

Supplemental Analysis

Infant gut 16S rDNA sequencing data from Koenig *et al.* [1] were downloaded from the NCBI Short Read Archive (SRA) website (<http://www.ncbi.nlm.nih.gov/sra>) along with their metadata. These data are a time-series of fecal bacterial communities from an infant subject, over the first 500 days of life. QIIME [2] was used to trim primer regions from these data. Clustering was performed using the unoise3 pipeline [3]; sequences were de-replicated at 100% identity using VSEARCH [4], zOTU centroid sequences were picked and chimeric sequences were removed using unoise3 [3], then all sequences were mapped onto zOTU seeds to create a zOTU table using VSEARCH. zOTU stands for “zero-radius operational taxonomic unit” [3]. Unlike traditional *de novo* clustered OTUs, zOTUs are exact sequence variants (ESVs) which are consistent and easily comparable across data. The SINA aligner [5] was used to align zOTU centroid sequences to the SILVA SSU Ref 128 database (available from <https://www.arb-silva.de/download/arb-files/>). We then used IQ-TREE [6] to build a phylogenetic tree from the aligned sequences.

The resulting zOTU table was rarefied to 1000 sequences per sample, and samples with fewer sequences were excluded. The last five timepoints were excluded as well because they were sampled at a much lower temporal resolution. This left 52 timepoints spread over the first 469 days of the infant subject’s life. The zOTU table was then split into two zOTU tables, one for timepoints before the infant started consuming baby formula, and one for those after. The “pre-formula” zOTU table contained ages 4 days through 146, and the “post-formula” zOTU table contained ages 161 days through 469 days. Each zOTU table was used to run our model as described above using 500 D values (Equation 1), ranging from underdispersed ($D = -5$) to overdispersed ($D = 5$), using zOTUs in lieu of species. zOTUs with zero phylogenetic distance between them were combined, because these zOTUs were uninformative for our statistical model (Equation 1; zero raised to a negative exponent is undefined). The 500 resulting surrogate datasets were compared to the empirical dataset as described above, using difference between phylogeny values at the middle timepoint. The logistic error model (Equation 2) was fit and bootstrapped, yielding an estimate for D and 95% confidence intervals for that estimate.

“Moving pictures” sequence data from Caporaso *et al.* [7] were downloaded from the MG-RAST database (<http://metagenomics.anl.gov/>). These are longitudinal data from one adult male subject and one adult female subject, over a period of several hundred days, across multiple sample sites (feces, both palms, tongue). timepoints were excluded which did not have sequence data for each of the 8 environments (left palm, right palm, mouth, and feces of the male and female subjects), and rarefied to 5000 sequences per sample. This left 107 timepoints, ranging from day 1 to day 185. Analysis for each dataset (e.g. female right palm) was carried out as described above, except raw sequences were trimmed to a length of 91 bp after the end of the forward PCR primer site in order to ensure that all raw sequences spanned the same region of the 16S rRNA gene. 91 bp was chosen as a length cutoff in order to keep 95% of the sequence data (5% of sequences were discarded because they were shorter).

Analysis of the “moving pictures” data was also done using two approaches that allowed us to test the importance of the set of species that are included in the species pool. In principal, the model may perform differently if a broader representation of what is in the environment is in the species pool compared to what eventually colonizes the individual over time, as the latter may result in a species pool that is overall constrained by factors such as competition for niche space. Thus in an alternate approach we included sequences in the species pool from the other individual living in the same household, as these would be in the environment but not competing for the same niche. We analyzed palm communities in a “meta” context, where surrogate datasets were generated assuming the species pool for a given palm was composed of all four palms in the dataset. In this case, the difference between the “self” D estimate (generated above) and the “meta” D estimate (estimated with a metapopulation of zOTUs) is related to the exclusivity of species detected in the community. In other words, if we were to estimate similar D values for both the “meta” and “self” analyses, the inclusion of extra species in the species pool would be of little importance to the model, and we would learn that it would make little difference to community assembly patterns if the species pool really was composed of the “meta” set. We also analyzed a section of samples from the male right palm data that were collected every day over a period of 19 days, using a sliding window approach. We ran our model as described above on each window of 5 continuous days (15 windows), in order to see how D varied over time. We only conducted this analysis for the section of samples that were sampled every day, so that comparisons between windows would not be confounded by window size.

Finnish infant sequence data from Yassour *et al.* [8] and associated metadata were downloaded from the

DIABIMMUNE Microbiome Project website (<https://pubs.broadinstitute.org/diabimmune>). These are longitudinal gut microbiome data from Finnish infants, collected over the first 36 months of life [8]. Roughly half of these infants were repeatedly treated with oral antibiotics, almost universally for ear infections. Metadata for this dataset were compiled in a different re-analysis of these data [9] and were downloaded from the authors' GitHub page (https://github.com/ShadeLab/microbiome_trait_succession). Subject datasets belonging to the groups "Antibiotic" ($n=18$) or "Control" ($n=15$) were each analyzed using our model, similar to above. These subjects had between 19 and 36 samples collected over 36 months, with a mean of 28 samples. Sequence data were rarefied to 5000 sequences, and our model was run per above. We compared the estimated D values between antibiotic and control groups using a Mann-Whitney test. Because this dataset had so many subjects, we used this analysis as an opportunity to analyze whether the number of zOTUs, total phylodiversity, or number of timepoints had an effect on estimated D values. This was done via correlation analysis of D estimates with the aforementioned potential covariates.

References

- [1] Koenig JE, Spor A, Scalfone N, Fricker AD, Stombaugh J, Knight R, et al. Succession of microbial consortia in the developing infant gut microbiome. *Proceedings of the National Academy of Sciences of the United States of America*. 2011 mar;108 Suppl(Supplement_1):4578–85.
- [2] Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et al. QIIME allows analysis of high-throughput community sequencing data. *Nature methods*. 2010 may;7(5):335–6.
- [3] Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon sequencing. *bioRxiv*. 2016; Available from: <http://www.biorxiv.org/content/early/2016/10/15/081257>.
- [4] Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for metagenomics. *PeerJ*. 2016;4:e2409v1.
- [5] Pruesse E, Peplies J, Glöckner FO. SINA: Accurate high-throughput multiple sequence alignment of ribosomal RNA genes. *Bioinformatics*. 2012;28(14):1823–1829.
- [6] Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ. IQ-TREE: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution*. 2015;.
- [7] Caporaso JG, Lauber CL, Costello EK, Berg-Lyons D, Gonzalez A, Stombaugh J, et al. Moving pictures of the human microbiome. *Genome biology*. 2011 jan;12(5):R50.
- [8] Yassour M, Vatanen T, Siljander H, Hämäläinen AM, Härkönen T, Ryhänen SJ, et al. Natural history of the infant gut microbiome and impact of antibiotic treatment on bacterial strain diversity and stability. *Science Translational Medicine*. 2016;.
- [9] Guittar J, Shade A, Litchman E. Trait-based community assembly and succession of the infant gut microbiome. *Nature Communications*. 2019;.