

# Scaling laws predict global microbial diversity

Kenneth J. Locey<sup>1\*</sup>, Jay T. Lennon<sup>1\*</sup>

## Affiliations

<sup>1</sup> Department of Biology, Indiana University, Bloomington, IN, 47405, USA.

\*Correspondence to: ken@weecology.org; lennonj@indiana.edu.

**Abstract:** Scaling laws underpin unifying theories of biodiversity and are among the most predictively powerful relationships in biology. However, scaling laws developed for plants and animals often go untested or fail to hold for microorganisms. As a result, it is unclear whether scaling laws of biodiversity will span evolutionarily distant domains of life that encompass all modes of metabolism and scales of abundance. Using a global-scale compilation of ~35,000 sites and  $\sim 5.6 \cdot 10^6$  species, including the largest ever inventory of high-throughput molecular data and one of the largest compilations of plant and animal community data, we demonstrate similar rates of scaling in commonness and rarity across microorganisms and macroscopic plants and animals. We document a universal dominance scaling law that holds across 30 orders of magnitude, an unprecedented expanse that predicts the abundance of dominant ocean bacteria. In combining this scaling law with the lognormal model of biodiversity, we predict that Earth is home to upwards one trillion ( $10^{12}$ ) microbial species. Microbial biodiversity seems greater than ever anticipated yet predictable from the smallest to the largest microbiome.

22

The understanding of microbial biodiversity has rapidly transformed over the past decade. High throughput sequencing and bioinformatics have expanded the catalog of microbial taxa by orders of magnitude, while the unearthing of new phyla is reshaping the tree of life (1-3). At the same time, discoveries of novel forms of metabolism have provided insight into how microbes persist in virtually all aquatic, terrestrial, engineered, and host-associated ecosystems (4, 5). However, this period of discovery has uncovered few, if any, general rules for predicting microbial biodiversity at scales of abundance that characterize, for example, the  $\sim 10^{14}$  cells of bacteria that inhabit a single human or the  $\sim 10^{30}$  cells of bacteria and archaea estimated to inhabit Earth (6, 7). Such findings would aid the estimation of global species richness and reveal whether theories of biodiversity hold across all scales of abundance and whether so-called law-like patterns of biodiversity span the tree of life.

34

A primary goal of ecology and biodiversity theory is to predict diversity, commonness, and rarity across evolutionarily distant taxa and scales of space, time, and abundance (8-10). This goal can hardly be achieved without accounting for the most abundant, widespread, and metabolically, taxonomically, and functionally diverse organisms on Earth, i.e., microorganisms. Yet tests of biodiversity theory rarely include both microbial and macrobial datasets. At the same time, the study of microbial ecology has yet to uncover quantitative relationships that predict diversity, commonness, and rarity at the scale of host microbiomes and beyond. These unexplored opportunities leave the understanding of biodiversity limited to the most conspicuous taxa and largely unresolved for microorganisms. This lack of synthesis has also resulted in the independent study of two phenomena that likely represent a single universal pattern, i.e., highly uneven distributions of abundance that underpin biodiversity theory and that are ubiquitous

among communities of plants and animals (11), and the universal pattern of microbial  
commonness and rarity known as the microbial “rare biosphere” (12).

Scaling laws provide a promising path to the unified understanding and prediction of  
biodiversity. Also referred to as power-laws, the forms of these relationships,  $y \approx x^z$ , predict  
linear rates of change under logarithmic transformation, i.e.,  $\log(y) \approx z\log(x)$  and hence,  
proportional changes across orders of magnitude. Scaling laws reveal how physiological,  
ecological, and evolutionary constraints hold across genomes, cells, organisms, and communities  
of greatly varying size (13-15). Among the most widely known are the scaling of metabolic rate  
( $B$ ) with body size ( $M$ ),  $B = B_0 M^{3/4}$  (13) and the rate at which numbers of species ( $S$ ) scale with  
area ( $A$ ),  $S = cA^z$  (16). These scaling laws are predicted by powerful ecological theories, though  
evidence suggests that they fail for microorganisms (17-19). Beyond area and body size there is  
an equally general constraint on biodiversity, i.e., the number of individuals ( $N$ ). Often referred  
to as total abundance,  $N$  can range from less than 10 individuals in a given area to the nearly  $10^{30}$   
cells of bacteria and archaea on Earth (6, 7). This expanse outstrips the 22 orders of magnitude  
that separate the mass of a *Prochlorococcus* cell ( $3 \cdot 10^{-16}$  kg) from a blue whale ( $1.9 \cdot 10^5$  kg), and  
the 26 orders of magnitude that result from measuring Earth’s surface area at a spatial grain  
equivalent to bacteria ( $5.1 \cdot 10^{26} \mu\text{m}^2$ ).

Here, we consider whether  $N$  may be one of the most powerful constraints on  
commonness and rarity, and one of the most expansive variables across which aspects of  
biodiversity could scale. While  $N$  imposes an obvious constraint on the number of species (i.e.,  $S$   
 $\leq N$ ), empirical and theoretical studies suggest that  $S$  scales with  $N$  at a rate of 0.25 to 0.5 (i.e.,  $S$   
 $\approx N^z$ ,  $0.25 \leq z \leq 0.5$ ) (20-22). Importantly, this relationship pertains to samples from different  
systems and not to cumulative patterns, e.g., collector’s curves, which are based on resampling

(20-22). Recent studies have also shown that  $N$  constrains universal patterns of commonness and rarity by imposing a numerical constraint on how abundance varies among species, across space, and through time (23-24). Most notably, greater  $N$  leads to increasingly uneven distributions and greater rarity. Hence, we expect greater  $N$  to correspond to an increasingly uneven distribution among a greater number of species, an increasing portion of which should be rare. However, the strength of the relationships, whether they differ between microbes and macrobes, and whether they conform to scaling-laws across orders of magnitude are virtually unknown.

If aspects of diversity, commonness, and rarity scale with  $N$ , then local-to-global scale predictions of microbial biodiversity could be within reach. Likewise, if these relationships are similar for microbes and macrobes, then we may be closer to a unified understanding of biodiversity than previously thought. To answer these questions, we compiled the largest publicly available microbial and macrobial datasets, to date. These data include 20,376 sites of bacterial, archaeal, and microscopic fungal communities and 14,862 sites of tree, bird, and mammal communities. We focused on taxonomic aspects of biodiversity including species richness ( $S$ ), similarity in abundance among species (evenness), the concentration of  $N$  among relatively low-abundant species (rarity), and the number of individuals belonging to the most abundant species (absolute dominance,  $N_{max}$ ). We use the resulting relationships to predict  $N_{max}$  and  $S$  in large microbiomes, and to make empirically supported and theoretically underpinned estimates for the number of microbial species on Earth.

## Results and Discussion

As predicted, greater  $N$  led to an increase in species richness, dominance, and rarity, and a decrease in species evenness. Rarity, evenness, and dominance scaled across eight orders of magnitude in  $N$  at rates that differed little, if at all, between microbes and macrobes (Fig. 1). We

found that richness ( $S$ ) scaled at a greater rate for microbes ( $z = 0.38$ ) than macrobes ( $z = 0.24$ ), but still near the expected range of  $0.25 \leq z \leq 0.5$  (Fig. 1). However, for a given  $N$ , microbes had greater rarity, less evenness, and more species than macrobes (Fig. 1). As a result, microbes and macrobes are similar in how commonness and rarity scale with  $N$ , but differ in ways that support the exceptional nature of the microbial rare biosphere. The most unifying relationship we observed was a nearly isometric (i.e.,  $0.9 < z < 1.0$ ) scaling of dominance ( $N_{max}$ ). When extended to global scales, this dominance scaling law closely predicts the abundance of dominant ocean bacteria. Using the lognormal model of biodiversity theory, published estimates of global microbial  $N$ , and published and predicted values of  $N_{max}$ , we predict that Earth is occupied by  $10^{11}$  to  $10^{12}$  microbial species. This estimate is also supported by the scaling of  $S$  with  $N$ .

**Scaling relationships point to an exceptional rare biosphere.** Across microbial and macrobial communities, increasing  $N$  led to greater rarity, greater absolute dominance, less evenness, and greater species richness (Fig. 1, See SI Appendix, Figs. S5 to S9). Bootstrapped multiple regressions revealed that the significance of differences between microbes and macrobes with regard to rarity and evenness, were dependent on sample size. Larger samples suggested significant differences but were less likely to pass the assumptions of multiple regression (see Methods, See SI Appendix, Fig. S5). Though based on disparate types of data (i.e., counts of individual organisms vs. environmental molecular surveys), absolute dominance scaled at similar rates for microbes and macrobes (Fig. 1). Each relationship was best fit by a power-law as opposed to linear, exponential, or semi-log relationships (See SI Appendix, Table S1).

Since being first described nearly a decade ago (25), the rare biosphere has become an intensively studied pattern of microbial commonness and rarity (12). While its general form reiterates the ubiquitously uneven nature of ecological communities, our results suggest that

microbial communities are exceptional in degrees of rarity and unevenness. While artifacts sometimes associated with molecular surveys may inflate disparities in abundance or generate false singletons, our findings suggest that relationships of rarity, dominance, evenness, and richness were robust to the inclusion or exclusion of singletons and different percent cutoffs in sequence similarity (See SI Appendix, Figs. S8-S9). Naturally, the inclusion of unclassified sequences led to higher taxonomic richness. As a result of this large-scale comparison, we suggest that the rare biosphere is driven by the unique biology and ecology of microorganisms. Examples are the ability of small populations to persist in suboptimal environments via resilient life stages, the ability of microbes to disperse long distance and colonize new habitats, the capacity of microbes to finely partition niche axes, and the greater ability of asexual organisms to maintain small population sizes (12).

**Predicted scaling of species richness ( $S$ ).** Scaling exponents ( $z$ ) for the relationship of species richness ( $S$ ) to  $N$  fell near or within the predicted range (i.e.,  $0.25 < z < 0.5$ ) (20-22) (Fig. 1; Table 1). Despite variation in the relationship among datasets (SI Appendix, Fig. S7 A-I), the error structure across datasets was largely symmetrical (See SI Appendix, Fig. S5). Across datasets,  $z$ -varied more greatly for macrobes (0.07 to 1.23) than microbes (0.20 to 0.46), which more closely resembled the expected relationship (Table 1; See SI Appendix, Fig. S7). However, pooling all data to make use of the full range of  $N$  and to average out idiosyncrasies across datasets provided a stronger overall relationship, and produced an exponent ( $z = 0.51$ ), nearly identical to that observed in other empirical studies (20-22).

**An expansive dominance scaling law.** While greater  $N$  naturally leads to greater absolute dominance ( $N_{max}$ ) (26), this relationship is rarely explored and, to our knowledge, has not been

studied as a scaling law. We found that  $N_{max}$  scaled with  $N$  at similar and nearly isometric rates (i.e.  $0.9 < z < 1.0$ ) for microbes and macrobes across eight orders of magnitude (Fig.1;  $R^2 = 0.94$ ). Based on the strength of this result, we tested whether this scaling law holds at greater scales of  $N$ . We used published estimates for  $N$  and  $N_{max}$  from the human gut (27, 28), the cow rumen (29, 30), the global ocean (non-sediment), and Earth (6, 7, 31, 32). In each case, we found that  $N_{max}$  fell within the 95% prediction intervals of the dominance scaling law (Fig. 2). Though derived from datasets where  $N < 10^8$ , the dominance scaling law predicted the global abundance of some of the most abundant bacteria on Earth (SAR11, *Prochlorococcus marinus*) within an order of magnitude of prior estimates (31, 32). As a result, this dominance scaling law appears to span an unprecedented 30 orders of magnitude in  $N$ , extending to the upper limits of abundance in nature. The only other biological scaling law that approaches this expanse is the  $3/4$  power scaling of metabolic power to mass, which holds across 27 orders of magnitude (33).

**Predicting global microbial  $S$  using  $N$  and  $N_{max}$ .** Knowing the number of species on Earth is among the greatest challenges in biology (34-37). Historically, scientists have estimated global richness ( $S$ ) by extrapolating rarefaction curves and rates of accumulation, and often without including microorganisms (36, 37, 38). Though estimates of global microbial  $S$  exist, they range from  $10^4$  to  $10^9$  and rely on cultured organisms, precede large-scale sequencing projects, and are often based on the extrapolation of statistical estimators (e.g., rarefaction, Chao). These approaches also lack the theoretical underpinnings that distinguish extrapolations of statistical estimates from predictions of biodiversity theory. As an alternative approach to estimating  $S$ , we leveraged our scaling relationships with well-established biodiversity theory.

Based on the scaling of  $S$  with  $N$  (Fig. 1), we would expect a global microbial  $S$  of  $2.12 \pm 0.138 \cdot 10^{11}$  species. However, this risky type of exercise would extrapolate 26 orders of



magnitude beyond the available data (Fig. 2). Instead, we used the dominance scaling law and one of the most successful models of biodiversity (i.e., the lognormal distribution) to make a theoretically underpinned prediction of global microbial  $S$  (35, 39). The lognormal predicts that the distribution of abundance among species is approximately normal when species abundances are log-transformed (20). An extension of the central limit theorem, the lognormal arises from the multiplicative interactions of many random variables (20, 39). Though historically used to predict patterns of commonness and rarity, the lognormal was later derived to predict  $S$  using  $N$  and  $N_{max}$  (35). This led to predictions of  $S$  for habitats ranging in size from a milliliter of water to an entire lake, and speculations of  $S$  for the entire ocean.

To our knowledge, the lognormal is the only general biodiversity model that has been derived to predict  $S$  using only  $N$  and  $N_{max}$  as inputs. We used the lognormal to predict microbial  $S$  in two ways. First, we used published estimates of  $N$  and predicted the values of  $N_{max}$  via our dominance scaling law (Fig. 2). Second, we made predictions of  $S$  using published estimates of both  $N$  and  $N_{max}$ <sup>6,7,31,32</sup>. Assuming that global microbial  $N$  ranges from  $9.2 \cdot 10^{29}$  to  $3.2 \cdot 10^{30}$  (6, 7), the lognormal predicts  $3.23 \pm 0.227 \cdot 10^{12}$  species when  $N_{max}$  is predicted from the dominance scaling law (see Methods). However, using published estimates for  $N_{max}$  ranging from  $2.9 \cdot 10^{27}$  to  $2.4 \cdot 10^{28}$  (31, 32), the lognormal model predicts a value of global microbial  $S$  that is on the same order of magnitude as the richness-abundance scaling relationship, i.e.,  $3.9 \pm 0.054 \cdot 10^{11}$  species (Fig. 3). The general agreement between the lognormal model and the richness scaling relationship is encouraging given the magnitude of these predictions.

Our predictions of  $S$  for large microbiomes are among the most rigorous to date, resulting from intersections of empirical scaling, ecological theory, and the largest ever molecular surveys of microbial communities. However, several caveats should be considered. First, observed  $S$  for

the Earth Microbiome Project (EMP) differed greatly depending on whether we used closed or open reference data (see Methods), where  $S$  was  $\sim 6.9 \cdot 10^4$  and  $5.6 \cdot 10^6$ , respectively. In our main study, we used the closed reference data owing to the greater accuracy of that approach and because 42% of all taxa in the open reference EMP dataset were only detected twice or less. Consequently, choices such as how to assign OTUs and which primers or gene regions to use, need to be made cautiously and deliberately. Second, estimates of  $S$  will be much greater than observed when many species are detected only once or twice, as with the EMP. Statistical estimators of  $S$  such as rarefaction, jackknife, Chao, ACE, etc., are driven by singletons and doubletons (26). Third, it is difficult to estimate the portion of species missed when only a miniscule fraction of all individuals are sampled. For example, the intersection of the lognormal model and the richness scaling relationship suggests that  $S$  for an individual human gut could range from  $10^5$  to  $10^6$  species (Fig. 3). However,  $S$  of human gut samples is often on the order of  $10^3$ , while  $N$  is often less than  $10^6$ . Yet these are vanishingly small fractions of the gut microbiome, even when many samples are compiled together. For example, compiling all 4,303 samples from the Human Microbiome Project (HMP) dataset yields only  $2.180 \cdot 10^7$  reads; hardly sufficient for detecting  $10^5$  to  $10^6$  species among  $10^{14}$  cells. Consequently, detecting the true  $S$  of microbiomes with large  $N$  is a profound challenge that requires many large samples.

## Conclusion

We estimate that Earth is inhabited by  $10^{11}$  to  $10^{12}$  microbial species. This prediction is based on ecological theory reformulated for large-scale predictions, a new and perhaps most expansive ecological scaling-law, a richness scaling relationship with empirical and theoretical support, and the largest molecular surveys compiled to date. The profound magnitude of our prediction for Earth's microbial diversity stresses the need for continued investigation. We

expect the dominance scaling law we uncovered to be valuable in predicting richness, commonness, and rarity across all scales of abundance. To move forward, biologists will need to push beyond current computational limits and increase their investment in collaborative sampling efforts to catalog Earth's microbial diversity. For context,  $\sim 10^4$  species have been cultured, less than  $10^5$  species are represented by classified sequences, and the entirety of the Earth Microbiome Project has cataloged less than  $10^7$  species, 29% of which were only detected twice. Powerful relationships like those documented here and a greater unified study of commonness and rarity will greatly contribute to finding the potentially 99.999% of microbial taxa that remain undiscovered.

## Materials and Methods

**Data.** Our macrobial datasets comprised 14,862 different sites of mammal, tree, and bird communities. We used a compilation of data that included species abundance data for communities distributed across all continents except Antarctica (40). This compilation is based, in part, on five continental- to global-scale surveys: USGS North American Breeding Bird Survey (41) (2,769 sites), citizen science Christmas Bird Count (42) (1,412 sites), Forest Inventory Analysis (43) (10,356 sites), Alwyn Gentry's Forest Transect Data Set (44) (222 sites), and one global-scale data compilation, the Mammal Community Database (45) (103 sites). We limited our CBC dataset to sites where  $N$  was no greater than  $10^4$ , i.e., the reported maximum  $N$  for the BBS. Above that, estimates of  $N$  are not likely based on counts of individuals. No site is represented more than once in our data. Greater detail can be found elsewhere (appendix of 40).

We used 20,376 sites of communities of bacteria, archaea, and microscopic fungi. 14,615 of these were from the Earth Microbiome Project (EMP) (1) obtained on 22 August, 2014.

Sample processing, sequencing and amplicon data are standardized and performed by the EMP

and all are publicly available at [www.microbio.me/emp](http://www.microbio.me/emp). The EMP data consist of open and closed reference datasets. The QIIME tutorial ([http://qiime.org/tutorials/otu\\_picking.html](http://qiime.org/tutorials/otu_picking.html)) defines closed-reference as a classification scheme where any rRNA reads that do not hit a sequence in a reference collection are excluded from analysis. In contrast, open-reference refers to a scheme where reads that do not hit a reference collection are subsequently clustered de novo and represent unique but unclassified taxonomic units. Our main results are based on closed-reference data due to the greater accuracy of that approach and because 13% of all taxa in the open reference EMP dataset were only detected once while 29% are only detected twice.

We also used 4,303 sites from the Data Analysis and Coordination Center (DACC) for the National Institutes of Health (NIH) Common Fund supported Human Microbiome Project (HMP) (46). These data consisted of samples taken from 15 or 18 locations (including the skin, gut, vagina, and oral cavity) on each of 300 healthy individuals. The V3-V5 region of the 16S rRNA gene was sequenced for each sample. We excluded sites from pilot phases of the HMP as well as time-series data. See [http://hmpdacc.org/micro\\_analysis/microbiome\\_analyses.php](http://hmpdacc.org/micro_analysis/microbiome_analyses.php) for details on HMP sequencing and sampling protocols. We also included the 139 “prokaryote enriched” samples from 68 pelagic and mesopelagic locations, representing all major oceanic regions (except the Arctic), gathered by the Tara Oceans expedition (47).

We included 1,319 non-experimental PCR-targeted rRNA amplicon sequencing projects from the Argonne National Laboratory metagenomics server MG-RAST (48). Represented in this compilation were samples from arctic aquatic systems (130 sites; MG-RAST id: mgp138), hydrothermal vents (123 sites; MG-RAST id: mgp327) (49), freshwater lakes in China (187 sites; MG-RAST id: mgp2758) (50), arctic soils (44 sites; MG-RAST id: mgp69) (51), temperate soils (84 sites; MG-RAST id: mgp68) (52), bovine fecal samples (16 sites; MG-RAST id:

mgp14132), human gut microbiome samples not part of HMP (529 sites; MG-RAST id:

mgp401) (53), a global-scale dataset of indoor fungal systems (128 sites) (54), and freshwater, marine, and intertidal river sediments (34 sites; MG-RAST id: mgp1829). Using MG-RAST

allowed us to choose common parameter values for sequence similarity (i.e. 97% for species-level) and taxa assignment including a maximum e-value (probability of observing an equal or

better match in a database of a given size) of  $10^{-5}$ , a minimum alignment length of 50 base pairs, and minimum percent sequence similarities of 95, 97, and 99% to the closest reference sequence

in MG-RAST's M5 rRNA database (48-55). Below, we analyze MG-RAST datasets with respect to these cutoffs and reveal no significant effect on scaling relationships. Among the taxa not

included in our analyses are reptiles, amphibians, fish, large mammals, invertebrates, and protists. These taxa were absent because large datasets do not exist for their communities or

because redistribution rights could not be gained for publication.

**Quantifying dominance, evenness, rarity, and richness.** We calculated or estimated aspects of

diversity (dominance, evenness, rarity, richness) for each site in our data compilation. All

analyses can be reproduced or modified for further exploration by using code, data, and

following directions provided here: <https://github.com/LennonLab/ScalingMicroBiodiversity>.

*Rarity:* Here, rarity quantifies the concentration of species at low abundance (26). Our

primary rarity metric was the skewness of the frequency distribution of arithmetic abundance

classes ( $R_{skew}$ ), which are almost always right-skewed distributions (26). Due to the inability to

take the logarithm of a negative skew,  $R_{skew}$  was given a modulo transformation. The log-modulo transformation adds a value of one to each measure of skewness and converts negative values to

positive values, making them all positive and able to be log-transformed. We also quantified

rarity using log-transformed abundances ( $R_{log-skew}$ ) (26). We present results for  $R_{log-skew}$  in the

Supplement (See SI Appendix, Fig. S4). *Dominance*: Dominance refers to the abundance of the most abundant species, the simplest measures of which is the abundance of the most abundant species (absolute dominance;  $N_{max}$ ) (26). Relative dominance is also a common measure, and is known the Berger-Parker index ( $N_{max}/N = D_{BP}$ ). We focus on  $N_{max}$  in the main body because of the previously undocumented scaling with  $N$  and the ability to predict  $S$  using  $N$ ,  $N_{max}$ , and the lognormal model. We also calculated dominance as the sum of the relative abundance of the two most abundant taxa (i.e., McNaughton's dominance) and as Simpson's diversity, which is more accurately interpreted as an index of dominance (26). We present results for dominance metrics other than  $N_{max}$  in the Supplement (See SI Appendix, Fig. S3).

*Evenness*: Species evenness captures similarity in abundance among species (26, 56). We used five evenness metrics that perform well according to a series of statistical requirements (56), including lacking a strong bias towards very large or very small abundances, independence of richness ( $S$ ), and scaling between 0 (no evenness) and 1 (perfect evenness). These metrics included Smith and Wilson's indices ( $E_{var}$ ,  $E_Q$ ), Simpson's evenness ( $E_{1/D}$ ), Bulla's index ( $O$ ), and Camargo's index ( $E'$ ) (26, 56). We present results for  $E_{1/D}$  in the main results and results for other four metrics in the Supplement (See SI Appendix, Fig. S2). *Richness*: Richness ( $S$ ) is the number of species observed or estimated from a sample. Estimates of  $S$  are designed to account for rare species that go undetected in unbiased surveys (26). We present results for observed  $S$  in the main body along with results for six estimators of  $S$  (Chao1, ACE, jackknife, rarefaction, Margelef, McHennick) in the Supplement (See SI Appendix, Fig. S1).

**Approximating ranges of  $N_{max}$  for large microbiomes.** *Cow rumen*: The most dominant taxonomic unit (based on 97% sequence similarity in 16S rRNA reads) in the cow rumen is typically a member of the *Provatella* genus and has been reported to account for about 1.5 to 2.0

% of 16S rRNA gene reads in a sample (29, 30). Assuming there are about  $10^{15}$  microbial cells in the cow rumen (29, 30), and if these percentages are reflective of community wide relative dominance ( $D_{BP}$ ), then  $N_{max}$  of the cow rumen would be in range of  $1.5 \cdot 10^{14}$  and  $2 \cdot 10^{14}$ . *Human gut*: Deep sequencing of the human gut reveals that the most dominant taxon (based on 97% 16S rRNA sequence similarity) accounts for 10.6% to 12.2% of 16S rRNA gene reads in a sample (6, 28). Assuming these percentages are reflective of the microbiome, at large, and that there are about  $10^{14}$  microbial cells in the human gut (5, 28, 46), then  $N_{max}$  would be in range of  $1.06 \cdot 10^{13}$  and  $1.22 \cdot 10^{13}$ . *Global ocean (non-sediment) and Earth*: The most abundant microbial species on Earth has yet been determined. Perhaps, the best genus-level candidates (based on 97% 16S rRNA sequence similarity) are the marine picocyanobacteria *Synechococcus* and *Prochlorococcus* with estimated global abundances of  $7.0 \pm 0.3 \cdot 10^{26}$  and  $2.9 \pm 0.1 \cdot 10^{27}$ , respectively (32). Members of the SAR11 clade (i.e., Pelagibacterales), have an estimated global abundance of  $2.0 \cdot 10^{28}$  and may also be candidate for the most abundant microorganisms on Earth (31). We used SAR11 as the upper limit for the most dominant microbial species on Earth, i.e., the most abundant species cannot be more abundant than the most dominant order-level clade. We used  $6.7 \cdot 10^{26}$  to  $3.0 \cdot 10^{27}$ , as the range for  $N_{max}$  of the non-sediment global ocean, and used  $2.9 \cdot 10^{27}$  to  $2.0 \cdot 10^{28}$ , as the range for  $N_{max}$  of Earth. We used the range of  $3.6 \cdot 10^{28}$  to  $1.2 \cdot 10^{29}$  as lower and upper range for the number of microbial cells in the open ocean (7) and  $9.2 \cdot 10^{29}$  to  $3.2 \cdot 10^{30}$  (6) as the lower and upper range for the number of microbial cells on Earth.

**Predictions of  $S$  for large microbiomes and Earth.** We used the method of Curtis et al. (35) to predict global microbial richness ( $S$ ) using the lognormal species abundance model of (39). Curtis et al. (35) used the lognormal to estimate microbial  $S$  in a gram of soil, a milliliter of water, an entire lake, and then speculate on what  $S$  may be for a ton of soil (many small

ecosystems) and the entire ocean (many large ecosystems). The lognormal prediction of  $S$  is based on the ratio of total abundance ( $N$ ) to the abundance of the most abundant species ( $N_{\max}$ ), and the assumption that the rarest species is a singleton,  $N_{\min} = 1$ . Equation 1 from Curtis et al. (35): According to the log-normal model, in a communities with  $S(N)$  species, the number of taxa that contain  $N$  individuals is:

$$S(N) = \frac{Sa}{\sqrt{\pi}} \exp \left\{ - (a \log_2(\frac{N}{N_0}))^2 \right\}$$

where  $a$  is an inverse measure of the width of the distribution whose standard deviation is  $\sigma^2$  ( $a = [2 \ln 2 \sigma^2]^{-1/2}$ ) and  $N_0$  is the most common (i.e., modal) abundance class. Equation 3 from Curtis et al (35): If it is assumed that the log-normal species abundance curve is not truncated and therefore is symmetric about  $N_0$ , then it can be shown that,

$$N_{\min} = \frac{N_0^2}{N_{\max}}$$

The second method for estimating the spread of the lognormal distribution,  $a$ , is by knowing or assuming  $N_{\min}$ . By using Equations 1 and 3 and the assumption that  $S(N_{\min}) = 1$ ,  $S$  can be expressed in terms of  $a$ ,  $N_{\min}$ , and  $N_{\max}$ . Curtis et al. (35) reason that  $S$  will not be sensitive to small deviations from the  $N_{\min} = 1$  assumption and hence, that knowledge of  $N_{\min}$ ,  $N_{\max}$ , and  $N$  allows Equation 11 to be solved numerically for Preston's  $a$  parameter and, subsequently, for  $S$  to be predicted using their Equation 10:

$$S(N) = \frac{\sqrt{\pi}}{a} \exp \left\{ (a \log_2(\sqrt{\frac{N_{\max}}{N_{\min}}}))^2 \right\}$$

The authors show that the above equation can be used to rewrite their Equation 5 as Equation 11:

$$N_T = \frac{\sqrt{\pi N_{\min} N_{\max}}}{2a} \exp \left\{ (a \log_2(\sqrt{\frac{N_{\max}}{N_{\min}}}))^2 \right\} \exp \left\{ \left( \frac{\ln(2)}{2a} \right)^2 \right\}$$



$$\cdot \left[ \operatorname{erf} \left( a \log_2 \left( \sqrt{\frac{N_{\max}}{N_{\min}}} - \frac{\ln(2)}{2a} \right) \right) + \operatorname{erf} \left( a \log_2 \left( \sqrt{\frac{N_{\max}}{N_{\min}}} + \frac{\ln(2)}{2a} \right) \right) \right]$$

Equation 11 can be numerically solved for  $a$ , which is then used in Equation 10 to solve for  $S$ .

We coded Equations 1, 3, 10, and 11 into a Python script that can be used to recreate the results

of Curtis et al. (35) under the functions “alpha2” to derive  $a$  and “s2” to estimate  $S$ . In predicting

$S$ , we accounted for variability in  $N$  and  $N_{\max}$  by randomly sampling within their published

ranges (See SI Appendix, Fig. S14). This allowed us to generate means and standard errors,

which are often lacking from large-scale predictions of  $S$ .

**Resampling and dependence on sample size and sequence similarity.** We examined

relationships of rarity, evenness, dominance, and richness to the number of individual organisms

or gene reads ( $N$ ) using 10,000 bootstrapped multiple regressions based on stratified random

sampling of microbial and macrobial datasets. We examined the sensitivity of our results to

sampling strategy, sample size, particular datasets, and the microbe/macrobe dummy variable,

results of which can be found in the Supplement (See SI Appendix, Figs. S5 to S13). To use

equal numbers of sites for macrobes and microbes in each multiple regression analysis, we used

100 sites from each macrobial dataset for a total of 500 randomly chosen sites. To obtain 500

sites from our microbial data, we used 50 randomly chosen sites from each microbial dataset

having more than 100 sites, and 20 randomly chosen sites from smaller datasets. We used the

mean values of coefficients and intercepts (accounting for whether differences between microbes

and macrobes were significant at  $p < 0.05$ ,  $\alpha = 0.05$ ) from multiple regressions to estimate

the relationships of rarity, evenness, dominance, and richness to  $N$ . We examined whether

scaling relationships for microbial data were sensitive to the percent cutoff in rRNA sequence

similarity, which is used to bin taxa into species-level units. These analyses were restricted to

datasets obtained from MG-RAST but reveal no statistical differences due to whether sequences were binned based on 95, 97, and 99% similarity.

**Power-law behavior vs. other functional forms.** We tested whether relationship of richness, evenness, rarity, and dominance were better fit by a power-law (log-log) than by linear, exponential, and semi-log relationships (See SI Appendix, Table. S1). The power-law model explained substantially greater variance or, in the one case where it was nearly tied in explanatory power, had a substantially lower AIC and BIC values than other models.

**Available code and data.** We used freely available open source computing and version control tools. Analyses and figures can be automatically regenerated using Python scripts and data files in the public GitHub repository <https://github.com/LennonLab/ScalingMicroBiodiversity>.

Analyses can be recreated step-by-step using the directions given in the repository.

**Author contributions.** K.J.L and J.T.L designed the research; K.J.L. performed the research; K.J.L. acquired and analyzed datasets; and K.J.L. and J.T.L. wrote the paper.

## Acknowledgments

We thank M. Muscarella, W. Shoemaker, N. Wisnoski, X. Xiao, S. Gibbons, E. Hall, and E. P. White for critical reviews. We thank the individuals involved in collecting and providing the data used in this paper including the essential citizen scientists who collect the Breeding Bird Survey and Christmas Bird Count data; researchers who collected, sequenced, and provided metagenomic data on MG-RAST as well as the individuals who maintain and provide the MG-RAST service; the Audubon Society; the US Forest Service; the Missouri Botanical Garden; and Alwyn H. Gentry. This work was supported by a National Science Foundation Dimensions of Biodiversity Grant (#1442246) and the U.S. Army Research Office (W911NF-14-1-0411).

## References

1. J. A. Gilbert, J. K. Jansson, R. Knight. The Earth Microbiome project: successes and aspirations. *BMC biology* **12**, 69 (2014).
2. C. T. Brown, *et al.* Unusual biology across a group comprising than 15% of domain Bacteria. *Nature* **523**, 208–211 (2015).
3. Luef, B., *et al.* (2015) Diverse uncultivated ultra-small bacterial cells in groundwater. *Nature communications* 6. doi:10.1038/ncomms7372.
4. Venter, J. Craig, *et al.* "Environmental genome shotgun sequencing of the Sargasso Sea." *science* 304.5667 (2004): 66-74.
5. Gill, Steven R., *et al.* (2006). Metagenomic analysis of the human distal gut microbiome. *Science*, 312: 1355-1359.
6. J. Kallmeyer, R. Pockalny, R. R. Adhikari, D. C. Smith, S. D'Hondt. Global distribution of microbial abundance and biomass in subseafloor sediment. *Proc Natl Acad Sci USA* **109**, 16213–16216 (2012).
7. W. B. Whitman, D. C. Coleman, W. J. Wiebe. Prokaryotes: The unseen majority. *Proc Natl Acad Sci USA* **95**, 6578–6583 (1998).
8. Hubbell, S. P. (2001). The Unified Neutral Theory of Biodiversity and Biogeography. Princeton, New Jersey, USA: Princeton University Press.
9. McGill, B. J. (2010). Towards a unification of unified theories of biodiversity. *Ecology Letters*, 13(5), 627–642.
10. Harte, J. (2011). Maximum Entropy and Ecology. New York, New York, USA: Oxford University Press.

- 408 11. B. J. McGill, *et al.* Species abundance distributions: moving beyond single prediction  
theories to an integration within an ecological framework. *Eco Lett.* **10**, 995–1015 (2007).
- 410 12. A. Reid, M. Buckley. (2011). The Rare Biosphere: A report from the American Academy of  
Microbiology. (Washington, DC: American Academy of Microbiology, 2011)
- 412 13. J. H. Brown, J. F. Gillooly, A. P. Allen, V. M. Savage, G. B. West. Toward a metabolic  
theory of ecology. *Ecology* **85**, 1771–1789 (2004).
- 414 14. M. Lynch, J. S. Conery. The origins of genome complexity. *Science* **302**, 1401–1404 (2007).
15. M. E. Richie, H. Olff. Spatial scaling laws yield a synthetic theory of biodiversity. *Nature*  
416 **400**, 557–560 (1999).
16. Lomolino, M. V. (2000). Ecology's most general, yet protean 1 pattern: the species-area  
418 relationship." *Journal of Biogeography* 27: 17-26.
17. J. P. DeLong, J. Okie, M. E. Moses, R. M. Sibly, J. H. Brown. Shifts in metabolic scaling,  
420 production, and efficiency across major evolutionary transition of life. *Proc Natl Acad  
Sci USA* **107**, 12941–12945 (2010).
- 422 18. M. C. Horner-Devine, M. Lage, J. B. Hughes, B. J. M. Bohannan. A taxa-area relationship  
for bacteria. *Nature* **432**, 750–753 (2004).
- 424 19. J. L. Green, J. L., *et al.* Spatial scaling of microbial eukaryote diversity. *Nature* **432**, 747-750  
(2004).
- 426 20. May, R.M. (1975). Patterns of species abundance and diversity. In M.L. Cody and J.M.  
Diamond, eds., *Ecology and Evolution of Communities*, pp. 81-120. Cambridge: Harvard  
428 University Press, Belknap Press.
21. May, R.M. (1978) The dynamics and diversity of insect faunas. In L. A. Mount and N.  
430 Waloff, eds., *Diversity of Insect Faunas*, pp. 188-204. Oxford, England: Blackwell.

22. Siemann, E., Tilman D., Haarstad J. (1996). Insect species diversity, abundance and body  
size relationships." *Nature* 380: 704-706.
23. Locey, K. J., & White, E. P. (2013). How species richness and total abundance constrain the  
distribution of abundance. *Ecology Letters*, 16(9), 1177–1185.
24. Xiao, X., Locey, K. J., & White, E. P. (2015). A Process-Independent Explanation for the  
General Form of Taylor’s Law. *The American Naturalist*. 186:E51-E60.
25. Sogin, Mitchell L., et al. "Microbial diversity in the deep sea and the underexplored “rare  
biosphere”." *Proceedings of the National Academy of Sciences* 103.32 (2006): 12115-  
12120.
26. Magurran, A. E., McGill B. J., eds. (2011). *Biological diversity: frontiers in measurement  
and assessment*. Vol. 12. Oxford: Oxford University Press.
27. R. D. Berg. The indigenous gastrointestinal microflora. *Trends Microbiol.* **4**, 430–435  
(1996).
28. L. Dethlefsen, S. Huse, M. Sogin, D. A. Relman. The Pervasive Effects of an Antibiotic  
on the Human Gut Microbiota, as Revealed by Deep 16S rRNA Sequencing. *PLoS Biol.*  
**6**, e280 (2008).
29. E. Jami, I. Mizrahi. Composition and similarity of bovine rumen microbiota across  
individual animals. *PLoS ONE* **7**, e33306 (2012).
30. D. M. Stevenson, P. J. Weimer. Dominance of *Prevotella* and low abundance of classical  
ruminal bacterial species in the bovine rumen revealed by relative quantification real-  
time PCR. *Appl Microbiol Biotechnol.* **75**, 165–174 (2007).
31. R. M. Morris, *et al.* SAR11 clade dominates ocean surface bacterioplankton communities.  
*Nature* **420**, 806–809 (2002).

- 454 32. P. Flombaum, *et al.* Present and future global distributions of the marine cyanobacteria  
*Prochlorococcus* and *Synechococcus*. *Proc Natl Acad Sci USA* **110**, 9824–9829 (2013).
- 456 33. West, G. B., Woodruff, W. H., & Brown, J. H. (2002). Allometric scaling of metabolic rate  
from molecules and mitochondria to cells and mammals. *Proceedings of the National*  
458 *Academy of Sciences*, 99(suppl 1), 2473-2478.
34. R. M. May. How many species are there on Earth? *Science* **241**, 1441–1449 (1988)
- 460 35. T. P. Curtis, W. T. Sloan, J. W. Scannell. Estimating prokaryotic diversity and its limits.  
*Proc Natl Acad Sci USA* **99**, 10494–104999 (2002).
- 462 36. C. Mora, D. P. Tittensor, S. Adl, A. G. B. Simpson, B. Worm. How many species are there  
on Earth and in the ocean? *PLoS Biol.* **9**, e1001127 (2011).
- 464 37. N. E. Stork, J. McBroom, C. Gely, A. J. Hamilton. New approaches narrow global species  
estimates for beetles, insects, and terrestrial arthropods. *Proc Natl Acad Sci USA* **112**,  
466 7519–7523 (2015).
38. P. D. Schloss, J. Handelsman. Status of the microbial census. *Microbiol Mol Biol R.* **68**,  
468 686–691 (2004).
39. F. W. Preston. The commonness, and rarity, of species. *Ecology* **29**, 254–283 (1948).
- 470 40. E. P. White, K. M. Thibault, X. Xiao. Characterizing species abundance distributions across  
taxa and ecosystems using a simple maximum entropy model. *Ecology* **93**, 1772–1778  
472 (2012).
41. J. R. Sauer, *et al.* The North American Breeding Bird Survey 1966–2009. Version 3.23.2011.  
474 (USGS Patuxent Wildlife Research Center, Laurel, MD., 2011).
42. National Audubon Society. The Christmas Bird Count historical results. Available at:  
476 <http://www.audubon.org/bird/cbc>. (2002).

43. U.S. Department of Agriculture. Forest inventory and analysis national core field guide  
 478 (Phase 2 and 3), version 4.0. U.S. Department of Agriculture Forest Service, Forest  
 Inventory and Analysis, Washington, DC. (2010).

44. O. Phillips, J. S. Miller. *Global patterns of plant diversity: Alwyn H. Gentry's forest  
 480 transect data set.* (Missouri Botanical Garden Press, 2002).

45. K. M. Thibault, S. R. Supp, M. Giffin, E. P. White, S. K. M. Ernest. Species composition  
 482 and abundance of mammalian communities. *Ecology* **92**, 2316. (2011).

46. P. J. Turnbaugh, *et al.* The human microbiome project. *Nature* **449**, 804–810 (2007)

47. S. Sunagawa, *et al.* Structure and function of the global ocean microbiome. *Science* **348**,  
 486 1261359 (2015)

48. F. Meyer, *et al.* The metagenomics RAST server - a public resource for the automatic  
 488 phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* **9**, 386  
 (2008).

49. G. E. Flores, *et al.* Microbial community structure of hydrothermal deposits from  
 490 geochemically different vent fields along the Mid-Atlantic Ridge. *Environ. Microbiol.* **13**,  
 492 2158–2171 (2011).

50. J. Wang, *et al.* Phylogenetic beta diversity in bacterial assemblages across ecosystems:  
 494 deterministic versus stochastic processes. *ISME J.* **7**, 1310-1321 (2014).

51. H. Chu, *et al.* Soil bacterial diversity in the Arctic is not fundamentally different from that  
 496 found in other biomes. *Environ. Microbiol.* **12**, 2998–3006 (2010).

52. N. Fierer, *et al.* Comparative metagenomic, phylogenetic and physiological analyses of soil  
 498 microbial communities across nitrogen gradients. *ISME J.* **6**, 1007–1017 (2012)

53. T. Yatsunenko, *et al.* Human gut microbiome viewed across age and geography. *Nature* **486**,  
222-227 (2012).

54. A. S. Amend, K. A. Seifert, R. Samson, T. D. Bruns. Indoor fungal composition is  
geographically patterned and more diverse in temperate zones than in the tropics. *Proc.*  
*Natl Acad. Sci. USA* **107**, 13748–13753 (2010).

55. J. Goris, *et al.* DNA–DNA hybridization values and their relationship to whole-genome  
sequence similarities. *Int J Syst Evol Micr* **57**, 81–91 (2007).

56. B. Smith, J. B. Wilson. A consumer’s guide to evenness indices. *Oikos* **76**, 70–82 (1996).

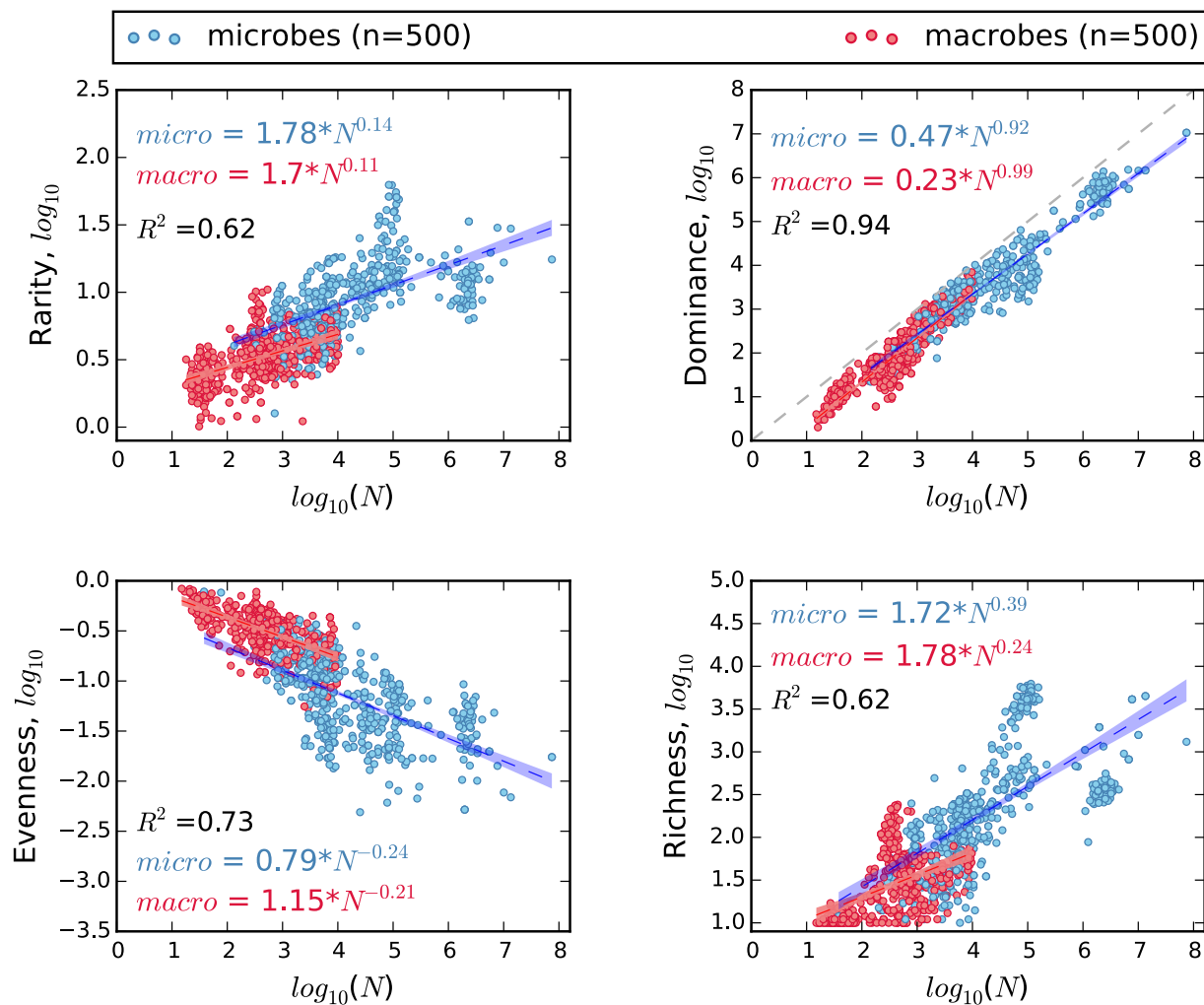


## Figure Legends

**Fig. 1.** Microbial communities (blue dots) and communities of macroscopic plants and animals (red dots) are similar in the rates at which rarity, absolute dominance, and species evenness scale with the number of individuals or genes reads ( $N$ ). However, for a given  $N$ , microbial communities have greater rarity, less evenness, and greater richness than those of macroorganisms. Coefficients and exponents of scaling equations are mean values from 10,000 bootstrapped multiple regressions, with each regression based on 500 microbial and 500 macrobial communities chosen by stratified random sampling. Each scatter plots represent a single random sample; hulls are 95% confidence intervals.

**Fig. 2.** The dominance-abundance scaling law (dashed red line) predicts the abundance of the most abundant microbial taxa ( $N_{max}$ ) up to global scales. The pink hull is the 95% prediction interval for the regression based on 3,000 sites chosen via stratified random sampling (red heat map) from our microbial data compilation. Gray cross-hairs are ranges of published estimates of  $N$  and  $N_{max}$  for large microbiomes including Earth<sup>6,7,31,32</sup> (see Materials and Methods: Approximating ranges of  $N_{max}$  for large microbiomes). The light gray dashed line is the 1:1 relationship. The scaling equation and  $r^2$  only pertain to the scatter plot data.

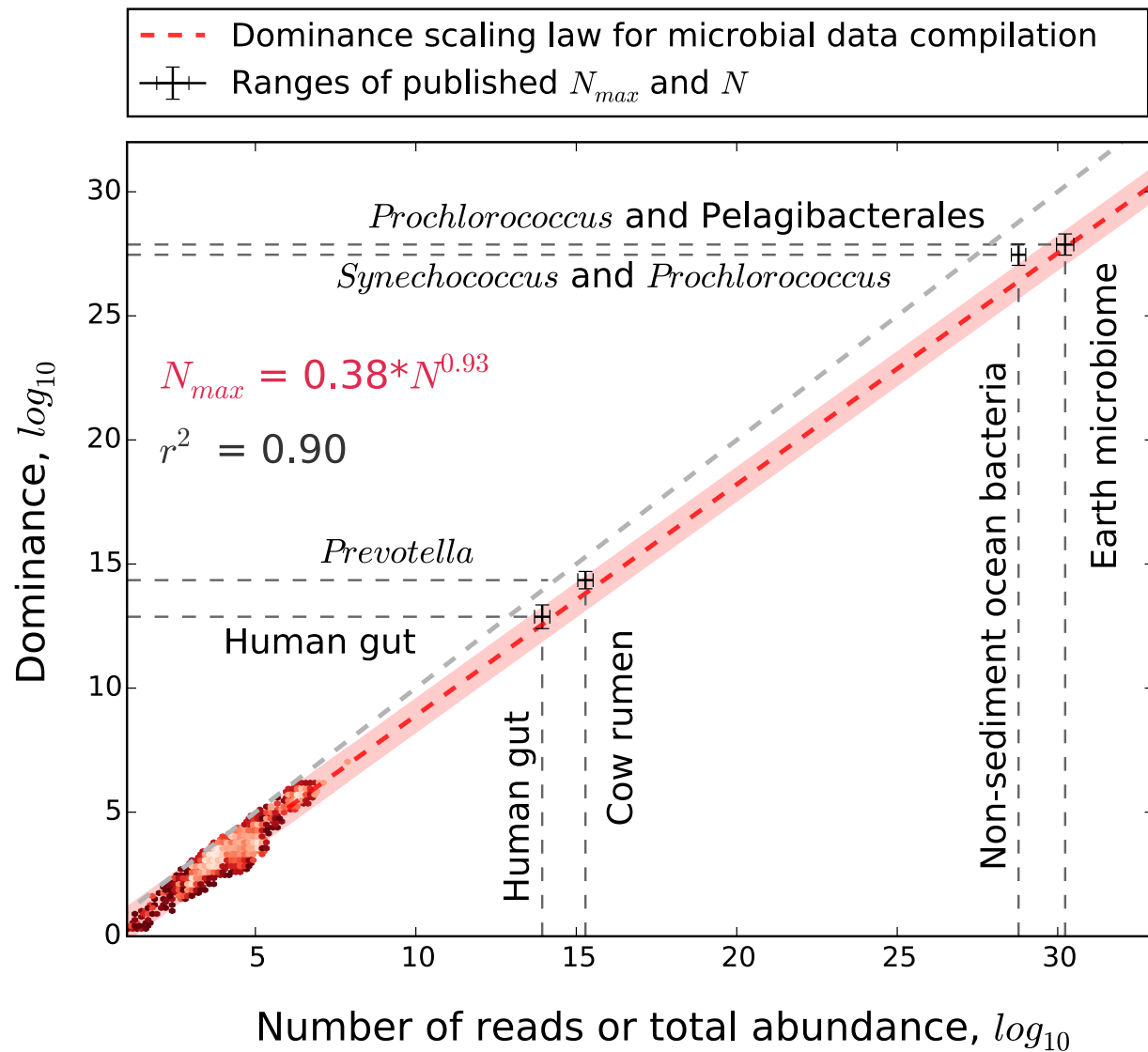
**Fig. 3.** The microbial richness-abundance scaling relationship (dashed red line) supports values of  $S$  predicted from the lognormal model using the published range of  $N$  and  $N_{max}$  (grey dots), as well as ranges of  $N_{max}$  predicted from the dominance scaling law (blue dots). The pink hull is the 95% prediction interval for the regression based on 3,000 sites chosen via stratified random sampling (red scatter plot). The scaling equation and r-square value are based solely on the red scatter plot data. Standard errors around predicted  $S$  are too small to illustrate.



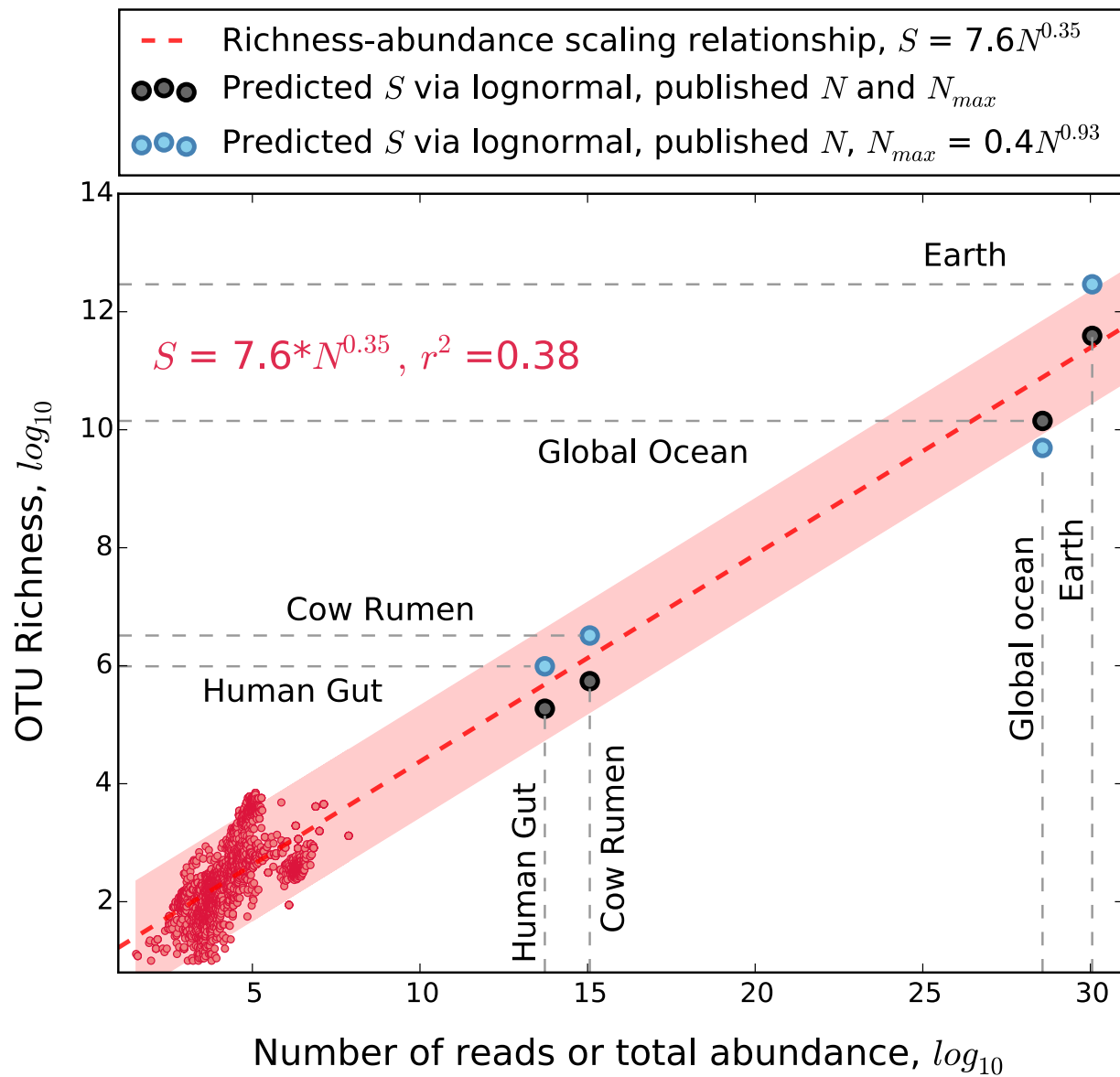
544

546

548



550



552

554

**Table 1.** Scaling relationships across datasets. with scaling exponents in bold and r-squared values in parentheses. Datasets are the Earth Microbiome Project (EMP), Human Microbiome Project (HMP), MG-RAST rRNA amplicon projects (MG-RAST), Tara Oceans Expedition (TARA), North American Breeding Bird Survey (BBS), Christmas Bird Count (CBC), Forest Inventory and Analysis (FIA), Gentry tree transects (GENTRY), and Mammal Community Database (MCDB). TARA was the only dataset where  $N$  ranged over less than an order of magnitude, leading results for TARA to be inconclusive. For most datasets,  $N_{max}$  scaled almost isometrically with  $N$ . For all datasets except TARA, evenness decreased with  $N$  while rarity increased. For birds and all microbe datasets,  $S$  scaled near the predicted range of 0.25 to 0.5.

Dataset	Rarity	Dominance	Evenness	Richness
EMP (n = 14,615)	<b>0.2</b> (0.30)	<b>1.01</b> (0.67)	<b>-0.44</b> (0.42)	<b>0.46</b> (0.42)
MG-RAST (n = 1,283)	<b>0.06</b> (0.20)	<b>0.98</b> (0.97)	<b>-0.17</b> (0.32)	<b>0.20</b> (0.45)
HMP (n = 4,303)	<b>0.14</b> (0.14)	<b>1.02</b> (0.70)	<b>-0.33</b> (0.18)	<b>0.29</b> (0.13)
TARA (n = 139)	<b>-0.26</b> (0.02)	<b>1.02</b> (0.13)	<b>0.06</b> (0.00)	<b>0.29</b> (0.13)
BBS (n = 2,769)	<b>0.16</b> (0.086)	<b>1.0</b> (0.54)	<b>-0.32</b> (0.22)	<b>0.32</b> (0.19)
CBC (n = 1,412)	<b>0.16</b> (0.39)	<b>1.07</b> (0.90)	<b>-0.35</b> (0.44)	<b>0.22</b> (0.48)
FIA (n = 10,355)	<b>0.07</b> (0.01)	<b>1.34</b> (0.68)	<b>-0.45</b> (0.27)	<b>0.07</b> (0.02)
GENTRY (n = 222)	<b>0.46</b> (0.27)	<b>0.29</b> (0.038)	<b>-0.19</b> (0.05)	<b>1.24</b> (0.46)
MCDB (n = 103)	<b>0.07</b> (0.07)	<b>1.07</b> (0.91)	<b>-0.16</b> (0.20)	<b>0.09</b> (0.19)