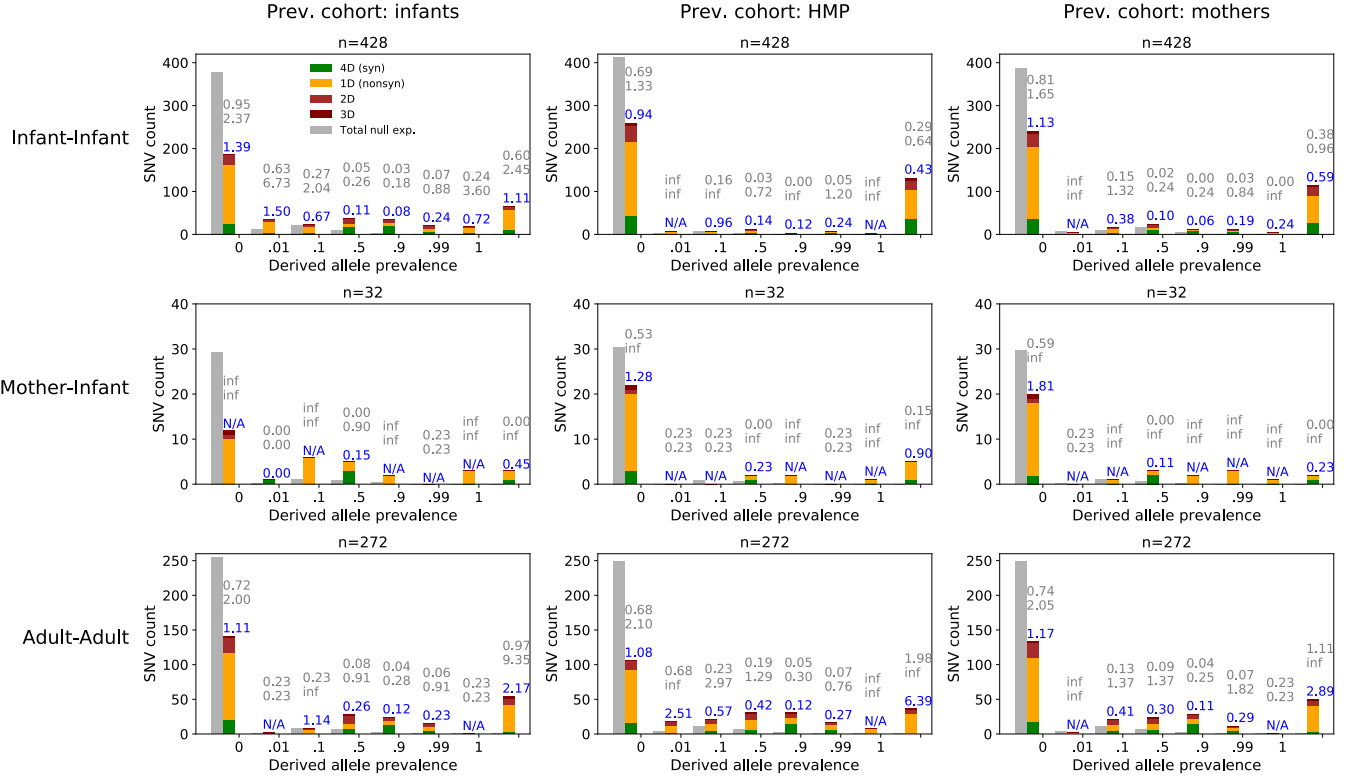


Figure S13: Prevalence of sweeping SNVs in the infant, HMP adult, and mother cohorts. The top left plot shows the prevalence of SNVs sweeping in infants with respect to a prevalence cohort defined by infants. By contrast, the top right plot shows the prevalence of same SNVs sweeping in infants, but with respect to a prevalence cohort defined by mothers. d_N/d_S of each prevalence bin are reported in blue with 95% confidence intervals reported in gray.

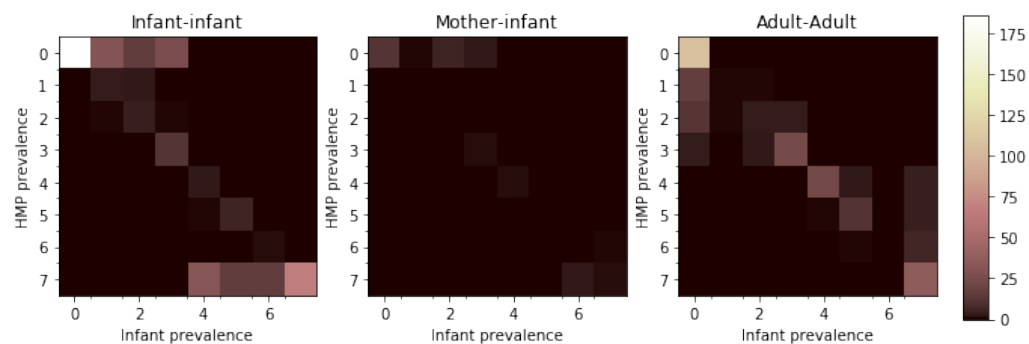


Figure S14: Comparison of HMP (adult) and infant prevalences of sweeping alleles involved in putative modification events. The cohort in which the sweeping allele was identified is indicated at top (e.g ‘Infant-Infant’, ‘Mother-Infant’, ‘Adult-Adult’). The colorbar indicates number of SNV changes.

Quasi-phasing

A major goal in this paper is to infer evolutionary changes from metagenomic samples. Here, an evolutionary change represents a change in allele frequency over time among lineages belonging to a strain. Unfortunately, two major potential confounders of evolutionary changes in metagenomic data are strain fluctuations and sampling error, since both can generate allele frequency changes.

To control for these potential confounders, we “quasi-phase” samples with sufficiently simple lineage structures. Quasi-phasing means that pairs of alleles can be confidently assigned to a single lineage’s genome. In doing so, we can identify evolutionary changes that accrue on the background of a single lineage. The approach we take is similar to that of Truong et al. 2017 in which a dominant allele is assigned to a dominant strain, but in Garud, Good et al. 2019 we put bounds on the error for phasing (for further statistical details, please see the supplement of Garud, Good et al. 2019).

To quasi-phase, we leverage knowledge about the lineage structure of a given species within a host. As described in Garud, Good et al. 2019, hosts are typically colonized by a small handful of genetically distinct lineages belonging to the same species. In **Figure S15**, we plot a distribution of allele frequencies for the common species *Bacteroides vulgatus* in three infant samples from Backhed et al. 2015, which as we describe below, illustrate a range of typical within-host lineage structures.

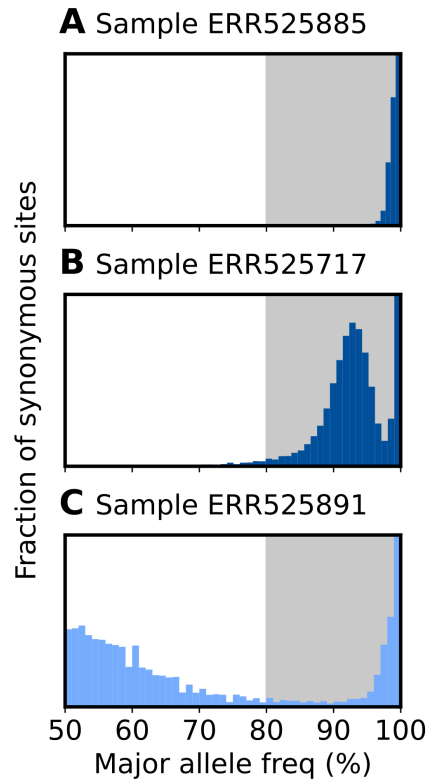


Figure S15: Allele frequency distributions for three infant hosts. Plotted are allele frequency distributions for the species *B. vulgatus* from three infant hosts from the Backhed et al. dataset. Only major allele frequencies are plotted. Figures (A) and (B) depict distributions from quasi-phaseable samples, in which the majority of polymorphic sites have a within-host frequency ≥ 0.8 . The sample in figure (C) is not quasi-phaseable since a large fraction of sites have a frequency < 0.8 .

In **Figure S15A**, there is a mass of sites with allele frequencies (f) close to 0. As described in Garud, Good et al. 2019, these sites are likely comprised of a mixture of sequencing errors and low frequency mutations that have arisen due to the expansion of a single lineage within a host. By contrast, in Figures S1B and C there are a mass of sites with allele frequencies close to 0 and another mass of sites with allele frequencies peaked at intermediate frequencies. As described in Garud, Good et al. 2019, these peaked distributions are inconsistent with a single lineage expanding within a host and instead represent multiple, divergent lineages present within a host. The allele frequencies at which these distributions are peaked are representative of the relative frequencies at which lineages are colonizing the host.

In scenario A, a new mutation that arises represents a true evolutionary modification rather than a strain fluctuation, because there is only one strain colonizing the host. However, in scenarios B and C, a shift in frequency of the multiple strains could also generate an allele frequency change. In scenarios B and C we attempt to solve this problem by identifying samples in which a nucleotide confidently to a single lineage's haplotype, or, in other words, can be 'quasi-phased'.

Scenarios A and B are fairly straightforward to quasi-phase. In this scenario, there is typically a single dominant allele (**Figure S16**), which can be assigned to the dominant lineage. Even in the case where an allele is mis-assigned, there will still be a lineage harboring both alleles, even if it is not the dominant lineage (**Figure 16**, Garud et al. 2019, SI text). However, in scenario C, two strains are present at roughly 50% frequency. In this scenario, there is no allele that is dominant and there is a ~50% chance of incorrectly assigning an allele to the dominant lineage (**Figure S16**).

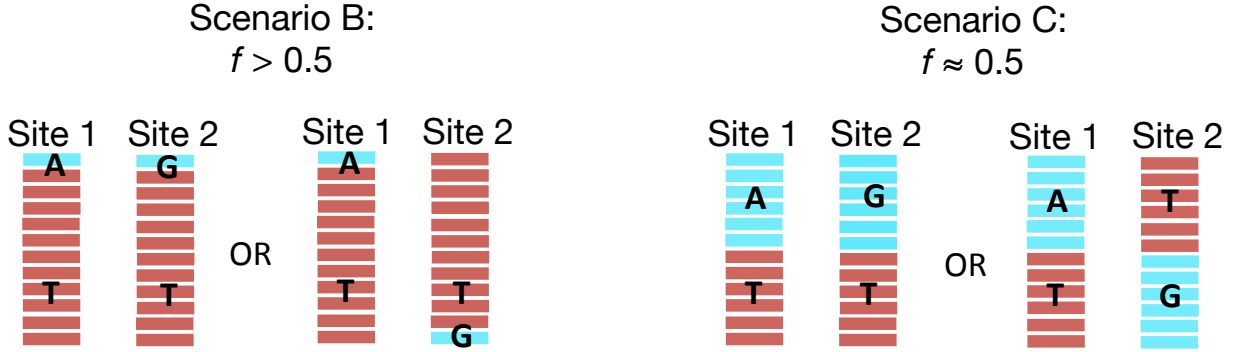


Figure S16: Quasi-phasing of two hypothetical samples with different allele frequencies. Hypothetical read pileups for two sites in Scenarios B and C. In Scenario B, there is a dominant lineage present at >80% frequency, whereas in Scenario C two lineages colonize the host at roughly 50% frequency (Fig S1). In Scenario B, the dominant alleles at both sites 1 and 2 can be assigned to the dominant lineage. If there is a phasing error, there still will be a fraction of cells that possess the pair of Ts on the same haplotype (see Garud, Good et al. 2019). By contrast, in Scenario C, there is a 50% chance that the alleles will be assigned to the incorrect lineage. In this scenario, it is unlikely that there will be any lineage that harbors both Ts.

Thus, samples with a large number of intermediate frequency alleles are more suspect to phasing errors. Quasi-phaseable samples are those that have few alleles at intermediate frequency. To identify quasi-phaseable samples, we wish to identify an allele frequency cutoff, f^* , that signifies the upper bound of what constitutes ‘intermediate frequency’. With such a cutoff, we can then assess the probability that an observed frequency, \hat{f} , is greater than f^* given k alternate alleles and D number of reads and a true allele frequency of f :

Eq 1:

$$\Pr[\hat{f} \geq f^* \mid D, f] = \sum_{k > f^* D} \binom{D}{k} f^k (1 - f)^{D-k}$$

As described in Garud, Good et al. 2019, this probability can be computed across the genome to obtain a genome-wide error rate of incorrectly phasing an allele. With sufficient depth D , which we assign to be a minimum of 20 in our analysis, and a sufficiently high f^* , which we set to be 0.8, sampling error is minimized. Quasi-phaseable samples are identified if they contain sufficiently low numbers of sites with $\hat{f} < f^*$, as described in greater detail in Garud, Good et al. 2019.

3. False positive rate for SNV changes

We next quantified SNV changes between quasi-phaseable time point pairs from the same host. To do so, we identified extreme allele frequency changes from ≤ 0.2 to ≥ 0.8 (or vice versa) between two samples, S1 and S2 (**Figure S17**).

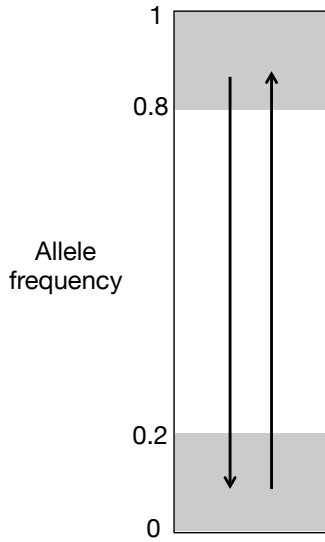


Figure S17: Schematic of allele frequency changes detected. To detect evolutionary changes within a host over time, we identified sites that changed allele frequency from ≤ 0.2 to ≥ 0.8 , or vice versa.

We compute the probability of observing an allele frequency due to sampling error (P_{err}) as follows:

Eq2:

$$P_{\text{err}} = \int (P(\hat{f}_1 \leq 0.2 | D_1, f) * (\hat{f}_2 \geq 0.8 | D_2, f) + P(\hat{f}_1 \geq 0.8 | D_1, f) * (\hat{f}_2 \leq 0.2 | D_2, f)) * P(D_1, D_2, f) dD_1 dD_2 df$$

Where \hat{f}_1 and \hat{f}_2 are allele frequencies in samples 1 and 2, respectively, and D_1 and D_2 are read depths in samples 1 and 2 respectively. f is assumed to be the same in both samples 1 and 2 under the null hypothesis where there is no evolutionary change. $P(\hat{f}_1 \geq f^*)$ is computed as eq1. $P(D_1, D_2, f) \sim P(D_1)P(D_2)P(f)$ and is estimated empirically from the data as described in Garud et al. 2019 SI text 1.

To compute a genome-wide false positive rate, we can multiply the per-site error rate by the length of the genome, L to estimate the total expected number of false positives:

$$N_{\text{err}} = P_{\text{err}} * L$$

For a depth of 20 in samples 1 and 2 (the minimum depth we require) and a true allele frequency $f=0.2$ in both samples, the probability of observing an allele frequency change by chance from 0.2 to 0.8 is 1.7×10^{-8} . Multiplying this by a mean genome size of 10^6 , the expected number of false positives for a given genome is 0.017, which is $\ll 1$.

A similar logic as in Eq2 is applied for inferring gene changes, as described in greater depth in Garud, Good et al. 2019.