

A census-based estimate of Earth's bacterial and archaeal diversity

- S1 Text -

Stilianos Louca^{1,2,3,4}, Florent Mazel^{3,5}, Michael Doebeli^{3,4,6} & Laura Wegener Parfrey^{3,4,5}

¹*Department of Biology, University of Oregon, Eugene, Oregon, USA*

²*Institute of Ecology and Evolution, University of Oregon, Eugene, Oregon, USA*

³*Biodiversity Research Centre, University of British Columbia, Vancouver, Canada*

⁴*Department of Zoology, University of British Columbia, Vancouver, Canada*

⁵*Department of Botany, University of British Columbia, Vancouver, Canada*

⁶*Department of Mathematics, University of British Columbia, Vancouver, Canada*

1 The pitfalls of extrapolating host-specific microbial diversity estimates

2 Estimating global species diversity has long been a challenging endeavor, with microorganisms constituting a
3 particularly cryptic aspect of Life [1, 2]. In a recent meta-analysis, Larsen et al. [3] attempted to estimate the
4 global number of host-associated bacterial operational taxonomic units (OTUs, a bacterial species analog),
5 by extrapolating an estimated number of 10.7 unique OTUs per insect species to an estimated 163.2 million
6 animal species worldwide, concluding that there exist billions of bacterial OTUs ($10.7 \times 163.2 \times 10^6$).
7 Their findings suggested the existence of an immense undiscovered bacterial diversity and triggered anew
8 discussions about the relative contribution of microbial taxa to the “Pie of Life” [3]. Here we explain why
9 [Larsen et al.](#)'s extrapolation is mathematically erroneous and likely led to a severe overestimate of global
10 host-associated bacterial OTU diversity, henceforth denoted N .

11 Based on pairwise comparisons of congeneric insect species, [Larsen et al.](#) calculated the average number
12 of OTUs per host species and the average number of OTUs unique to each host species when compared
13 to another randomly chosen congeneric species. For example, using previously published data [4], [Larsen](#)
14 [et al.](#) calculated that each species in the *Cephalotes* genus (turtle ants) has on average 19 bacterial OTUs
15 and that on average 9 OTUs were shared in any given comparison of two *Cephalotes* species, concluding
16 that each *Cephalotes* species has 10 unique bacterial OTUs. [Larsen et al.](#) obtained comparable results for
17 the genera *Drosophila* and *Nasonia*, and thus concluded that each insect species has on average $U = 10.7$
18 unique bacterial OTUs, and hence that $S = 163.2 \times 10^6$ animal species must host $\sim S \times U$ distinct OTUs.
19 [Larsen et al.](#)'s assumption that U generally applies to animal species other than insects is itself questionable,
20 and so is their relatively high estimate of 163 million animal species (e.g., Mora et al. [5] predicts only ~ 7.7
21 million animal species). Nevertheless, here we shall focus on a more fundamental (mathematical) issue in
22 [Larsen et al.](#)'s reasoning, by showing that the product $S \times U$ is under no circumstances a proper estimator
23 for N , for two important reasons.

24 First, [Larsen et al.](#) estimated U solely by comparing congeneric host species, and it is known that substantial
25 overlap exists between the microbiota of different host genera and even of distantly related animal taxa. This
26 overlap is partly due to host trait convergences [6–9], and partly due to the fact that some bacteria found

27 in hosts are not host-specific but merely taken up temporarily from the environment. For example, the gut
28 microbiota of fish overlap substantially with those of mammals and insects [7]. The extent of such overlaps
29 cannot possibly be inferred solely from congeneric host species comparisons, but it is needed for correctly
30 extrapolating to all animal taxa. To illustrate our point, in the (admittedly extreme) hypothetical scenario
31 where each animal family has the same set of bacterial OTUs, N would be equal to the OTU diversity within
32 a single animal family; purely congeneric host species comparisons could never rule out such a scenario.

33 Second, even if there was no overlap between the microbiota of distinct animal genera (an unrealistic
34 scenario, which would result in the highest possible N , given Larsen et al.’s U and S), the proper estimate
35 for N would be $G \times N_g$, where G is the number of animal genera and N_g is the number of bacterial OTUs
36 per animal genus. Assuming that Larsen et al.’s pairwise congeneric host species comparisons are indeed
37 representative, N_g may be estimated for a particular animal genus using the Chao2 incidence-based diversity
38 estimator [10, Eq. 11a]. The Chao2 estimator was originally designed for estimating total (observed +
39 unobserved) OTU diversity in a “region” or “community” (here, an animal genus) based on the number of
40 OTUs observed once or twice in a set of independent “sampling units” (here, two congeneric host species),
41 and is thus particularly suited for interpreting Larsen et al.’s data. In the latter case, the Chao2 estimator
42 takes the simple form

$$N_g = \frac{B^2}{B - U}, \quad (1)$$

43 where B is the average number of bacterial OTUs found in a single host species. Note that the above estima-
44 tor is mathematically analogous to mark-recapture approaches conventionally used to estimate the size of a
45 population [11], with B being analogous to the number of individuals marked in the first survey and $B - U$
46 being analogous to the number of marked individuals recaptured in a subsequent survey.

47 Taking Larsen et al.’s calculations for *Cephalotes* as an example ($B = 19$, $U = 10$), Eq. (1) would predict
48 only $N_g = 40.1$ OTUs for the entire *Cephalotes* genus. We point out that an approach analogous to Larsen
49 et al. [3], i.e. estimating N_g as $S_g \times U$ (where S_g is the number of *Cephalotes* species), would fail even for
50 a single genus. The reason is that the number of OTUs unique to a host species, when compared to a single
51 random congeneric species, is generally greater than the number of OTUs unique to a host species when
52 compared to all other congeneric species together. In other words, the number of OTUs that are truly unique
53 to a host species (i.e., not found in any other host), is generally much smaller than the average number of
54 OTUs unique to a host in pairwise comparisons. The *Cephalotes* genus contains about 130 described species
55 [3], and Mora et al. [5, Fig. S4] also estimate that there are about 100 species per animal genus. Assuming
56 ~ 100 species per animal genus, and assuming that Larsen et al.’s estimate of global animal species diversity
57 is correct, there are $G \sim S/100 \approx 1,632,000$ animal genera, and hence at most 65,443,200 bacterial OTUs
58 globally ($1,632,000 \times 40.1$).

59 In conclusion, even if Larsen et al.’s estimates of U and B for insect genera can be generalized to all
60 animal genera, and even if there was no overlap between the microbiota of distinct animal genera (evidently
61 a strongly unrealistic scenario [6, 7]), and even if Larsen et al.’s unusually high estimate of $S = 163.2 \times 10^6$
62 animal species was accurate, a mathematically correct use of Larsen et al.’s U , B and S would predict a global
63 bacterial diversity 25 times lower than claimed by Larsen et al. [3]. Using the animal diversity estimate by
64 Mora et al. [5], the estimate further reduces to 3,087,700 bacterial OTUs. Taking into account the substantial
65 overlap between animal genera would further reduce the estimated N . For example, even at a conservative
66 overlap of only 0.1% between any two randomly chosen genera, the global host-associated bacterial diversity
67 estimate would drop to 40,100 OTUs (applying the Chao2 estimator for $B = 40.1$ and $U = 0.999 \cdot B$).

References

- 68 1 Whitman WB, Coleman DC, Wiebe WJ. Prokaryotes: The unseen majority. *Proc Natl Acad Sci USA*.
69 1998; 95(12):6578–6583.
- 71 2 Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of the tree
72 of life. *Nat Microbiol*. 2016; 1:16048 EP. doi:<https://doi.org/10.1038/nmicrobiol.2016.48>.
- 73 3 Larsen BB, Miller EC, Rhodes MK, Wiens JJ. Inordinate fondness multiplied and redistributed:
74 the number of species on Earth and the new pie of life. *Q Rev Biol*. 2017; 92(3):229–265.
75 doi:<https://doi.org/10.1086/693564>.
- 76 4 Sanders JG, Powell S, Kronauer DJC, Vasconcelos HL, Frederickson ME, Pierce NE. Stability and
77 phylogenetic correlation in gut microbiota: lessons from ants and apes. *Mol Ecol*. 2014; 23(6):1268–
78 1283. doi:<https://doi.org/10.1111/mec.12611>.
- 79 5 Mora C, Tittensor DP, Adl S, Simpson AG, Worm B. How many species are there on Earth and in the
80 ocean? *PLoS Biol*. 2011; 9(8):e1001127.
- 81 6 Muegge BD, Kuczynski J, Knights D, Clemente JC, González A, Fontana L, et al. Diet drives conver-
82 gence in gut microbiome functions across mammalian phylogeny and within humans. *Science*. 2011;
83 332(6032):970–974.
- 84 7 Sullam KE, Essinger SD, Lozupone CA, O’connor MP, Rosen GL, Knight R, et al. Environmental and
85 ecological factors that shape the gut bacterial communities of fish: a meta-analysis. *Mol Ecol*. 2012;
86 21(13):3363–3378. doi:<https://doi.org/10.1111/j.1365-294X.2012.05552.x>.
- 87 8 Aylward FO, Suen G, Biedermann PHW, Adams AS, Scott JJ, Malfatti SA, et al. Convergent
88 bacterial microbiotas in the fungal agricultural systems of insects. *mBio*. 2014; 5(6):e02077–14.
89 doi:<https://doi.org/10.1128/mBio.02077-14>.
- 90 9 Salzman S, Whitaker M, Pierce NE. Cycad-feeding insects share a core gut microbiome. *Biol J Linn*
91 *Soc*. 2018; p. bly017. doi:<https://doi.org/10.1093/biolinnean/bly017>.
- 92 10 Chao A, Chiu CH. Species richness: estimation and comparison. *Wiley StatsRef: Statistics Reference*
93 *Online*. 2016; .
- 94 11 Krebs CJ. *Ecological Methodology*. Life Sciences. San Francisco, USA: Benjamin/Cummings; 1999.