

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

Supplementary Information for  
**A General Framework for Quantitatively Assessing Ecological Stochasticity**

Daliang Ning, Ye Deng, James M. Tiedje, and Jizhong Zhou

Jizhong Zhou  
Email: [jzhou@ou.edu](mailto:jzhou@ou.edu)

**This PDF file includes:**

- Supplementary Text
- Figs. S1 to S6
- Tables S1 to S4
- Supplementary References for SI reference citations

## 26 Supplementary Text

### 27 A. Community similarity/dissimilarity metrics

28 Various similarity/dissimilarity metrics have been applied in ecological research to measure  $\beta$  diversity,  
29 either by taxonomic or phylogenetic  $\beta$ -diversity metrics. Many metrics have both incidence-based  
30 (qualitative) and abundance-based (quantitative) formats. Table S3 summarizes commonly used taxonomic  
31 (1-14) and phylogenetic (15-20) similarity/dissimilarity metrics. Please see Parks and Beiko's paper (17)  
32 for more phylogenetic  $\beta$ -diversity metrics (39 different indexes). The taxonomic diversity metrics used in  
33 our study were calculated using function "vegdist" and/or "designdist" in R package "vegan" (21).

34 In principle, our method requires the value of the metric ranging from 0 to 1, and the complementarity  
35 of similarity (C) and dissimilarity (D) indexes ( $C = 1 - D$ ). Many taxonomic and phylogenetic metrics  
36 satisfy this requirement, such as Jaccard, Ružička, Sørensen, Bray-Curtis, Kulczynski, Gower, Canberra,  
37 Morisita-Horn, unweighted Unifrac, and Phylosor, which can be directly implemented with our method.  
38 However, quite a few metrics do not have fixed upper limit and/or not have clear defined similarity measure,  
39 such as Euclidean, Manhattan, Binomial, Cao, modified Gower,  $\beta$ MNTD,  $\beta$ MPPD, and most of weighted  
40 phylogenetic dissimilarity metrics. These metrics need to be standardized to meet the requirement before  
41 being applied to our method. Inspired by what Cao et al (14) did to define similarity from Cao dissimilarity  
42 index, we proposed a general method as follows.

$$D_{ij} = \frac{D'_{ij} - \min\{D'\}}{\max\{D'\} - \min\{D'\}} = \frac{D'_{ij}}{D'_{max}} \quad \text{Eq. S1}$$

43 where  $\min\{D'\} = 0$ , considering the lower limit of dissimilarity metrics is always zero.

$$D'_{max} = \begin{cases} D'_{up} & \text{if upper limit is fixed} \\ \max_{ij}\{D'_{ij}, G'_{maxij}\} & \text{if upper limit is not fixed} \end{cases} \quad \text{Eq. S2}$$

$$C_{ij} = 1 - D_{ij} = 1 - \frac{D'_{ij}}{D'_{max}} \quad \text{Eq. S3}$$

$$G_{ij} = \frac{G'_{ij}}{D'_{max}} \quad \text{Eq. S4}$$

$$E_{ij} = 1 - G_{ij} = 1 - \frac{G'_{ij}}{D'_{max}} \quad \text{Eq. S5}$$

- 44  $D_{ij}$  Standardized dissimilarity between community  $i$  and  $j$ .  
45  $D'_{ij}$  Unstandardized (original) dissimilarity between community  $i$  and  $j$ .  
46  $C_{ij}$  Standardized similarity between community  $i$  and  $j$ .  
47  $D'_{max}$  Probable maximum unstandardized dissimilarity.  
48  $D'_{up}$  Defined upper limit of original dissimilarity, e.g. for Bray-Curtis,  $D'_{up} = 1$ .  
49  $G_{ij}$  Standardized null (randomly expected) dissimilarity between community  $i$  and  $j$ .  
50  $G'_{ij}$  Unstandardized null dissimilarity calculated between community  $i$  and  $j$ .

- 51  $E_{ij}$  Standardized null similarity calculated between community  $i$  and  $j$ .  
 52  $E'_{ij}$  Unstandardized null similarity calculated between community  $i$  and  $j$ .  
 53  $G'_{max\ ij}$  The estimated maximum value of null dissimilarity between community  $i$  and  $j$ . It is calculated as  
 54 3 times of smoothing bandwidth beyond the maximum simulated value of the dissimilarity  $D'_{ij}$ ,  
 55 using function “density” with Gaussian model in R package “stats”.  
 56

57 Our method described in the main text is based on similarity and dissimilarity metrics ranging from 0  
 58 to 1 or standardized according to above equations. If the upper limit of dissimilarity is not fixed,  $D'_{max}$  is  
 59 still not a fixed value but depends on community data matrix and null model algorithms, which leads to the  
 60 uncertainty of metrics standardization. Therefore, the original metrics ranging from 0 to 1 (i.e. there is no  
 61 need for standardization) are preferred if they have the same performance in terms of accuracy and precision  
 62 as standardized metrics.  
 63

64 We tested thirteen incidence-based metrics, which can be classified into three major categories:

- 65 i. Unique-ratio metrics: the dissimilarity is measured by the ratio of unique taxa (i.e. the taxa only  
 66 observed in one of the two samples). Jaccard is the unique taxa number divided by total observed  
 67 taxa number in two samples, which is exactly the same as incidence-based Canberra, modified  
 68 Gower (mGower), and modified Manhattan (mManhattan), and will be incidence-based Cao if  
 69 multiplied with a constant. Sørensen is the unique taxa number divided by the sum of observed  
 70 taxa number in two samples, which is exactly the same as incidence-based Morisita-Horn.  
 71 Kulczynski is the mean of unique taxa percentage in each sample, and Gower is the unique taxa  
 72 number divided by observed taxa number across all samples. Thus, these nine metrics can be  
 73 classified into the same type named unique-ratio metrics.
- 74 ii. Unique-number metrics: the dissimilarity is directly measured by the number of unique taxa in two  
 75 samples. Incidence-based Manhattan is defined as this, and incidence-based Binomial is Manhattan  
 76 multiplied with a constant.
- 77 iii. Squared-root metrics: the metrics are calculated from squared root of unique taxa number. This  
 78 type includes incidence-based Euclidean and modified Euclidean (mEuclidean) which is Euclidean  
 79 divided by total observed taxa number in the two samples.  
 80

81 We also tested fifteen abundance-based metrics, which can be divided into four major groups as below:

- 82 i. Relative-difference metrics: the abundance difference (or the relatively smaller abundance) of each  
 83 taxon between two samples is divided by the abundances of the taxon in the samples before or after  
 84 summed up. Ružička, Bray-Curtis, and Kulczynski definitely belong to this type. In Chao’s formula,  
 85 the total number of individuals in the taxa shared by the two samples ( $C_i$ ) is divided by the total  
 86 number of individuals in a sample (i.e. abundance sum in a sample), thus Chao is also classified  
 87 into this type.
- 88 ii. Average-relative-difference metrics: the sum of relative difference between two samples (or other  
 89 value represent relative difference) is further divided by total taxa number in the samples. Canberra  
 90 is a typical metric defined as average relative difference of taxon abundance between two samples.  
 91 In the equation of mGower, the numerator is calculated from the difference of logarithmic  
 92 transformed abundances which is equal to the ratio between larger and smaller abundance before

logarithmic transformation, thus can be regarded as relative difference. And the denominator is total taxa number in the two samples, thus mGower can be classified as average relative difference. Cao also has a numerator related to the ratio of each taxon's abundance between two samples and the total taxa number as the denominator, thus belongs to this type.

- iii. Absolute-difference metrics: the abundance difference between two samples is not divided by the taxa abundances in the two samples. Manhattan is defined as the sum of absolute abundance difference. mManhattan is Manhattan divided by total taxa number in the two samples rather than any abundance-related value, thus classified into this type. Gower and Binomial appear like relative-abundance metrics, but usually show stronger correlation with Manhattan or mManhattan than other relative-abundance metrics. For example, in the empirical data used in this study, Gower and Binomial showed obviously higher correlation coefficients with Manhattan ( $r=0.964$  and  $0.897$ ) than with Bray-Curtis ( $r=0.359$  and  $0.558$ ). Thus, they are classified into this type.
- iv. Squared-sum metrics: the metrics are calculated from the sum of squared abundance difference or product of abundances in two samples. Euclidean is squared root of the squared difference sum and mEuclidean is Euclidean divided by total taxa number in the two samples. Morisita and Morisita-horn are calculated from the product of abundances of each taxon in two samples. By some mathematical deviation, Morisita-horn is actually the squared sum of each taxon's proportion difference divided by the sum of squared proportions, i.e.  $[\sum_k (p_{ik} - p_{jk})^2] / (\sum_k p_{ik}^2 + \sum_k p_{jk}^2)$ , where  $p_{ik}$  and  $p_{jk}$  are the proportions of taxon  $k$  in sample  $i$  and  $j$ , respectively. and Morisita has some minor difference. Thus, these four metrics are classified as a type of squared-sum.

## B. Normalization of stochasticity ratio

Intuitively, the indexes measuring stochasticity and determinism are expected to range from 0% to 100%, and reach the extreme values when community assembly is completely deterministic or stochastic. We defined stochasticity ratio ( $ST_{ij}$ ) as the ratio of average null expectation ( $\overline{E_{ij}}$  or  $\overline{G_{ij}}$ ) to observed similarity ( $C_{ij}$ ) or dissimilarity ( $D_{ij}$ ). Because null expectation is calculated from null model which simulates stochastic assembly, when community assembly is highly stochastic, the average observed similarity or dissimilarity can be very close to the average null expectation, and  $ST$  can approach the accurate value of stochasticity (i.e. 100%). However, also because null model simulates stochastic assembly, the average null expectation always has substantial deviations from 0, no matter how deterministic the observed similarity or dissimilarity is. Therefore, when the community assembly is highly deterministic, the expected stochasticity approaches 0%, but the values of  $ST_{ij}$  always has substantial deviations from 0%. It means that  $ST$  would obviously overestimate stochasticity when expected stochasticity is very low, although it could be relatively accurate when expected stochasticity is high. Thus, we applied the following formula to obtain normalized selection strength ( $NSS$ ) and normalized stochasticity ratio ( $NST$ ).

$$NSS = \frac{SS - TSS}{DSS - TSS} \quad \text{Eq. S6}$$

$$NST = 1 - NSS = \frac{DSS - SS}{DSS - TSS} \quad \text{Eq. S7}$$

128 where  $SS$  is the observed selection strength,  $DSS$  and  $TSS$  are the theoretical extreme values of  $SS$  under  
 129 completely deterministic and stochastic assembly, respectively. After normalization, when community  
 130 assembly is completely deterministic,  $SS$  is equal to  $DSS$ ,  $NSS$  will be 100%, and  $NST$  will be 0%. When  
 131 community assembly is completely stochastic,  $SS$  will be equal to  $TSS$ ,  $NSS$  will be 0% and  $NST$  will be  
 132 100%. Thus,  $NSS$  and  $NST$  are theoretically better than  $SS$  and  $ST$  for measuring determinism and  
 133 stochasticity in community assembly.

134 Before further derivation, we introduce a generalized function  $\xi$  to make equations simpler.

$$\xi(x, y) = \frac{x - y}{x - \delta} \quad \delta = \begin{cases} 0 & x \geq y \\ 1 & x < y \end{cases} \quad \text{Eq. S8}$$

135 If we set  $x$  as the observed similarity between community  $i$  and  $j$  ( $C_{ij}$ ), and set  $y$  as the average null  
 136 expectation of the similarity between community  $i$  and  $j$  ( $\overline{E_{ij}}$ ),  $x \geq y$  means type A situation,  $x < y$  means  
 137 type B situation, and  $\xi(C_{ij}, \overline{E_{ij}})$  is the same as our definition of  $SS$  between community  $i$  and  $j$ . Thus, we  
 138 can simplify Eq. 1 and Eq. 3 in the main text into one equation as below.

$$SS_{ij} = \xi(C_{ij}, \overline{E_{ij}}) \quad \text{Eq. S9}$$

139 and

$$\overline{E_{ij}} = \frac{\sum_{k=1}^{N_r} E_{ij}^{(k)}}{N_r} \quad \text{Eq. S10}$$

140 where  $E_{ij}^{(k)}$  is the null similarity between community  $i$  and  $j$  at the  $k^{\text{th}}$  randomization time of null model  
 141 analysis, and  $N_r$  is the randomization time of null model, which is usually set as 1000 times.

142 To estimate  $DSS$ , we consider two extreme situations. If the deterministic factors lead to more similar  
 143 community structure ( $C_{ij} \geq \overline{E_{ij}}$ , type A situation), the extremely deterministic assembly should have the  
 144 similarity,  ${}^D C_{ij}$ , approaching to the maximum value of 1. In contrast, if the deterministic factors lead to  
 145 more dissimilar community structure ( $C_{ij} < \overline{E_{ij}}$ , type B situation), the extremely deterministic assembly  
 146 should have the dissimilarity close to the maximum and the similarity,  ${}^D C_{ij}$ , close to 0. Thus,  ${}^D SS_{ij}$  and  
 147  ${}^D SS$  can be estimated by following equations.

$${}^D C_{ij} = \begin{cases} 1 & C_{ij} \geq \overline{E_{ij}} \\ 0 & C_{ij} < \overline{E_{ij}} \end{cases} \quad \text{Eq. S11}$$

$${}^D SS_{ij} = \xi({}^D C_{ij}, \overline{E_{ij}}) \quad \text{Eq. S12}$$

$${}^D SS = \frac{\sum_{ij} {}^D SS_{ij}}{n} = \frac{\sum_{ij} \xi({}^D C_{ij}, \overline{E_{ij}})}{n} \quad \text{Eq. S13}$$

148 where  $n$  is the number of pairwise comparisons.

149 Before estimating  ${}^TSS$ , we would like to explain why  ${}^TSS$  cannot be simply set as zero. We need to  
 150 consider the “uncertainty” of similarity/dissimilarity of communities when they are under completely  
 151 stochastic assembly. Here “uncertainty” means that, similarity of each pairwise comparison under  
 152 completely stochastic assembly ( ${}^TC_{ij}$ ) has probability to be any value within the range of null expectation,  
 153 because of the randomness of stochastic assembly. The completely stochastic assembly can be simulated  
 154 by null model.  ${}^TC_{ij}$  can be estimated as null similarity  $E_{ij}$  which is not a certain value but a distribution  
 155  $\{E_{ij}^{(k)}\}_k$  with highest probability usually at the average null expectation  $\overline{E_{ij}}$ . The estimated  $SS$  under  
 156 complete stochastic assembly  ${}^TSS$  is the average relative deviation of  $E_{ij}$  from the mean  $\overline{E_{ij}}$ . Similar to  
 157 standard deviation,  ${}^TSS$  value depends on variance of  $E_{ij}$  and could be equal to zero only if the variances  
 158 of  $E_{ij}$  in every pairwise comparison are all equal to zero. Due to the randomness of stochastic assembly,  
 159  $E_{ij}$  always has variance larger than zero, thus  ${}^TSS$  can never be zero.

160 To estimate  ${}^TSS$ , we simulate stochastic assembly by randomizing the observed community structure  
 161 with a null model algorithm for as many times as necessary (usually  $N_r=1000$  times). At each time of  
 162 randomization, the  $SS$  value of each null pairwise comparison  ${}^TSS_{ij}^{(k)}$  can be calculated from the null  
 163 similarity  $E_{ij}^{(k)}$ . Then, we can obtain the average  $SS$  value of the null communities at each randomization  
 164 time  ${}^TSS^{(k)}$ . To ensure the index  $NSS$  will not exceed 100%,  ${}^TSS$  is calculated as the minimum value of  
 165  $\{{}^TSS^{(k)}\}_k$ . Altogether,  ${}^TSS$  can be estimated as following equations.

$${}^TSS_{ij}^{(k)} = \xi(E_{ij}^{(k)}, \overline{E_{ij}}) \quad \text{Eq. S14}$$

$${}^TSS = \min_k \{ {}^TSS^{(k)} \} \quad \text{Eq. S15}$$

166  
 167 Altogether,  $NSS$  and  $NST$  are calculated as below.

$$NSS = \frac{SS - {}^TSS}{{}^DSS - {}^TSS} = \frac{\sum_{ij} \xi(C_{ij}, \overline{E_{ij}}) - \min_k \{ \sum_{ij} \xi(E_{ij}^{(k)}, \overline{E_{ij}}) \}}{\sum_{ij} \xi({}^DC_{ij}, \overline{E_{ij}}) - \min_k \{ \sum_{ij} \xi(E_{ij}^{(k)}, \overline{E_{ij}}) \}} \quad \text{Eq. S16}$$

$$NST = \frac{{}^DSS - SS}{{}^DSS - {}^TSS} = \frac{\sum_{ij} \xi({}^DC_{ij}, \overline{E_{ij}}) - \sum_{ij} \xi(C_{ij}, \overline{E_{ij}})}{\sum_{ij} \xi({}^DC_{ij}, \overline{E_{ij}}) - \min_k \{ \sum_{ij} \xi(E_{ij}^{(k)}, \overline{E_{ij}}) \}} \quad \text{Eq. S17}$$

168 Because such indexes are originally derived from every pairwise comparison, they are not independent.  
 169 The distribution of  $NSS$  or  $NST$  is unknown and probably not normal. Therefore, the nonparametric  
 170 permutation test, permutational multivariate analysis of variance (PERMANOVA), is used to examine  
 171 whether the communities under different conditions differ in their  $NSS$  and  $NST$ . The  $ST$  and  $NST$   
 172 calculation and PERMANOVA test can be performed using the function “NST” on a web-based pipeline  
 173 (<http://ieg3.rccc.ou.edu:8080/>) built on Galaxy platform (22) or a R package “NST”.

174

## 175 C. Estimating stochasticity in simulated communities

### 176 C1. Simulation models

#### 177 (a) Spatially implicit simulation model

178 We built a spatially implicit simulation model to obtain a total of **21 datasets** with the expected abundance-  
179 based stochasticity ranging from 0% to 100% (5% interval, scenario A in Table S1 and Fig. S1a). Each  
180 dataset has two groups of local communities from **2 plots** under distinct environments (e.g. very hot and  
181 cold environments). The two plots share the same **metacommunity**. Each plot has **12 local communities**  
182 as biological replicates. In each local community, the total richness and total abundances of deterministic  
183 and stochastic species are set according to a certain expected stochasticity. The total abundance of  
184 microorganisms in each local community is set as 20,000, which is a normal sequencing depth of 16S rRNA  
185 gene in many microbial community studies.

186  
187 The **metacommunity structure** (i.e. abundance of each species in the metacommunity) is generated  
188 according to metacommunity zero-sum multinomial distribution (mZSM) (23) derived from Hubbell's  
189 Unified Neutral Theory Model (24), using R package "sads" (25) with  $J=10^8$ ,  $\theta=5000$ , and 10,000 species  
190 sampled. In each **local community**, the **stochastic species** are simulated as a random draw of 100 species  
191 (i.e. the assigned richness of stochastic species in the local community) from metacommunity, with  
192 probabilities proportional to their regional frequencies. The regional frequency of a stochastic species is  
193 calculated from its regional relative abundance according to Sloan's Neutral Model (26) which was also  
194 derived from Hubbell's neutral theory and particularly developed for microbial communities. The dispersal  
195 rate ( $m$ ) is set as 0.1. The abundances of stochastic species in a local community are simulated as a random  
196 draw of a certain number of individuals (i.e. the assigned total abundance of stochastic species in the local  
197 community) from metacommunity into the stochastic species in this community, with probabilities  
198 proportional to the regional relative abundances of the species. We set only two types of **deterministic**  
199 **species**: one is thermophilic, and the other is psychrophilic. The local communities from hot environment  
200 have equal abundances of the thermophilic species, but no psychrophilic species. The communities from  
201 cold environment are under exactly opposite situation, such that the similarity of deterministic species is  
202 100% within group and 0% between groups.

203 The expected stochasticity in a simulated community can be defined as incidence-based or abundance-  
204 based measures as below.

$$ST_{exp.in} = \frac{S_t}{S_t + S_d} \quad \text{Eq. S18}$$

$$ST_{exp.ab} = \frac{J_t}{J_t + J_d} \quad \text{Eq. S19}$$

205	$ST_{exp.in}$	Incidence-based expected stochasticity in a simulated local community.
206	$ST_{exp.ab}$	Abundance-based expected stochasticity in a simulated local community.
207	$S_d$	Richness of deterministic species in a simulated local community.
208	$J_d$	Total abundance of deterministic species in a simulated local community.
209	$S_t$	Richness of stochastic species in a simulated local community.
210	$J_t$	Total abundance of stochastic species in a simulated local community.

211

## 212 (b) Spatially explicit simulation model

213 To examine scale dependence of stochasticity estimation, we built a spatially explicit simulation model  
214 (Scenario B-F in Table S1, Fig. 2a, and Fig. S1b). The model has four-level metacommunities, including  
215 local (for each site), regional, continental, and global metacommunities. In the model, an area of 16,384  
216 (128×128) cells are divided into 4 (2×2) continents, each continent is divided into 4 (2×2) regions, and each  
217 region is divided into 4 (2 × 2) sites. Each site has 4 (2×2) plots, sharing the same local metacommunity.  
218 Each plot has 64 (8×8) cells, and each cell represents a local community with 20,000 individuals (Fig. 2a).  
219 We take all individuals from a single cell as a sample, and a certain number of samples from each plot (6  
220 samples/plot unless specified) to get a simulated dataset. Ecological stochasticity was estimated with  
221 different indexes based on the pairwise comparisons of all samples within each unit at different spatial  
222 scales, i.e. plot, site, region, continent, or global.

223

224 To investigate more complicated deterministic forces, **deterministic species** were simulated under three  
225 types of scenarios. The **first** scenario (Scenario B in Table S1 and Fig. S1b) is simple abiotic filtering  
226 without environmental noise. Plots at the same row (like latitude) have the same temperature, while  
227 temperature increases by 2°C per plot along each column (like longitude), from 0°C at the top (northmost)  
228 plot to 30°C at the bottom (southmost) plot (Fig. 2a). The temperature is homogeneous within each plot.  
229 All local communities (cells) under each temperature have equal abundance of the only deterministic  
230 species which prefer this temperature.

231 The **second** type of scenarios (Scenario C-E in Table S1) is abiotic filtering with environmental noise.  
232 The mean temperature of each plot is the same as that in the first scenario, but the temperature in each cell  
233 is a random value from a normal distribution with a certain standard deviation (temperature deviation,  $\sigma_t$ ).  
234 Temperature within each cell is still set homogenous. The abundances of deterministic species in each cell  
235 are determined by a Gaussian function as below (Eq. S20). The temperature deviation in each plot is set at  
236 different level comparing to fitness deviation ( $\sigma_f$ , defined in Eq. S20), to simulate low ( $\sigma_t=5\%\sigma_f$ , Table S1  
237 scenario C), medium ( $\sigma_t=25\%\sigma_f$ , Table S1 scenario D), and high ( $\sigma_t=200\%\sigma_f$ , Table S1 scenario E)  
238 environmental noise.

$$A_{ij} = J_{d0} \exp \left[ -\frac{(T_j - T_i)^2}{2\sigma_f^2} \right] \quad \text{Eq. S20}$$

239	$A_{ij}$	Abundance of species $i$ in local community $j$ .
240	$J_{d0}$	Expected maximum abundance of deterministic species $i$ in a local community.
241	$T_j$	Temperature of local community $j$ .
242	$T_i$	Optimum temperature of species $i$ .
243	$\sigma_f$	Fitness deviation, set as 0.4 in this study.

244

245 The **third** type of scenario is to consider biotic competition (Table S1 scenario F). Each of the 256  
246 competitors randomly occupies one cell at the very beginning. Then, the competitors randomly disperse to  
247 an adjacent cell at each time step, with equal probabilities to all four directions, until all cells are occupied



248 by competitors. In each cell (i.e. each local community), the first-arrived competitor excludes other  
249 competitor(s) and stops them passing through the cell.

250 The **fourth** type of scenario is to investigate community under complex deterministic forces (Table S1  
251 scenario G). In each simulated community with deterministic part, deterministic species controlled by  
252 abiotic filtering without environmental noise are simulated as in the first scenario, and then combined at a  
253 certain abundance ratio with species controlled by competition which are simulated as in the third types of  
254 scenarios.

255 In different scenarios, **stochastic species** were simulated in the same way as below. First, a global  
256 metacommunity was generated in the same way as that in spatially implicit model, with  $J=10^9$ ,  $\theta=5000$ , and  
257 10,000 species sampled. Second, we developed a two-step random assembly model to simulate stochastic  
258 assembly in the spatially explicit model. At the first step, a certain number of species ( $S_i$ ) are randomly  
259 drawn from the higher-level metacommunity to a lower-level (meta)community, according to the expected  
260 occurrence frequencies and relative abundance of all species in the higher-level metacommunity as  
261 described in the spatially implicit model. The expected occurrence frequencies are calculated according to  
262 Sloan's neutral model with a certain dispersal rate ( $m_1$ ). At the second step, the species ( $S_i$ ) for each lower-  
263 level (meta)community are randomly drawn from three sources, the higher-level metacommunity (with a  
264 dispersal rate of  $m_1$ ), first-step pool of this lower-level (meta)community ( $m_2$ ), and all first-step pools of  
265 adjacent (meta)communities ( $m_3$ ), to simulate dispersal from higher-level species pool and adjacent  
266 communities, respectively. Third, we applied this two-step random assembly model to simulate the  
267 (meta)communities at each level. Each continental metacommunity (5,200 species and  $8 \times 10^7$  individuals)  
268 is simulated as two-step random draw from global metacommunity with the dispersal rates of  $m_1=0.001$ ,  
269  $m_1=0.997$ , and  $m_2=0.002$ . In the same way, we simulated each regional metacommunity (2,700 species,  
270  $2 \times 10^7$  individuals,  $m_1=0.05$ ,  $m_1=0.8$ ,  $m_2=0.15$ ), local metacommunity (each site, 1,400 species,  $5 \times 10^6$   
271 individuals,  $m_1=0.1$ ,  $m_1=0.5$ ,  $m_2=0.4$ ), and local community (each cell, 100 species,  $m_1=0.2$ ,  $m_1=0.2$ ,  
272  $m_2=0.6$ ). In each local community (cell), the total individual number of stochastic species depends on the  
273 expected abundance-based stochasticity ( $ST_{exp.ab}$ ). The dispersal rates from the adjacent (meta)community  
274 pool ( $m_2$ ) are higher at lower spatial scales because dispersal is easier at smaller spatial scales.

275 For each scenario, we combined deterministic and stochastic species at different abundance ratios to  
276 generate 11 datasets with expected abundance-based stochasticity ranging from 0% to 100% (10% interval).  
277 For each dataset, we estimated stochasticity at different spatial scales with the three indexes and evaluated  
278 their accuracy and precision as described below.

279

## 280 **C2. Stochasticity indexes**

281 In each dataset, the stochasticity within each group of simulated communities was estimated by  $ST$  and  $NST$ ,  
282 and the neutral species percentage ( $NP$ ).  $NST$  and  $ST$  in simulated communities were calculated based on  
283 the null model algorithm "PF" (described in part D and Table S4) and various similarity metrics (Table S2).

284 Sloan et al. (26) developed a neutral model about the relationship between occurrence frequency and  
285 relative abundance in source community for microbial communities. We applied Sloan's neutral model to  
286 fit the occurrence frequency of each species in a group of communities and the relative abundance in the  
287 whole dataset. The species within the 95% confidence interval of Sloan neutral model are defined as neutral  
288 species (27). The abundance-weighted and unweighted percentage of neutral species ( $NP$ ) were used to

289 estimate abundance-based and incidence-based stochasticity, respectively. *NP* was calculated using the R  
 290 codes reported by Burns et al. (27).

291 Modified Roup-Crick metrics (RC) and standardized effect size (SES) were also applied to  
 292 communities simulated by spatial implicit model, calculated as previously reported (28-30). The percentage  
 293 of turnovers with  $|\text{SES}| < 2$  and that with  $|\text{RC}| < 0.95$  were counted as stochastic turnover ratio (SR) based on  
 294 SES ( $\text{SR}_{\text{SES}}$ ) and SR based on RC ( $\text{SR}_{\text{RC}}$ ), respectively (29). These indexes showed obviously worse  
 295 quantitative performance than NST and ST for data of spatial implicit model, when calculated based on  
 296 Bray-Curtis and Ružička (Fig. S6). Thus, we did not further apply them to other simulated data or test  
 297 various metrics.

298

### 299 *C3. Evaluating the accuracy and precision of different stochasticity indexes*

300 We evaluated the performance of each stochasticity index quantitatively by accuracy and precision  
 301 coefficients. Concordance Correlation Coefficient (CCC) was developed as a measure of agreement  
 302 between two methods (31). It has meaningful components of accuracy ( $\chi_a$ , Eq. S21) and precision ( $\rho$ , Eq.  
 303 S22) (32), in which the precision coefficient is the same as Pearson correlation coefficient. Thus, we applied  
 304 these two coefficients to evaluate the accuracy and precision of stochasticity values estimated by different  
 305 methods (Table S2). Based on the equation, high accuracy coefficient value means the estimated values  
 306 have very similar mean and variance as true values. In contrast, high precision coefficient means the  
 307 variation of estimated values have very similar trend as true values, thus can precisely reflect the relative  
 308 change of true values. Therefore, a qualified stochasticity index should have high scores in both accuracy  
 309 and precision coefficients. For example, we assume the true values are 20%, 40%, 60%, 80%, 90% in  
 310 sequence. When the estimated values are 90%, 60%, 80%, 40%, 20%, the accuracy is very high ( $\chi_a=1$ ) but  
 311 the precision is very low ( $\rho=-0.86$ , negative), thus the index is useless. When the estimated values are 2%,  
 312 4%, 6%, 8%, 9%, the accuracy is very low ( $\chi_a=0.04$ ) but the precision is very high ( $\rho=1$ ), thus the index  
 313 cannot reflect the magnitude of true value but can be used to estimate the relative changes of true values.  
 314 When the estimated values are 19%, 41%, 60%, 79%, 91%, the accuracy and precision are both very high  
 315 ( $>0.99$ ), thus the index can be used to estimate the true values.

$$\chi_a = \frac{2\sigma_x\sigma_y}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2} \quad \text{Eq. S21}$$

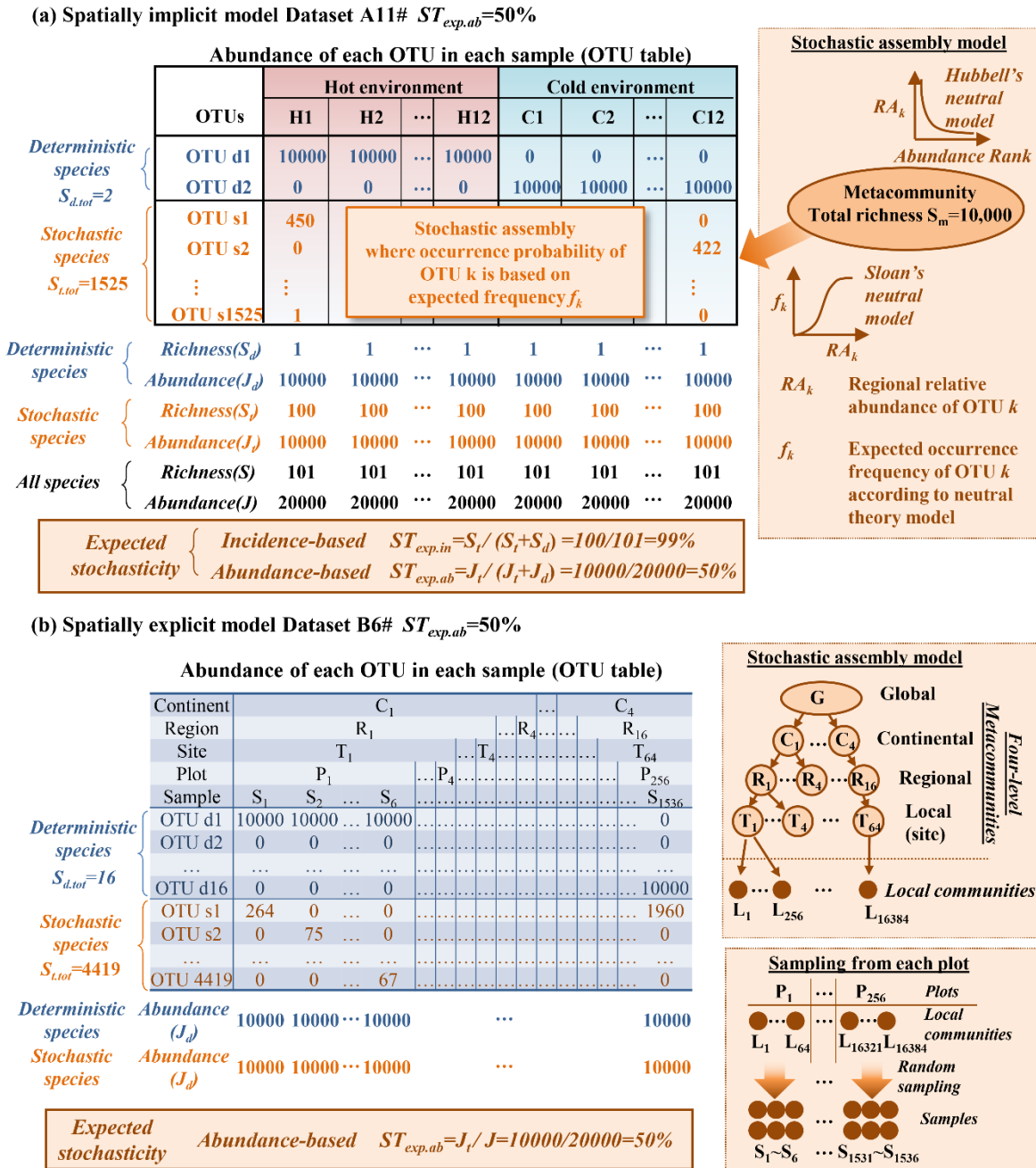
$$\rho = \frac{\sigma_{yx}}{\sigma_x\sigma_y} \quad \text{Eq. S22}$$

316

317	$\sigma_{yx}$	Covariance of $x$ and $y$ . In our study, $x$ is expected stochasticity, and $y$ is estimated stochasticity.
318	$\sigma_x^2$	Variance of $x$ .
319	$\sigma_y^2$	Variance of $y$ .
320	$\mu_x$	Mean of $x$ .
321	$\mu_y$	Mean of $y$ .

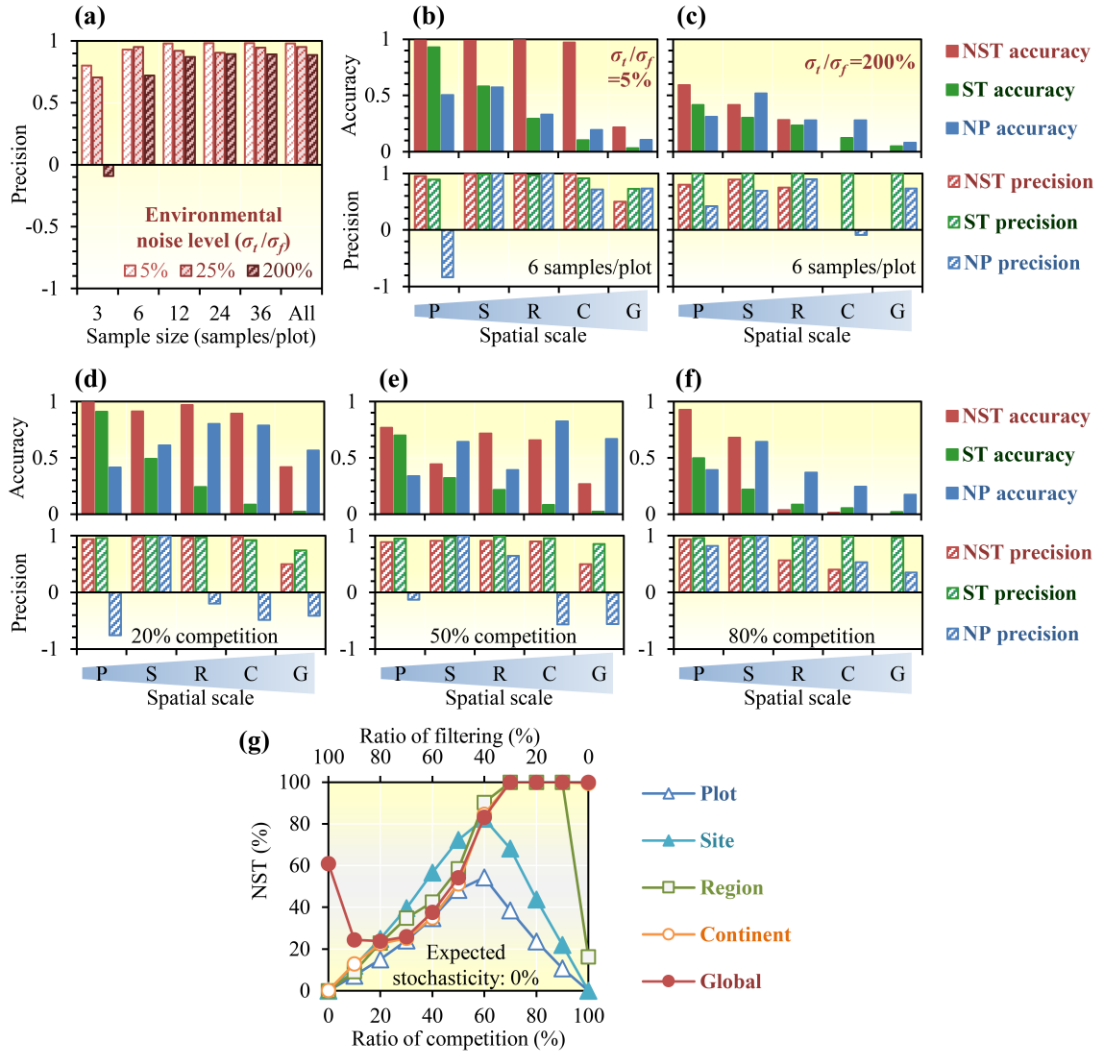
#### 322 **D. Null model algorithms**

323 In general, there are nine major types of null model algorithms for species co-occurrence analysis,  
324 previously elucidated by Gotelli (33) (Table S4). When randomizing the observed communities, different  
325 null model algorithms use different ways to constrain the occurrence frequency of each taxon and taxon  
326 richness in each sample. We listed the abbreviation and formula to calculate the probability of a taxon  
327 present in a sample in each algorithm in Table S4. If abundance weighted metrics are used, after getting  
328 occurrence data matrix, abundance can be assigned as random draw of individuals with probabilities  
329 proportional to the regional relative abundances of the taxa as previously described by Stegen et al (29).  
330 All samples of the empirical dataset were considered as from the same regional species pool, thus  
331 randomization was performed across all samples. ST and NST can be calculated based on different null  
332 model algorithms and different metrics using the function “NST” on the pipeline  
333 (<http://ieg3.rccc.ou.edu:8080>), or using the function “tNST” in a R package “NST”.  
334



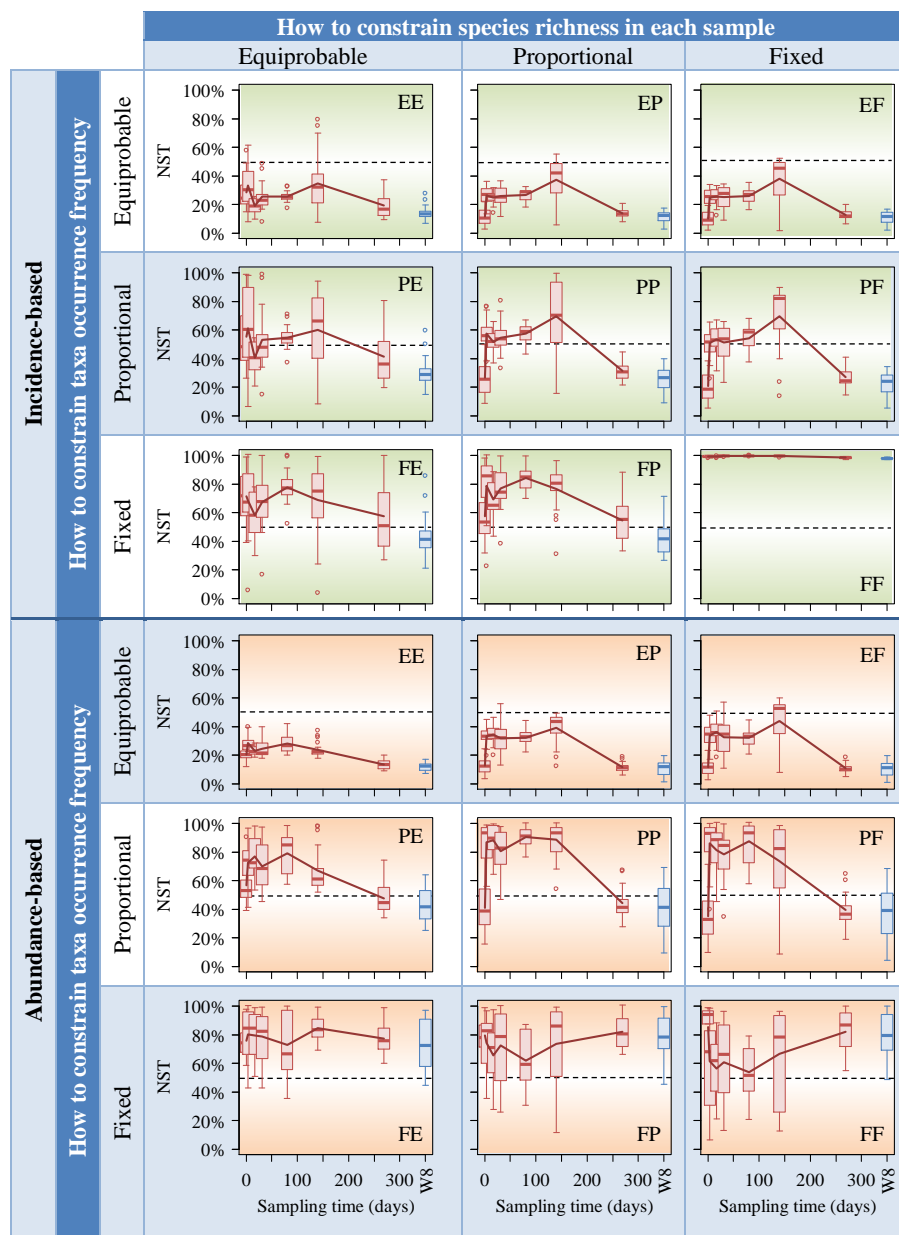
336  
 337 **Fig. S1. Community composition, stochastic assembly model and expected stochasticity in the**  
 338 **example datasets from (a) the spatially implicit model and (b) the spatially explicit model.** The OTU  
 339 tables and associated annotation (left part of each panel) show the abundances, richness, and abundance  
 340 sum of deterministic (blue) and stochastic (orange) species in each sample, while  $S_{d,tot}$  and  $S_{s,tot}$  represent  
 341 the total richness of deterministic and stochastic species across all samples in a dataset. Each column of  
 342 the OTU tables represents a sample, and each OTU represents a species. The bottom box (brown) in each  
 343 panel shows how the expected stochasticity is calculated. The box of stochastic assembly model on the

344 right panel shows how the stochastic species are simulated. In both spatially implicit (a) and explicit (b)  
345 models, the top-level metacommunity is simulated according to Hubbell's neutral theory model, and each  
346 local community is generated as random draw from local metacommunity based on Sloan's neutral  
347 model. In the spatially explicit model with four-level metacommunities, lower-level metacommunities are  
348 simulated as random draw of species from higher-level metacommunities using a two-step random  
349 assembly method based on Sloan's neutral model. In the spatially implicit model, all simulated local  
350 communities (12/plot) are taken as samples. In the spatially explicit model, a certain number (6 in this  
351 example) of local communities are taken as samples from each plot (the box about sampling). See  
352 supplementary text part C and Table S1 for details.



355  
 356 **Fig. S2. Accuracy and precision of stochasticity estimation in simulated communities under abiotic**  
 357 **filtering with noise from deterministic environmental factors (i.e. environmental noise) or with**  
 358 **competition. (a) Influence of sample size on precision of normalized stochasticity ratio (NST) at plot scale**  
 359 **with different degrees of environmental noise. The precision obviously decreased when sample size was**  
 360 **not large enough ( $\leq 6$  samples/plot), which was more obvious when environmental noise was higher. The**  
 361 **accuracy did not show obvious trend, thus not showed here. (b) Accuracy and precision of stochasticity**  
 362 **estimation under low environmental noise ( $\sigma_t/\sigma_f=5\%$ , Table S1 scenario C) and (c) high environmental**  
 363 **noise ( $\sigma_t/\sigma_f=200\%$ , Table S1 scenario E) across different spatial scales (P, plot; S, site; R, region; C,**  
 364 **continent; G, global). NST can have high accuracy and precision when environmental noise was not too**  
 365 **high. All indexes had very low accuracy especially at large scales, although stochasticity ratio (ST) still**  
 366 **showed high precision across different spatial scales. The community was simulated by spatially explicit**  
 367 **model (see Supplementary text C and Table S1 for details).  $\sigma_t$ , the standard deviation of temperature in each**  
 368 **plot;  $\sigma_f$ , the fitness deviation defined in Eq. S20. (d, e, f) Accuracy and precision of stochasticity estimation**  
 369 **when deterministic species include some controlled by abiotic filtering and the others controlled by**  
 370 **competition with the ratio of (d) 20%, (e) 50%, and (f) 80%. NST (red bars), normalized stochasticity ratio;**  
 371 **ST (green bars), stochasticity ratio; NP (blue bars), neutral species percentage. NST and ST were based on**

372 Ružička similarity index (Table S3) and the null model “PF” (Table S4). Accuracy (solid color bars) and  
373 precision (diagonal strip bars) were evaluated by the coefficients derived from concordance correlation  
374 coefficient (Eq. S21-22). (g) *NST* of simulated communities controlled by abiotic filtering and  
375 competition without stochastic assembly, estimated across different scales. Although the expected  
376 stochasticity is zero, *NST* still overestimated stochasticity. The overestimation is more obvious  
377 when filtering and competition are comparable (e.g. *NST*>50% when ratio of competition is  
378 50~60%). The overestimation is the lowest at plot level. In contrast, *NST* became up to 100% at  
379 regional to global scales when the ratio of competition is 70~90%.  
380

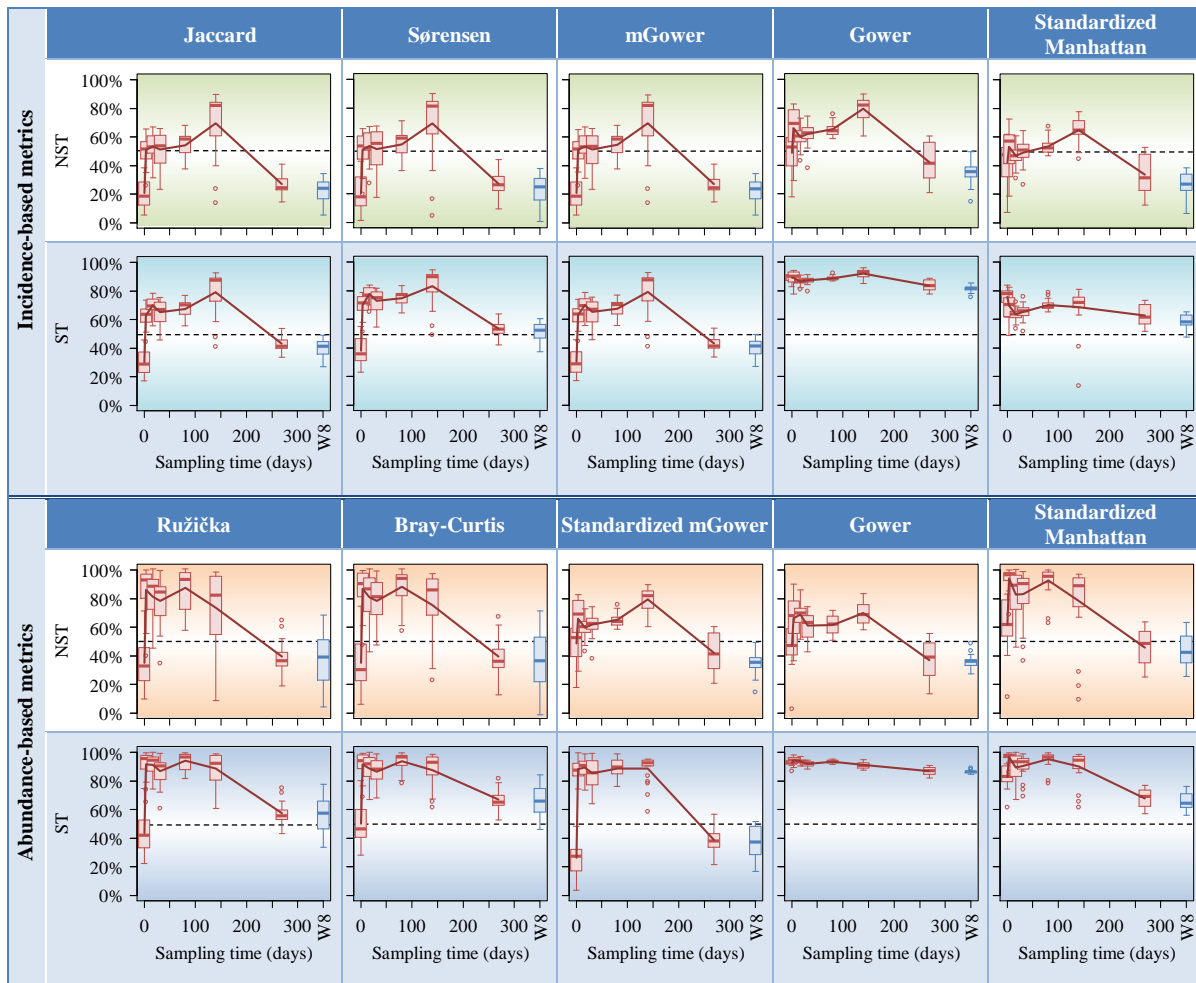


381 **Fig. S3. Effects of null model algorithms on NST estimation with incidence-based Jaccard (upper)**  
 382 **and abundance-based Ružička (lower) similarity metric.** The emulsified vegetable oil was injected at  
 383 Day 1 and almost exhausted at Day 269, and had minimal impact on the control well (W8). Therefore, at  
 384 Day 0, Day 269, and W8, the microbial communities were under very high selection pressure caused by  
 385 high concentrations of pollutants (e.g. heavy metals, nitrate) and carbon poor (34, 35), thus they should be  
 386 under more deterministic assembly with low stochasticity. The vegetable oil injection significantly  
 387 increased carbon resources (electron donors) and decreased some pollutants (34), thus should reduce the  
 388 impact of selection and increase stochasticity. The null model PP and PF showed more significant and  
 389 expected variations of stochasticity along time. Null model EP and EF showed similar trend but much less  
 390 estimated stochasticity than expected. Other null models did not show consistent or clear trend. The NST  
 391 values based on Ružička were obviously higher than those based on Jaccard although the trend is very  
 392 similar. See Table S4 for the detailed algorithms of null models.



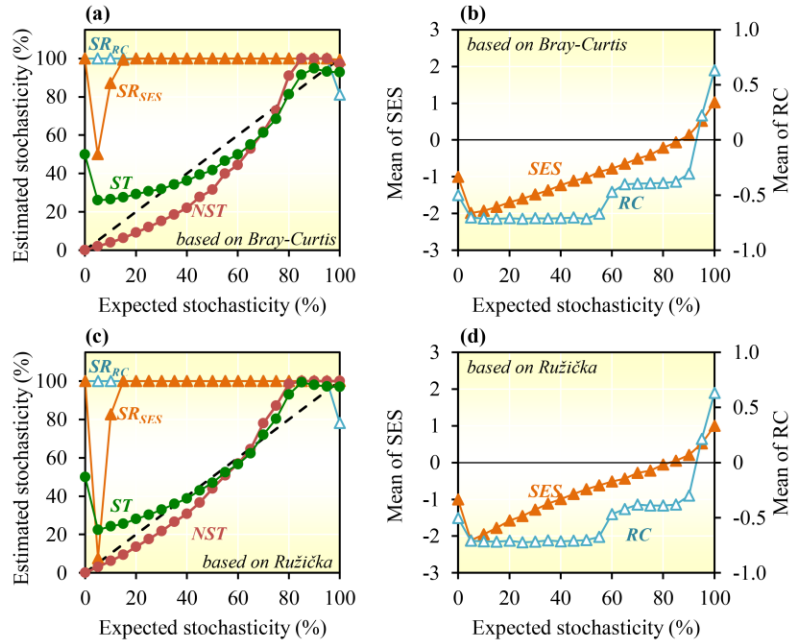


393 **Fig. S4. Effects of similarity metrics on NST estimation.** Most metrics showed very similar trend of  
 394 stochasticity variation, but the magnitude of *NST* obviously varied among some metrics. Most abundance-  
 395 based metrics showed similar trend of stochasticity variation but obviously higher magnitude of *NST*  
 396 comparing to their corresponding incidence-based metrics. Null model algorithm used was “PF” (Table  
 397 S4). See Table S3 for detailed definition of each similarity metric, and supplementary text part A for  
 398 metrics standardization method.



400 **Fig. S5. Comparison between *NST* and *ST* estimated with different similarity metrics.** Although *NST*  
 401 and *ST* basically showed consistent trend, *NST* exhibited much less variation when based on different  
 402 similarity metrics, i.e. *NST* is less sensitive to metric selection than *ST*.

403



404

405 **Fig. S6. Comparison of different null-model-based indexes applied to the simulated**  
 406 **communities with various levels of expected stochasticity. (a)** Estimated stochasticity with  
 407 different indexes based on Bray-Curtis; **(b)** Mean of standardized effect size (SES) and modified  
 408 Raup-Crick metrics (RC) based on Bray-Curtis; **(c)** Estimated stochasticity with different indexes  
 409 based on Ružička; **(d)** Mean of SES and RC based on Ružička. The simulation model was spatially  
 410 implicit. *NST* (red), normalized stochasticity ratio; *ST* (green), stochasticity ratio; SES (orange,  
 411 panel b and d), standardized effect size;  $SR_{SES}$  (orange, panel a and c), stochastic turnover ratio  
 412 based on SES, i.e. percentage of turnovers with  $|SES| < 2$ ; RC (aqua, panel b and d), modified Raup-  
 413 Crick metrics;  $SR_{RC}$  (aqua, panel a and c), stochastic turnover ratio based on RC, i.e. percentage  
 414 of turnovers with  $|RC| < 0.95$ .  
 415

416 **Supplementary Tables**

417 **Table S1. The richness ( $S$ ) and abundance ( $J$ ) of species under stochastic and deterministic**  
 418 **assembly in simulated communities.** Seven scenarios (A-G) are considered. Scenario A has 21 datasets,  
 419 while the others have 11 datasets.  
 420

Dataset #	Expected stochasticity		In each local community						In all samples		
	Incidence -based $ST_{exp.in}$	Abundance -based $ST_{exp.ab}$	Stochastic species		Deterministic species		All species		Stochastic species $S_{t.tot}$	Deterministic species $S_{d.tot}$	All species $S_{tot}$
			$S_t^{[1]}$	$J_t$	$S_d$	$J_d$	$S$	$J$			
<b>Scenario A: Spatially implicit model, abiotic filtering without environmental noise</b> (12 samples/plot × 2 plots = 24 samples)											
<b>A1</b>	0%	0%	0	0	1	20000	1	20000	0	2	2
<b>A2</b>	99%	5%	100	1000	1	19000	101	20000	1353	2	1355
<b>A3</b>	99%	10%	100	2000	1	18000	101	20000	1392	2	1394
<b>A4</b>	99%	15%	100	3000	1	17000	101	20000	1433	2	1435
<b>A5</b>	99%	20%	100	4000	1	16000	101	20000	1462	2	1464
<b>A6</b>	99%	25%	100	5000	1	15000	101	20000	1484	2	1486
<b>A7</b>	99%	30%	100	6000	1	14000	101	20000	1498	2	1500
<b>A8</b>	99%	35%	100	7000	1	13000	101	20000	1528	2	1530
<b>A9</b>	99%	40%	100	8000	1	12000	101	20000	1532	2	1534
<b>A10</b>	99%	45%	100	9000	1	11000	101	20000	1505	2	1507
<b>A11</b>	99%	50%	100	10000	1	10000	101	20000	1525	2	1527
<b>A12</b>	99%	55%	100	11000	1	9000	101	20000	1577	2	1579
<b>A13</b>	99%	60%	100	12000	1	8000	101	20000	1560	2	1562
<b>A14</b>	99%	65%	100	13000	1	7000	101	20000	1573	2	1575
<b>A15</b>	99%	70%	100	14000	1	6000	101	20000	1566	2	1568
<b>A16</b>	99%	75%	100	15000	1	5000	101	20000	1572	2	1574
<b>A17</b>	99%	80%	100	16000	1	4000	101	20000	1591	2	1593
<b>A18</b>	99%	85%	100	17000	1	3000	101	20000	1593	2	1595
<b>A19</b>	99%	90%	100	18000	1	2000	101	20000	1591	2	1593
<b>A20</b>	99%	95%	100	19000	1	1000	101	20000	1599	2	1601
<b>A21</b>	100%	100%	100	20000	0	0	100	20000	1559	0	1559
<b>Scenario B: Spatially explicit model, abiotic filtering without environmental noise</b> (6 samples/plot × 256 plots = 1536 samples)											
<b>B1</b>	0%	0%	0	0	1	20000	1	20000	0	16	16
<b>B2</b>	99%	10%	100	2000	1	18000	101	20000	3720	16	3736
<b>B3</b>	99%	20%	100	4000	1	16000	101	20000	3965	16	3981
<b>B4</b>	99%	30%	100	6000	1	14000	101	20000	4199	16	4215
<b>B5</b>	99%	40%	100	8000	1	12000	101	20000	4299	16	4315
<b>B6</b>	99%	50%	100	10000	1	10000	101	20000	4419	16	4435
<b>B7</b>	99%	60%	100	12000	1	8000	101	20000	4516	16	4532
<b>B8</b>	99%	70%	100	14000	1	6000	101	20000	4670	16	4686
<b>B9</b>	99%	80%	100	16000	1	4000	101	20000	4688	16	4704
<b>B10</b>	99%	90%	100	18000	1	2000	101	20000	4832	16	4848
<b>B11</b>	100%	100%	100	20000	0	0	100	20000	4807	0	4807

421 **Table S1. Continued**

Dataset #	Expected stochasticity		In each local community						In all samples		
	$ST_{exp.in}$	$ST_{exp.ab}$	$S_t$	$J_t$	$S_d$	$J_d$	$S$	$J$	$S_{t.tot}$	$S_{d.tot}$	$S_{tot}$
<b>Scenario C: Spatially explicit model, abiotic filtering with low environmental noise</b> ( $\sigma_i/\sigma_f=5\%^{[2]}$ , 6 samples/plot $\times$ 256 plots = 1536 samples)											
<b>C1</b>	0%	0%	0	0	1	19628~ 19999	1	19628~ 19999	0	16	16
<b>C2</b>	99%	10.0%~ 10.2%	100	2000	1	17682~ 17999	101	19682~ 19999	3701	16	3717
<b>C3</b>	99%	20.0%~ 20.2%	100	4000	1	15773~ 15999	101	19773~ 19999	3980	16	3996
<b>C4</b>	99%	30.0%~ 30.4%	100	6000	1	13752~ 13999	101	19752~ 19999	4161	16	4177
<b>C5</b>	99%	40.0%~ 40.4%	100	8000	1	11814~ 11999	101	19814~ 19999	4349	16	4365
<b>C6</b>	99%	50.0%~ 50.4%	100	10000	1	9847~ 9999	101	19847~ 19999	4454	16	4470
<b>C7</b>	99%	60.0%~ 60.5%	100	12000	1	7851~ 7999	101	19851~ 19999	4549	16	4565
<b>C8</b>	99%	70.0%~ 70.4%	100	14000	1	5888~ 5999	101	19888~ 19999	4592	16	4608
<b>C9</b>	99%	80.0%~ 80.2%	100	16000	1	3945~ 3999	101	19945~ 19999	4699	16	4715
<b>C10</b>	99%	90.0%~ 90.1%	100	18000	1	1969~ 1999	101	19969~ 19999	4804	16	4820
<b>C11</b>	100%	100%	100	20000	0	0	100	20000	4798	0	4798
<b>Scenario D: Spatially explicit model, abiotic filtering with medium environmental noise</b> ( $\sigma_i/\sigma_f=25\%$ , 6 samples/plot $\times$ 256 plots = 1536 samples)											
<b>D1</b>	0%	0%	0	0	1~2	15066~ 19999	1~2	15066~ 19999	0	16	16
<b>D2</b>	99%	10.0%~ 13.0%	100	2000	1~2	13358~ 17999	101~ 102	15358~ 19999	3735	16	3751
<b>D3</b>	99%	20.0%~ 25.4%	100	4000	1~2	11721~ 15999	101~ 102	15721~ 19999	3956	16	3972
<b>D4</b>	99%	30.0%~ 40.3%	100	6000	1~2	8870~ 13999	101~ 102	14870~ 19999	4147	16	4163
<b>D5</b>	99%	40.0%~ 49.3%	100	8000	1~2	8228~ 11999	101~ 102	16228~ 19999	4308	16	4324
<b>D6</b>	99%	50.0%~ 58.3%	100	10000	1~2	7145~ 9999	101~ 102	17145~ 19999	4436	16	4452
<b>D7</b>	99%	60.0%~ 71.0%	100	12000	1~2	4890~ 7999	101~ 102	16890~ 19999	4523	16	4539
<b>D8</b>	99%	70.0%~ 76.6%	100	14000	1	4286~ 5999	101	18286~ 19999	4598	16	4614
<b>D9</b>	99%	80.0%~ 86.3%	100	16000	1~2	2534~ 3999	101~ 102	18534~ 19999	4725	16	4741
<b>D10</b>	99%	90.0%~ 93.3%	100	18000	1	1289~ 1999	101	19289~ 19999	4745	16	4761
<b>D11</b>	100%	100%	100	20000	0	0	100	20000	4849	0	4849

423 Table S1. Continued

Dataset #	Expected stochasticity		In each local community						In all samples		
	$ST_{exp.in}$	$ST_{exp.ab}$	$S_t$	$J_t$	$S_d$	$J_d$	$S$	$J$	$S_{t.tot}$	$S_{d.tot}$	$S_{tot}$
<b>Scenario E: Spatially explicit model, abiotic filtering with high environmental noise</b> ( $\sigma_i/\sigma_f=200\%$ , 6 samples/plot $\times$ 256 plots = 1536 samples)											
E1	0%	0%	0	0	1~2	3~ 19999	1~2	3~ 19999	0	16	16
E2	99%	10.0%~ 100.0%	100	2000	0~2	0~ 17999	100~ 102	2000~ 19999	3747	16	3763
E3	99%	20.0%~ 100.0%	100	4000	0~2	0~ 15999	100~ 102	4000~ 19999	3958	16	3974
E4	99%	30.0%~ 100.0%	100	6000	1~2	2~ 13999	101~ 102	6002~ 19999	4153	16	4169
E5	99%	40.0%~ 100.0%	100	8000	0~2	0~ 11999	100~ 102	8000~ 19999	4300	16	4316
E6	99%	50.0%~ 100.0%	100	10000	0~2	0~ 9999	100~ 102	10000~ 19999	4455	16	4471
E7	99%	60.0%~ 100.0%	100	12000	0~2	0~ 7999	100~ 102	12000~ 19999	4522	16	4538
E8	99%	70.0%~ 100.0%	100	14000	1~2	1~ 5999	101~ 102	14001~ 19999	4643	16	4659
E9	99%	80.0%~ 100.0%	100	16000	0~2	0~ 3999	100~ 102	16000~ 19999	4720	16	4736
E10	99%	90.0%~ 100.0%	100	18000	0~2	0~ 1999	100~ 102	18000~ 19999	4779	16	4795
E11	100%	100%	100	20000	0	0	100	20000	4887	0	4887
<b>Scenario F: Spatially explicit model, biotic interspecies competition</b> (6 samples/plot $\times$ 256 plots = 1536 samples)											
F1	0%	0%	0	0	1	20000	1	20000	0	249	249
F2	99%	10%	100	2000	1	18000	101	20000	3721	251	3972
F3	99%	20%	100	4000	1	16000	101	20000	3964	247	4211
F4	99%	30%	100	6000	1	14000	101	20000	4203	248	4451
F5	99%	40%	100	8000	1	12000	101	20000	4351	249	4600
F6	99%	50%	100	10000	1	10000	101	20000	4438	244	4682
F7	99%	60%	100	12000	1	8000	101	20000	4533	253	4786
F8	99%	70%	100	14000	1	6000	101	20000	4629	248	4877
F9	99%	80%	100	16000	1	4000	101	20000	4714	250	4964
F10	99%	90%	100	18000	1	2000	101	20000	4787	251	5038
F11	100%	100%	100	20000	0	0	100	20000	4780	0	4780
<b>Scenario G: Spatially explicit model, abiotic filtering and competition</b> (6 samples/plot $\times$ 256 plots = 1536 samples)											
G1	0%	0%	0	0	2	20000	2	20000	0	272	272
G2	99%	10%	100	2000	2	18000	102	20000	4702	272	4974
G3	99%	20%	100	4000	2	16000	102	20000	4923	272	5195
G4	99%	30%	100	6000	2	14000	102	20000	5107	272	5379
G5	99%	40%	100	8000	2	12000	102	20000	5234	272	5506
G6	99%	50%	100	10000	2	10000	102	20000	5393	272	5665
G7	99%	60%	100	12000	2	8000	102	20000	5447	272	5719

Dataset #	Expected stochasticity		In each local community						In all samples		
	$ST_{exp.in}$	$ST_{exp.ab}$	$S_t$	$J_t$	$S_d$	$J_d$	$S$	$J$	$S_{t.tot}$	$S_{d.tot}$	$S_{tot}$
<b>G8</b>	99%	70%	100	14000	2	6000	102	20000	5557	272	5829
<b>G9</b>	99%	80%	100	16000	2	4000	102	20000	5588	272	5860
<b>G10</b>	99%	90%	100	18000	2	2000	102	20000	5715	272	5987
<b>G11</b>	100%	100%	100	20000	0	0	100	20000	5723	0	5723

424

425

426

427

428

<sup>[1]</sup>  $S_t$ ,  $S_d$ , and  $S$  are the richness of stochastic, deterministic, and all species in each local community;  $J_t$ ,  $J_d$ , and  $J$  are the abundance of stochastic, deterministic, and all species in each local community;  $S_{t.tot}$ ,  $S_{d.tot}$ , and  $S_{tot}$  are the overall richness of stochastic, deterministic, and all species in all samples.

<sup>[2]</sup>  $\sigma_t$  is the standard deviation of temperature in each plot,  $\sigma_f$  is fitness deviation defined in Eq. S20.

429 **Table S2. Accuracy and precision of stochasticity in simulated communities estimated by different**  
 430 **indexes based on various similarity metrics.**  
 431

Types	Similarity metrics	Accuracy coefficient <sup>[1]</sup>			Precision coefficient <sup>[1]</sup>		
		<i>NST</i> <sup>[2]</sup>	<i>ST</i> <sup>[2]</sup>	<i>NP</i> <sup>[2]</sup>	<i>NST</i>	<i>ST</i>	<i>NP</i>
Incidence-based	Jaccard <sup>[3]</sup>	0.999	0.779	0.744 <sup>[6]</sup>	1.000	0.999	0.118
	Sørensen <sup>[3]</sup>	0.999	0.764		1.000	0.999	
	Kulczynski	0.999	0.762		1.000	0.999	
	Gower <sup>[3]</sup>	0.994	0.773		1.000	0.999	
	Manhattan	0.999	0.607		1.000	0.999	
	Euclidean ( <i>S</i> <sup>[4]</sup> )	0.999	0.625		1.000	0.999	
	mEuclidean <sup>[5]</sup> ( <i>S</i> )	0.999	0.639		1.000	0.999	
Abundance-based	Ružička	0.985	0.968	0.462	0.985	0.925	0.275
	Bray-Curtis	0.966	0.969		0.969	0.897	
	Kulczynski	0.965	0.969		0.968	0.896	
	Canberra	0.416	0.255		0.532	0.811	
	Gower	0.616	0.308		0.660	0.865	
	mGower <sup>[5]</sup> ( <i>S</i> )	0.989	0.716		0.989	0.920	
	Morisita	0.629	0.809		0.724	0.606	
	Morisita-Horn	0.631	0.810		0.723	0.609	
	Manhattan ( <i>S</i> )	0.986	0.754		0.987	0.917	
	mManhattan <sup>[5]</sup> ( <i>S</i> )	0.989	0.714		0.988	0.920	
	Euclidean ( <i>S</i> )	0.652	0.323		0.928	0.822	
	mEuclidean <sup>[5]</sup> ( <i>S</i> )	0.637	0.275		0.936	0.829	
	Binomial ( <i>S</i> )	0.380	0.148		0.370	0.411	
	Chao	0.967	0.919		0.981	0.940	
Cao ( <i>S</i> )	0.753	0.164	-0.484	0.304			

432  
 433 <sup>[1]</sup> Communities are simulated by the spatially implicit model described in supplementary text C. Accuracy and  
 434 precision coefficients are derived from concordance correlation coefficient according to Lin et al. (31, 32).  
 435 <sup>[2]</sup> Stochasticity indexes: *NST*, normalized stochasticity ratio; *ST*, stochasticity ratio; *NP*, abundance-based or  
 436 incidence-based percentage of species fitting neutral model.  
 437 <sup>[3]</sup> The incidence-based similarity metrics Canberra, modified Gower (mGower), Cao, and modified Manhattan  
 438 showed exactly the same results as Jaccard metric. The incidence-based Morista-Horn metric showed the same  
 439 results as Sørensen metric. The incidence-based Binomial metrics showed the same results as Gower metric. See  
 440 Table S3 for the detailed definition of each similarity metric.  
 441 <sup>[4]</sup> “(S)” means the metrics need to be standardized as described in Supplementary text A before applied to *ST* and  
 442 *NST*.  
 443 <sup>[5]</sup> mEuclidean, mGower, and mManhattan indicate modified Euclidean, Gower, and Manhattan indexes,  
 444 respectively.  
 445 <sup>[6]</sup> *NP* does not depend on similarity metrics at all, thus only has one value here.  
 446



447  
448

**Table S3. List of similarity and dissimilarity metrics.**

No.	Methods	Refs	Dissimilarity ( $D$ )				Similarity ( $C$ )	
			Qualitative ( $D_{uw}$ )		Quantitative ( $D_w$ )		Formula	Upper limit
			Formula	Upper limit	Formula	Upper limit		
<b>Taxonomic measures</b>								
1	Jaccard & Ružička	(1, 2)	$\frac{A + B - 2J}{A + B - J}$	1	$\frac{\sum_k  x_{ik} - x_{jk} }{\sum_k \max\{x_{ik}, x_{jk}\}}$	1	$1 - D$	1
2	Sørensen & Bray-Curtis	(3, 4)	$\frac{A + B - 2J}{A + B}$	1	$\frac{\sum_k  x_{ik} - x_{jk} }{\sum_k (x_{ik} + x_{jk})}$	1	$1 - D$	1
3	Kulczynski	(5)	$1 - \frac{1}{2} \cdot \left( \frac{J}{A} + \frac{J}{B} \right)$	1	$1 - \frac{1}{2} \cdot \left( \frac{\sum_k \min\{x_{ik}, x_{jk}\}}{\sum_k x_{ik}} + \frac{\sum_k \min\{x_{ik}, x_{jk}\}}{\sum_k x_{jk}} \right)$	1	$1 - D$	1
4	Canberra	(6)	$\frac{A + B - 2J}{A + B - J}$	1	$\frac{1}{A + B - J} \cdot \sum_k \frac{ x_{ik} - x_{jk} }{x_{ik} + x_{jk}}$	1	$1 - D$	1
5	Gower	(7)	$\frac{A + B - 2J}{M}$	1	$\frac{1}{M} \cdot \sum_k \frac{ x_{ik} - x_{jk} }{\max\{x_k\} - \min\{x_k\}}$	1	$1 - D$	1
6	Modified Gower	(8)	$\frac{A + B - 2J}{A + B - J}$	1	$\frac{\sum_k  x'_{ik} - x'_{jk} }{A + B - J}$ where $x'_{ik} = \log_{10}(x_{ik}) + 1$ , unless $x_{ik} = 0$ , in which $x'_{ik} = 0$	Unfix.	$C_{uw} = 1 - D_{uw}$	1
							$C_w = \frac{\sum_i \min\{x'_{ij}, x'_{ik}\}}{A + B - J}$	Unfix.
7	Morisita	(9)	NA.	-	$1 - \frac{2 \sum_k x_{ik} x_{jk}}{(\lambda_i + \lambda_j) \sum_k x_{ik} \sum_k x_{jk}}$ where $\lambda_i = \frac{\sum_k [x_{ik}(x_{ik}-1)]}{\sum_k x_{ik}[(\sum_k x_{ik})-1]}$	1	$1 - D$	1
8	Morisita-Horn	(10)	$\frac{A + B - 2J}{A + B}$	1	$1 - \frac{2 \sum_k x_{ik} x_{jk}}{(\lambda'_i + \lambda'_j) \sum_k x_{ik} \sum_k x_{jk}}$ where $\lambda'_i = \frac{\sum_k x_{ik}^2}{(\sum_k x_{ik})^2}$	1	$1 - D$	1

449

450 **Table S3. Continued**

No.	Methods	Refs	Dissimilarity ( <i>D</i> )				Similarity ( <i>C</i> )	
			Qualitative ( <i>D<sub>uw</sub></i> )		Quantitative ( <i>D<sub>w</sub></i> )		Formula	Upper limit
			Formula	Upper limit	Formula	Upper limit		
9	Manhattan	(11)	$A + B - 2J$	Unfix.	$\sum_k  x_{ik} - x_{jk} $	Unfix.	NA.	-
10	Modified Manhattan	(8)	$\frac{A + B - 2J}{A + B - J}$	1	$\frac{\sum_k  x_{ik} - x_{jk} }{A + B - J}$	Unfix.	NA.	-
11	Euclidean	(11)	$\sqrt{A + B - 2J}$	Unfix.	$\sqrt{\sum_k (x_{ik} - x_{jk})^2}$	Unfix.	NA.	-
12	Modified Euclidean	(8)	$\frac{\sqrt{A + B - 2J}}{A + B - J}$	Unfix.	$\frac{\sqrt{\sum_k (x_{ik} - x_{jk})^2}}{A + B - J}$	Unfix.	NA.	-
13	Binomial	(12)	$(A + B - 2J)\log 2$	Unfix.	$\sum_k \left[ \frac{x_{ik}}{x_{ik} + x_{jk}} \log \left( \frac{x_{ik}}{x_{ik} + x_{jk}} \right) + \frac{x_{jk}}{x_{ik} + x_{jk}} \log \left( \frac{x_{jk}}{x_{ik} + x_{jk}} \right) - \log \frac{1}{2} \right]$	Unfix.	NA.	-
14	Chao	(13)	NA.	-	$1 - \frac{U_i U_j}{U_i + U_j - U_i U_j}$ where $U_i = \frac{c_i}{N_i} + \frac{N_j - 1}{N_j} \cdot \frac{q_1}{2q_2} \cdot \frac{s_{1i}}{N_i}$ , similar for $U_j$	1	$1 - D$	1
15	Cao	(14)	$\frac{A + B - 2J}{A + B - J} \rho$ where $\rho = 1.4954$ if using natural logarithms	$\rho$	$\frac{\sum_k \left( \log \left( \frac{x_{ik} + x_{jk}}{2} \right) - \frac{x_{ik} \log x_{jk} + x_{jk} \log x_{ik}}{x_{ik} + x_{jk}} \right)}{A + B - J}$ where if $x_{ik} = 0$ or $x_{jk} = 0$ , 0.1 is assigned	Unfix.	$1 - \frac{D}{\max\{D\}}$	1

451  
452

No.	Methods	Refs	Dissimilarity ( <i>D</i> )				Similarity ( <i>C</i> )	
			Qualitative ( <i>D<sub>uv</sub></i> )		Quantitative ( <i>D<sub>w</sub></i> )		Formula	Upper limit
			Formula	Upper limit	Formula	Upper limit		
<b>Phylogenetic measures</b>								
16	Phylogenetic analogue of Jaccard & Ružička	(15-17)	$\frac{a + b - 2c}{a + b - c}$ called Unifrac	1	$\sum_n  p_{in} - p_{jn}  W_n$ called weighted Unifrac	Unfix.	NA.	-
					$\frac{\sum_n  p_{in} - p_{jn}  W_n}{\sum_n \max(p_{in}, p_{jn}) W_n}$ Unnamed	1	1 - <i>D</i>	1
17	Phylogenetic analogue of Sørensen & Bray-Curtis	(16, 18)	$\frac{a + b - 2c}{a + b}$	1	$\frac{\sum_n  p_{in} - p_{jn}  W_n}{\sum_n (p_{in} + p_{jn}) W_n}$ called normalized weighted Unifrac	1	1 - <i>D</i>	1
18	$\beta$ MPD	(19)	$\frac{1}{AB} \sum_{k=1}^A \sum_{m=1}^B \delta_{km}$	Unfix.	$\frac{\sum_{k=1}^A \sum_{m=1}^B p_{ik} p_{jm} \delta_{km}}{\sum_{k=1}^A \sum_{m=1}^B p_{ik} p_{jm}}$	Unfix.	NA.	-
19	$\beta$ MNTD	(19, 20)	$\frac{1}{2} \left[ \frac{\sum_{k=1}^A \min(\delta_{km})}{A} + \frac{\sum_{m=1}^B \min(\delta_{km})}{B} \right]$	Unfix.	$\frac{1}{2} \left[ \sum_{k=1}^A p_{ik} \min(\delta_{km}) + \sum_{m=1}^B p_{jm} \min(\delta_{km}) \right]$	Unfix.	NA.	-

454 <sup>[1]</sup> *A* is the richness (number of taxa) in community *i*, while *B* is the richness in sample *j*, and *J* is the number of taxa that occur on both sample *i* and *j*.  
 455 <sup>[2]</sup> *x<sub>ik</sub>* is the abundance of taxon *k* in sample *i*, while *x<sub>jk</sub>* is the abundance of taxon *k* in sample *j*.  
 456 <sup>[3]</sup> *p<sub>ik</sub>* is the proportion of taxon *k* in sample *i*, while *p<sub>im</sub>* is the proportion of taxon *m* in sample *j*.  
 457 <sup>[4]</sup> {*x<sub>k</sub>*} is the set of abundances of taxon *k* in all samples.  
 458 <sup>[5]</sup> *M* is the number of taxa in all samples.  
 459 <sup>[6]</sup> For Chao index, *C<sub>i</sub>* is the total number of individuals in the taxa of sample *i* that are shared with sample *j*; *N<sub>i</sub>* is the total number of individuals in sample *i*,  
 460 *N<sub>j</sub>* is the total number of individuals in sample *j*; *q1* (and *q2*) are the number of species occurring in sample *i* that have only one (or two) individuals in  
 461 sample *j*; *s1<sub>i</sub>* is the total number of individuals in the species present in sample *i* that occur with only one individual in sample *j*.  
 462 <sup>[7]</sup> *a* is the amount of phylogenetic tree branch length in community *i*, *b* is the amount of branch length in community *j*, and *c* is the amount of branch length  
 463 shared between community *i* and *j*.  
 464 <sup>[8]</sup> *p<sub>in</sub>* is the proportion of sequences (taxa) from community *i* descendant from branch *n*; *W<sub>n</sub>* is the weight or length of branch *n*.  
 465 <sup>[9]</sup>  $\delta_{km}$  is the phylogenetic distance from sequence (taxon) *k* to sequence (taxon) *m*.  
 466 <sup>[10]</sup> *X* is the set of taxa in community *i*, while *Y* is the set of taxa in community *k*.

467  
468

**Table S4. Summary of null model algorithms for species co-occurrence analysis.**

No.	Abbreviation in this paper	Abbreviation in Gotelli (33)	Ways to constrain taxa occurrence frequency <sup>[1]</sup>	Ways to constrain richness in each sample <sup>[2]</sup>	Probability of taxon <i>i</i> present in sample <i>j</i> <sup>[3]</sup>
1	EE	SIM1	Equiprobable	Equiprobable	$P_{ij} = \frac{1}{N} \cdot \frac{1}{M}$
2	EP	SIM6	Equiprobable	Proportional	$P_{ij} = \frac{1}{N} \cdot \frac{S_j}{F}$
3	EF	SIM3	Equiprobable	Fixed	$P_{ij} = \frac{1}{N}$
4	PE	SIM7	Proportional	Equiprobable	$P_{ij} = \frac{f_i}{F} \cdot \frac{1}{M}$
5	PP	SIM8	Proportional	Proportional	$P_{ij} = \frac{f_i}{F} \cdot \frac{S_j}{F}$
6	PF	SIM5	Proportional	Fixed	$P_{ij} = \frac{f_i}{F}$
7	FE	SIM2	Fixed	Equiprobable	$P_{ij} = \frac{1}{M}$
8	FP	SIM4	Fixed	Proportional	$P_{ij} = \frac{S_j}{F}$
9	FF	SIM9	Fixed	Fixed	Not applicable

469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480

- <sup>[1]</sup> As to occurrence frequency, “Equiprobable” means that all taxa have equal probability to occur; “Proportional” means that the occurrence probability of a taxon is proportional to its observed occurrence frequency; “Fixed” means that the occurrence frequency of a taxon is fixed as observed.
- <sup>[2]</sup> As to species richness in each sample, “Equiprobable” means that all samples have equal probability to contain a taxon; “Proportional” means the occurrence probability in a sample is proportional to the observed richness in this sample; “Fixed” means the occurrence frequency of a taxon is fixed as observed.
- <sup>[3]</sup>  $P_{ij}$  is the probability of taxon *i* present in sample *j* in a null model.  
 $S_j$  is the observed richness in sample *j*,  $N$  is the total number of taxa,  $M$  is the total number of samples.  
 $f_i$  is the observed occurrence frequency of taxon *i*,  $F$  is the total number of occurrences.  
 $A_i$  is the regional abundance of taxon *i*,  $J$  is the total abundance of all taxa in all samples.

481 **Supplementary References**

- 482 1. Jaccard P (1912) The distribution of the flora in the alpine zone. *New Phytol* 11(2):37-50.  
 483 2. Ružička M (1958) Anwendung mathematisch-statistischer methoden in der geobotanik  
 484 (Synthetische bearbeitung von aufnahmen). *Biológia, Bratislava* 13:647–661.  
 485 3. Sørensen T (1948) A method of establishing groups of equal amplitude in plant sociology  
 486 based on similarity of species content. *Kongelige Danske Videnskabernes Selskab.*  
 487 *Biologiske Skrifter.* 4(1-34).  
 488 4. Bray JR & Curtis JT (1957) An Ordination of the Upland Forest Communities of Southern  
 489 Wisconsin. *Ecol Monogr* 27(4):326-349.  
 490 5. Kulczynski S (1928) Die Pflanzenassoziationen der Pieninen. *Bulletin de l'Académie*  
 491 *polonaise des sciences. Série des sciences biologiques.*  
 492 6. Lance GN & Williams WT (1966) Computer Programs for Hierarchical Polythetic  
 493 Classification (“Similarity Analyses”). *Comput J* 9(1):60-64.  
 494 7. Gower JC (1971) A General Coefficient of Similarity and Some of Its Properties.  
 495 *Biometrics* 27(4):857-871.  
 496 8. Anderson MJ, Ellingsen KE, & McArdle BH (2006) Multivariate dispersion as a measure  
 497 of beta diversity. *Ecol Lett* 9(6):683-693.  
 498 9. Morisita M (1959) Measuring of the dispersion and analysis of distribution patterns.  
 499 *Memoires of the Faculty of Science, Kyushu University, Series E. Biology.* 2:215-235.  
 500 10. Horn HS (1966) Measurement of "Overlap" in Comparative Ecological Studies. *Am Nat*  
 501 100(914):419-424.  
 502 11. Krebs CJ (1999) *Ecological Methodology* (A. Wesley Longman, NY, USA) 2 Ed.  
 503 12. Anderson MJ & Millar RB (2004) Spatial variation and effects of habitat on temperate reef  
 504 fish assemblages in northeastern New Zealand. *J Exp Mar Bio Ecol* 305(2):191-221.  
 505 13. Chao A, Chazdon RL, Colwell RK, & Shen T-J (2005) A new statistical approach for  
 506 assessing similarity of species composition with incidence and abundance data. *Ecol Lett*  
 507 8(2):148-159.  
 508 14. Cao Y, Williams WP, & Bark AW (1997) Similarity measure bias in river benthic  
 509 aufwuchs community analysis. *Water Environ Res* 69(1):95-106.  
 510 15. Lozupone C & Knight R (2005) UniFrac: a new phylogenetic method for comparing  
 511 microbial communities. *Appl Environ Microbiol* 71(12):8228-8235.  
 512 16. Lozupone CA, Hamady M, Kelley ST, & Knight R (2007) Quantitative and qualitative  $\beta$   
 513 diversity measures lead to different insights into factors that structure microbial  
 514 communities. *Appl Environ Microbiol* 73(5):1576-1585.  
 515 17. Parks DH & Beiko RG (2013) Measures of phylogenetic differentiation provide robust and  
 516 complementary insights into microbial communities. *ISME J* 7(1):173-183.  
 517 18. Bryant JA, *et al.* (2008) Microbes on mountainsides: Contrasting elevational patterns of  
 518 bacterial and plant diversity. *Proc Natl Acad Sci U S A* 105:11505-11511.  
 519 19. Webb CO, Ackerly DD, & Kembel SW (2008) Phylocom: software for the analysis of  
 520 phylogenetic community structure and trait evolution. *Bioinformatics* 24(18):2098-2100.  
 521 20. Stegen JC, Lin X, Konopka AE, & Fredrickson JK (2012) Stochastic and deterministic  
 522 assembly processes in subsurface microbial communities. *ISME J* 6:1653–1664.  
 523 21. Oksanen J, *et al.* (2017) *vegan: Community Ecology Package*.  
 524 22. Afgan E, *et al.* (2018) The Galaxy platform for accessible, reproducible and collaborative  
 525 biomedical analyses: 2018 update. *Nucleic Acids Res* 46(W1):W537-W544.

- 526 23. Alonso D & McKane AJ (2004) Sampling Hubbell's neutral theory of biodiversity. *Ecol*  
527 *Lett* 7(10):901-910.
- 528 24. Hubbell SP (2001) *The unified neutral theory of biodiversity and biogeography* (Princeton  
529 University Press, Princeton, NJ, USA) p 375.
- 530 25. Prado PI, Miranda MD, & Chalom A (2017) sads: Maximum Likelihood Models for  
531 Species Abundance Distributions. *R package version 0.4.1*.
- 532 26. Sloan WT, *et al.* (2006) Quantifying the roles of immigration and chance in shaping  
533 prokaryote community structure. *Environ Microbiol* 8(4):732-740.
- 534 27. Burns AR, *et al.* (2016) Contribution of neutral processes to the assembly of gut microbial  
535 communities in the zebrafish over host development. *ISME J* 10(3):655-664.
- 536 28. Chase JM, Kraft NJB, Smith KG, Vellend M, & Inouye BD (2011) Using null models to  
537 disentangle variation in community dissimilarity from variation in alpha-diversity.  
538 *Ecosphere* 2(2).
- 539 29. Stegen JC, *et al.* (2013) Quantifying community assembly processes and identifying  
540 features that impose them. *ISME J* 7:2069-2079.
- 541 30. Kraft NJB, *et al.* (2011) Disentangling the drivers of beta diversity along latitudinal and  
542 elevational gradients. *Science* 333(6050):1755-1758.
- 543 31. Lin LI (1989) A concordance correlation-coefficient to evaluate reproducibility.  
544 *Biometrics* 45(1):255-268.
- 545 32. Lin L, Hedayat AS, Sinha B, & Yang M (2002) Statistical methods in assessing agreement:  
546 Models, issues, and tools. *J Am Stat Assoc* 97(457):257-270.
- 547 33. Gotelli NJ (2000) Null model analysis of species co-occurrence patterns. *Ecology*  
548 81(9):2606-2621.
- 549 34. Zhang P, *et al.* (2017) Dynamic succession of groundwater sulfate-reducing communities  
550 during prolonged reduction of uranium in a contaminated aquifer. *Environ Sci Technol*  
551 51(7):3609-3620.
- 552 35. He Z, *et al.* (2018) Microbial functional gene diversity predicts groundwater contamination  
553 and ecosystem functioning. *MBio* 9(1):e02435-02417.  
554