

1 **Title:** Maximum entropy models elucidate the contribution of metabolic traits to patterns of
2 community assembly

3 **Running head:** Metabolic traits in maximum entropy models

4

5 **Authors:**

6 Jason Bertram^{1*}, Erica A. Newman^{2,3}, Roderick C. Dewar⁴

7 ¹Environmental Resilience Institute, Indiana University, Bloomington, IN 47401 USA

8 ²School of Natural Resources and the Environment, University of Arizona, Tucson, AZ 85721

9 USA

10 ³USDA Forest Service, Pacific Wildland Fire Sciences Lab, Seattle, WA 98103 USA

11 ⁴Plant Sciences Division, Research School of Biology, The Australian National University,

12 Canberra, ACT 2601 Australia

13 *Corresponding author email: jxb@iu.edu

14

15 **Number of words in the Abstract: 268**

16 **Number of words in main text (including appendices): 4323**

17 **Number of references: 33**

18

19 **ACKNOWLEDGEMENTS**

20 Funding for JB was provided by the Environmental Resilience Institute. Postdoctoral funding for

21 EAN was provided by University of Arizona Bridging Biodiversity and Conservation Science

22 program, and the USDA Forest Service. We thank John Harte for comments on earlier versions

23 of this manuscript, and Dan Gruner for use of his dataset.

24

25 **BIOSKETCHES**

26 **Jason Bertram** is a mathematical biologist who develops ecological and evolutionary theory at
27 Indiana University.

28 **Erica A. Newman** is postdoctoral researcher focusing on connections between disturbance
29 ecology and biodiversity patterns, including applied topics in wildlife ecology and wildfire
30 management.

31 **Roderick C. Dewar** is a professor at the Australian National University, researching entropy-
32 based principles of organization in biological and physical systems.

33

34

35 **ABSTRACT**

36 **Aim**

37 Maximum entropy (MaxEnt) models promise a novel approach for understanding community
38 assembly and species abundance patterns. One of these models, the “Maximum Entropy Theory
39 of Ecology” (METE) reproduces many observed species abundance patterns, but is based on an
40 aggregated representation of community structure that does not resolve species identity or
41 explicitly represent species-specific functional traits. In this paper, METE is compared to “Very
42 Entropic Growth” (VEG), a MaxEnt model with a less aggregated representation of community
43 structure that represents species (more correctly, functional types) in terms of their per capita
44 metabolic rates. We examine the contribution of metabolic traits to the patterns of community
45 assembly predicted by VEG and, through aggregation, compare the results with METE
46 predictions in order to gain insight into the biological factors underlying observed patterns of
47 community assembly.

48

49 **Innovation**

50 We formally compare two MaxEnt-based community models, METE and VEG, that differ as to
51 whether or not they represent species-specific functional traits. We empirically test and compare
52 the metabolic predictions of both models, thereby elucidating the role of metabolic traits in
53 patterns of community assembly.

54

55 **Main Conclusions**

56 Our analysis reveals that a key determinant of community metabolic patterns is the “density of
57 species” distribution $\rho(\epsilon)$, where $\rho(\epsilon)d\epsilon$ is the intrinsic number of species with metabolic rates
58 in the range $(\epsilon, \epsilon + d\epsilon)$ that are available to a community prior to filtering by environmental

59 constraints. Our analysis suggests that appropriate choice of $\rho(\epsilon)$ in VEG may lead to more
60 realistic predictions than METE, for which $\rho(\epsilon)$ is not defined, and thus opens up new ways to
61 understanding the link between functional traits and patterns of community assembly.

62

63 **Key words:** community assembly, functional traits, macroecology, metabolic requirements,
64 resource partitioning, species-abundance distribution, statistical aggregation

65

66 1 INTRODUCTION

67 One of the central aims of ecology is to understand the determinants of community
68 assembly. Many studies of community assembly involve summaries of community structure such
69 as the species abundance distribution (SAD), species-area relationship (SAR), and analogous
70 metabolic-rate distributions. We will refer to these summary distributions collectively as
71 community structure distributions (CSDs). CSDs, particularly the SAD, have attracted a lot of
72 attention because their shapes are strikingly similar across different communities, representing a
73 rare example of “universality” in community ecology (McGill *et al.* 2007).

74 The existence of universal features in CSDs is intriguing because these could reflect
75 universal aspects of the biological processes responsible for structuring communities. However,
76 CSDs could also be universal for statistical reasons (Tokeshi 1993; Ulrich *et al.*, 2010). Similar
77 to how the normal distribution is ubiquitous because many measured quantities involve statistical
78 averaging (the central limit theorem), CSDs could be universal simply because community-
79 specific details disappear in aggregating patterns to the level of species counts, or other forms of
80 averaging. This would make CSDs considerably less valuable for understanding the biological
81 determinants of community assembly, such as how community structure depends on the

82 functional traits of the organisms in the community (McGill *et al.*, 2006; Díaz *et al.* 2013). It is
83 therefore important to disentangle the contributions of biological versus statistical factors to
84 CSDs. This issue is closely related to the long-running debate on the relative roles of
85 “mechanism” and “drift” in ecology (McGill and Nekola, 2010; Vellend, 2010), and on
86 ecosystem stability and the role of disturbance (Newman *et al.* 2018).

87 A promising recent approach for disentangling biological from statistical factors in
88 ecological models is to use the statistical principle of maximum entropy (MaxEnt). MaxEnt
89 models are “top-down” in that they seek to identify a minimal set of biological assumptions
90 required to reproduce a given empirical pattern (such as a CSD). Once these assumptions have
91 been specified, the MaxEnt principle predicts statistical patterns of community structure by
92 effectively treating all other mechanistic details statistically as unbiased random noise. By
93 empirically testing predictions based on different assumptions, MaxEnt provides a means to
94 resolve the partitioning of biological versus statistical factors in driving observed ecological
95 patterns.

96 MaxEnt models have had some success at predicting CSDs, but the ecological
97 interpretation of these successes has not been straightforward. A number of MaxEnt models have
98 appeared in the ecological literature with a variety of different assumptions and justifications
99 (*e.g.* Shipley, Vile and Garnier, 2006; Pueyo, He and Zillio, 2007; Harte, Zillio, Conlisk and
100 Smith, 2008; Dewar and Porté, 2008; Banavar, Maritan and Volkov, 2010; Bertram and Dewar,
101 2015). This had led to extensive debates about the prospects and pitfalls of this approach.

102 Two key issues in these debates may be identified: one conceptual, the other more
103 technical. The conceptual issue concerns the interpretation of the MaxEnt procedure itself,
104 including the challenge of connecting MaxEnt to familiar ecological processes such as dispersal,

105 disturbance, and interactions between organisms. This issue has been discussed at length
106 elsewhere (Dewar, 2009; McGill and Nekola, 2010; Shipley, 2010; Supp, Xiao, Ernest and
107 White, 2012; Harte and Newman, 2014; Supp and Ernest 2014; Bertram and Dewar, 2015;
108 Newman et al., 2018), and will thus not be our focus here.

109 Rather, our focus will be on the technical issue, which concerns the level of detail at
110 which the community is described in the model before MaxEnt is even applied (He, 2010;
111 Favretti, 2017). Changing the variables used to describe a community (*e.g.* resolving a
112 community in greater detail) can dramatically alter the predictions that MaxEnt makes about the
113 community, and yet there is no apparent *a priori* reason to prefer one choice of variables over
114 another. As a result, a variety of choices have appeared in different models, usually with little
115 justification. A comparison of these disparate approaches is required in order to better guide the
116 application of MaxEnt models in ecology.

117 Here we present a comparison of two MaxEnt-based models in ecology which have both
118 successfully reproduced observed CSDs: METE (Maximum Entropy Theory of Ecology; Harte
119 et al. 2008; Harte et al. 2009; Harte 2011; Harte and Newman, 2014) and VEG (Very Entropic
120 Growth; Dewar and Porté, 2008; Bertram and Dewar, 2013; Bertram and Dewar 2015). These
121 models are well suited for our objective of comparison because METE describes communities at
122 the same coarse-grained level of detail as the SAD, whereas VEG is more detailed in that it
123 resolves the abundance of each separate species.

124 Crucially, this difference in community description allow us to explore the biological
125 determinants of patterns of community assembly. Specifically, VEG distinguishes species by
126 their per capita metabolic traits. By contrast, METE only distinguishes separate species by their
127 abundances, and requires the total number of species present in the community as an input rather

128 than as a prediction. METE then predicts a distribution of metabolic rates for the individuals in a
129 species as a function of its abundance, imparting functional traits statistically by abundance,
130 rather than by species identity (Section 2.1). It is therefore not possible to investigate how
131 different species-specific functional traits might modify community structure using METE. Thus,
132 by comparing METE and VEG, we are able to investigate more transparently what sorts of
133 functional trait assumptions are necessary for reproducing observed patterns.

134 2. METE AND VEG: TWO MAXENT MODELS OF COMMUNITY ASSEMBLY

135 2.1 METE

136 The central quantity predicted by METE is a joint probability distribution $R_M(n, \epsilon)$ called
137 the “ecosystem structure function.” By definition, $R_M(n, \epsilon)d\epsilon$ is the joint probability that a
138 species selected at random from a community has abundance n , and that an individual selected at
139 random from a species with abundance n has a metabolic requirement between ϵ and $\epsilon + d\epsilon$
140 (Harte et al. 2008; Harte, 2011; Appendix A). The ecosystem structure function is closely related
141 to the SAD: if we add together all of the possible metabolic requirements ϵ we obtain the
142 probability distribution for the abundance of a randomly selected species, $R_M(n) =$
143 $\int R_M(n, \epsilon)d\epsilon$; the SAD is then simply $SR_M(n)$ where S is the total number of species present in
144 the community. Thus, $R_M(n, \epsilon)$ is a SAD that has been extended to also incorporate information
145 about community metabolic structure.

146 METE assumes that $R_M(n, \epsilon)$ satisfies two constraints

$$147 \quad \sum_{n=1}^N \int_{\epsilon=1}^E n R_M(n, \epsilon) d\epsilon = N/S \quad (1)$$

$$148 \quad \sum_{n=1}^N \int_{\epsilon=1}^E n \epsilon R_M(n, \epsilon) d\epsilon = E/S. \quad (2)$$

149 In words, these constraints say that the total number of individuals in the community

150 $S \sum_{n=1}^N \int_{\epsilon=1}^E n R_M(n, \epsilon) d\epsilon$ is equal to N , and the total community metabolic requirement

151 $S \sum_{n=1}^N \int_{\epsilon=1}^E n \epsilon R_M(n, \epsilon) d\epsilon$ is equal to E . $R_M(n, \epsilon)$ is then obtained by maximizing the Shannon
152 entropy $-\sum_n \int R_M \ln R_M d\epsilon$ subject to constraints (1) and (2), as well the constraint that
153 $R_M(n, \epsilon)$ sums to 1 (since it is a probability distribution). This maximization procedure gives

154
$$R_M(n, \epsilon) \propto e^{-\lambda_1 n - \lambda_2 n \epsilon} \quad (3)$$

155 where λ_1 and λ_2 are constants (Lagrange multipliers) with values chosen such that constraints (1)
156 and (2) hold (for details, see Harte, 2011). The triplet of values N , E and S are the inputs to
157 METE (note that we will not consider the area-scaling component of the full METE theory;
158 Harte, 2011).

159

160 **2.2 VEG**

161 VEG is similar to METE in that it uses MaxEnt to infer community properties from a few
162 constraints. Moreover, the VEG constraints are similar to METE's (see below). The major
163 feature that differentiates VEG is that it represents community structure in more detail. In VEG,
164 species are distinguishable, whereas METE only specifies the proportion of species with each
165 abundance n via the ecosystem structure function (Fig. 1).

166 **[FIGURE 1]**

167 In contrast to METE (which uses MaxEnt to infer the ecosystem structure function
168 directly), VEG uses MaxEnt to predict the probability $p(\mathbf{n})$ that, when we take a snapshot of the
169 community, we observe the species abundances $\mathbf{n} = (n_1, n_2, \dots)$ (*i.e.* the species labeled 1 has
170 abundance n_1 , and so on). In VEG, species' abundances may be zero; the number of species
171 actually present in a snapshot is the number of nonzero elements of \mathbf{n} . VEG therefore predicts
172 probabilities for the abundance of each species separately; consequently, VEG also predicts the
173 expected number species that are present in the community. Species in VEG are also assigned

174 distinct functional traits: the individuals of species i are assumed to have a metabolic
175 requirement of ϵ_i , where the species labels are chosen such that $\epsilon_1 \leq \epsilon_2 \leq \epsilon_3 \dots$, and so on.

176 Similar to METE, VEG assumes total abundance and total metabolic requirement
177 constraints

$$178 \quad \sum_{\mathbf{n}} \sum_i n_i p(\mathbf{n}) = N \quad (4)$$

$$179 \quad \sum_{\mathbf{n}} \sum_i n_i \epsilon_i p(\mathbf{n}) = E \quad (5)$$

180 Since $p(\mathbf{n})$ represents the probability of observing the “snapshots” \mathbf{n} , the probabilities can be
181 interpreted as sample frequencies representing the proportion of time that the community spends
182 with different abundance compositions \mathbf{n} . Consequently, constraints (4) and (5) have a clear
183 ecological interpretation in VEG as fixing the time-averaged total abundance and total metabolic
184 requirement of the community to have the values N and E respectively; the latter can be
185 interpreted as an expression of the long-term steady-state ecological balance between resource
186 use (left-hand side of Eq. (5)) and supply (E , right-hand side of Eq. (5)). In contrast, the METE
187 constraints (Eqs. (1) and (2)), which are statements about “information”, do not have a similarly
188 straightforward ecological interpretation.

189 Again similarly to METE, $p(\mathbf{n})$ is obtained by maximizing the Shannon entropy
190 $-\sum_{\mathbf{n}} p(\mathbf{n}) \ln p(\mathbf{n})$ subject to constraints (4), (5), and the constraint that $p(\mathbf{n})$ sums to 1. This
191 maximization procedure gives

$$192 \quad p(\mathbf{n}) \propto e^{-\sum_i (\mu_1 + \mu_2 \epsilon_i) n_i} \quad (6)$$

193 where μ_1 and μ_2 are the Lagrange multipliers corresponding to constraints (4) and (5)
194 respectively. Note that in Eq. (6), $p(\mathbf{n})$ depends on the spectrum of metabolic requirements
195 present in the community $\epsilon_1 \leq \epsilon_2 \leq \dots$. Thus the inputs of VEG are N , E and the spectrum of

196 values ϵ_i . In contrast to METE, the number of species (S) present in the community is an output
197 of VEG, rather than an input.

198 **3 COMPARING METE AND VEG**

199 **3.1 THE VEG ECOSYSTEM STRUCTURE FUNCTION**

200 In this section we give an intuitive derivation of the ecosystem structure function implied
201 by VEG, which will be denoted $R_V(n, \epsilon)$ (a more rigorous mathematical derivation is given in
202 Appendix B). This will allow us to directly compare the predictions of METE and VEG.

203 When we sample a species at random from a community, all species present have the
204 same probability of being selected. However, the metabolic requirement ϵ of the selected species
205 is more likely to take some values than others due to two effects: (1) *Trait availability*. Among
206 the species currently inhabiting the community's broader geographic region, some values of ϵ
207 are more likely to occur than others due to intrinsic biophysical constraints on the traits
208 determining ϵ , and the region's evolutionary history; (2) *Environmental filtering* (Shipley *et al.*,
209 2006). From the distribution of possible metabolic rates, some values of ϵ are more likely to be
210 actually present in the community due to additional bias imposed by local environmental
211 constraints (such as Eqs. (4) and (5)).

212 VEG represents a special case in which there is no trait variation within species: all
213 individuals in species i have the same metabolic requirement ϵ_i (thus a VEG "species" is more
214 appropriately interpreted as a functional type rather than a taxonomic unit; Bertram and Dewar,
215 2013). Thus, the first effect above (trait availability) is represented by the fact that the metabolic
216 spectrum $\epsilon_1 \leq \epsilon_2 \leq \dots$ may be more densely packed at some values of ϵ than at others. To
217 represent this effect mathematically, we introduce the "density of species" distribution $\rho(\epsilon)$;
218 $\rho(\epsilon)d\epsilon$ counts the number of metabolic requirement values ("species") contained in the interval

219 $(\epsilon, \epsilon + d\epsilon)$. For comparison with METE, in which ϵ is a continuous variable, we assume that the
220 metabolic requirement spectrum is sufficiently dense that we can approximate $\rho(\epsilon)$ as a
221 continuous function of ϵ . Intuitively, the shape of $\rho(\epsilon)$ represents the relative probabilities that a
222 species selected at random out of all possible species that could be present in the community has
223 a metabolic requirement within a given interval (Fig. 2).

224 [FIGURE 2]

225 Once a species has been sampled out of all possible species and its metabolic requirement
226 has been found to be ϵ , the probability that it has abundance n , denoted $p(n|\epsilon)$, can then be
227 straightforwardly calculated in VEG from Eq. (6) (from Appendix B, $p(n|\epsilon) \propto e^{-(\mu_1 + \mu_2 \epsilon)n}$).
228 VEG also explicitly accounts for the second effect above (environmental filtering), through the
229 Lagrange multipliers μ_1 and μ_2 that reflect the environmental constraints of Eqs. (4) and (5).

230 To construct $R_V(n, \epsilon)$, which only refers to species that are actually present, we restrict
231 our attention to $n \geq 1$. Thus, the joint probability of sampling a species with abundance n from
232 the community, and an individual from such a species with metabolic requirement ϵ , is
233 proportional to $\rho(\epsilon)p(n|\epsilon)$, where $n \geq 1$. This gives

$$234 \quad R_V(n, \epsilon) \propto \rho(\epsilon)p(n|\epsilon) \quad (7)$$

235 The above argument leading to Eq. (7) for R_V (and the more rigorous argument given in
236 Appendix B) is quite general. It can be applied to obtain the ecosystem structure function for any
237 model in which we know the density of species $\rho(\epsilon)$ (which need not be restricted to a species-
238 specific trait spectrum as in VEG), and which predicts abundance probabilities conditional on the
239 trait values $p(n|\epsilon)$ (whether those probabilities are predicted using MaxEnt or by other means).

240 3.2 SPECIES ABUNDANCE DISTRIBUTIONS

241 A large number of ecological models have reproduced realistic SADs, including METE
242 (Harte et al., 2008) and VEG (Bertram and Dewar, 2015). SAD comparisons consequently only
243 have weak power to discriminate the predictions of different ecological theories (they are “weak
244 tests”; McGill 2003, McGill et al. 2007). In particular, the SAD predictions of METE and VEG
245 will not tell us much about their differences. It is interesting to demonstrate this “weak test”
246 property of SADs explicitly in terms of the METE and VEG ecosystem structure functions.

247 As noted in section 2.1, the SAD is obtained by integrating the ecosystem structure
248 function over ϵ (the SAD is proportional to $R(n) = \int R(n, \epsilon) d\epsilon$). We therefore expect that the
249 SAD will be to some extent insensitive to the exact manner in which $R(n, \epsilon)$ depends on ϵ .

250 In the case of METE, the predicted SAD is almost entirely independent of the value of
251 E/S in the metabolic constraint Eq. (2) for many of the most heavily studied SAD datasets (i.e.
252 $R_M(n)$ is independent of λ_2 ; Harte et al. 2008). This behavior represents the limiting case of
253 large E , corresponding to resource-rich communities. Thus, in many cases of interest, METE
254 produces SADs that are insensitive to the value of E/S (note, however, that the existence of the
255 metabolic constraint Eq. (2) is necessary to get a Fisher log-series form for $R_M(n) =$

256 $\int R_M(n, \epsilon) d\epsilon \propto \frac{e^{-\lambda_1 n}}{n}$).

257 VEG allows us investigate the “weak test” property in greater depth because we can
258 independently change the form of $\rho(\epsilon)$ and check if this appreciably changes the VEG SAD.
259 Suppose for illustrative purposes that the metabolic requirement spectrum has a power law form
260 (Dewar and Porté, 2008)

261
$$\rho(\epsilon) \propto \epsilon^\alpha \quad (8)$$

262 where α is a free parameter. Our motivation for Eq. (8) is to have a simple one-parameter
263 function in which we can control the relative density of species at low versus high ϵ ($\alpha = 0$
264 corresponds to a uniformly spaced spectrum; Fig. 2). Using Eq. (8), it can be shown that

$$265 \quad R_V(n) = \int R_V(n, \epsilon) d\epsilon \propto \frac{e^{-\mu_1 n}}{n^{\alpha+1}} \quad (9)$$

266 for all but the lowest abundance species (see Appendix C). Thus, although the exact quantitative
267 shape of $R_V(n)$ does depend on α (both explicitly in Eq. (9) and implicitly via the fact that μ_1 and
268 μ_2 depend on α), $R_V(n)$ will qualitatively have the familiar “hollow curve” shape (McGill et al.
269 2007) regardless of the particular choice of α . In particular, $\alpha = 0$ gives the Fisher log-series
270 (similar to METE). Thus, since $\rho(\epsilon)$ represents the spectrum of functional traits, we can
271 conclude that the shape of the VEG SAD is only marginally sensitive to the metabolic trait
272 values of the species present.

273 However, recall that VEG predicts the total number of species/functional types S in the
274 community (Sec. 2.2). This predicted S is more sensitive to the assumed metabolic trait values
275 than the SAD shape, and could differ from the observed value of S for given observed values of
276 N and E . By contrast, METE uses the empirically observed value of S to construct the METE
277 SAD.

278 3.3 METABOLIC-RANK DISTRIBUTIONS

279 In this section we compare the metabolic dependence of the two structure functions
280 $R_V(n, \epsilon)$ and $R_M(n, \epsilon)$. We do this in two ways: via the marginal distribution for individual
281 metabolic rates $R(\epsilon) = \sum_n R(n, \epsilon)$, and via the individual-level energy distribution (IED)
282 defined by $\Psi(\epsilon) = \frac{S}{N} \sum_n nR(n, \epsilon)$ (Harte 2011; Newman *et al.*, 2014). $R(\epsilon)$ is the probability
283 that a species sampled at random from the community has metabolic rate ϵ , while $\Psi(\epsilon)$ is the
284 probability that an individual sampled at random from the community has metabolic requirement

285 ϵ . In contrast to the SAD (Section 3.1), $R(\epsilon)$ and $\Psi(\epsilon)$ are both sensitive to the shape of $\rho(\epsilon)$ in
286 VEG. We can thus ask, what shape does $\rho(\epsilon)$ need to be to match metabolic data, and how does
287 this $\rho(\epsilon)$ compare to the predictions of METE?

288 Following Harte et al. (2017), we calculate and plot $\Psi(\epsilon)$ cumulatively such that log
289 metabolic rate appears on the vertical axis, and the horizontal axis is the proportion of the
290 population with metabolic rate greater than or equal to a given ϵ (*i.e.* the rank of the
291 corresponding individual). We assumed a power law spectral density as in Eq. (8), taking α as a
292 free parameter to be fitted, and then minimized the least-squares difference between measured
293 log metabolic rates and the predictions of VEG. We repeated this procedure for the three datasets
294 considered in Harte et al. (2017): Barro Colorado Island trees (Hubbell et al. 2005), Hawaiian
295 island arthropods (Gruner, 2007), and Rocky Mountain subalpine meadow plants (Newman et al.
296 2014).

297 In all three datasets we found values of α that give superior $\Psi(\epsilon)$ fits to METE (bottom
298 three panels of Fig. 3; note the logarithmic horizontal axis). This is no great victory given that we
299 have introduced a free parameter α that is not available to METE, but it confirms that the power
300 law form for $\rho(\epsilon)$ gives plausible metabolic predictions. METE and VEG both track the middle
301 and higher ranks closely, but at lower ranks the VEG metabolic rates are too low whereas the
302 METE predictions are too high. The corresponding marginal metabolic distributions $R_M(\epsilon)$ and
303 $R_V(\epsilon)$ (upper panels in Fig. 3) confirm that METE assigns higher probabilities to the highest
304 values of ϵ ($R_M(\epsilon)$ has a longer tail).

305

306 **[FIGURE 3]**

307

308 4 DISCUSSION

309 A key insight of the above analysis is that the ecosystem structure $R(n, \epsilon)$ is sensitive to
310 the shape of the density of species distribution represented mathematically by $\rho(\epsilon)$. In the case
311 of VEG, Eq. (7) implies $R_V(\epsilon) \propto \rho(\epsilon) \sum_{n \geq 1} p(n|\epsilon) = \rho(\epsilon)[1 - p(0|\epsilon)]$. This expression clearly
312 shows the two effects introduced at the start of Sec. 3.1: $R(\epsilon)$ is the density of species $\rho(\epsilon)$
313 multiplied by the probability $1 - p(0|\epsilon)$ that a species with metabolic rate ϵ is actually present
314 in the community.

315 Whereas VEG requires us to specify the form of $\rho(\epsilon)$, METE infers $R_M(n, \epsilon)$ using only
316 MaxEnt and the constraint equations (1) and (2). In this sense METE is a null model for the
317 contribution of functional traits to community patterns, treating functional traits as “random
318 noise” within the community constraints imposed by S , N , and E . However, METE only infers
319 the trait distribution as would be observed in already-assembled communities. METE refers only
320 to species that are already present in the community, and does not give an expression for $p(0|\epsilon)$;
321 it is therefore not possible to compute the density of species $\rho(\epsilon)$ implicitly inferred by METE.
322 Nonetheless, observed ecological communities generally have a large proportion of individuals
323 with low metabolic requirement. This implies $p(0|\epsilon) \approx 0$ and thus $R(\epsilon) \approx \rho(\epsilon)$ for low ϵ (see
324 the convergence of $R(\epsilon)$ and $\rho(\epsilon)$ in VEG in the upper panels of Fig. 3), giving us a glimpse of
325 the $\rho(\epsilon)$ predictions of METE.

326 VEG explicitly separates the trait values that are possible from the trait values that are
327 actually observed post-assembly. Since $\rho(\epsilon)$ is an input, VEG represents an explicit model for
328 the contribution of functional traits to CSDs. This begs the question of what then determines
329 $\rho(\epsilon)$ as the appropriate choice in VEG. There are at least two answers:

- 330 (i) On short timescales, $\rho(\epsilon)$ may simply express the mix of potential species that are
331 available to the community at any given time, as in biodiversity manipulation
332 experiments where a given restricted set of species is thrown together and left to self-
333 organize. This short-term $\rho(\epsilon)$ could be highly contingent on the community's recent
334 history, and could have a strong effect on metabolic patterns following disturbance.
- 335 (ii) On longer timescales, $\rho(\epsilon)$ may express the totality of conceivable species that might
336 be available to the community. In this case $\rho(\epsilon)$ would depend on how we define
337 species in the first place. With reference to Eq. (7), the choice $\alpha = 0$ corresponds to
338 *defining* “species” by their metabolic requirement, *i.e.* discretize ϵ -space into equal
339 intervals of width $\Delta\epsilon$ and define species i to be the set of individuals whose metabolic
340 requirement e lies between $(i - 1)\Delta e$ and $i\Delta e$. Alternatively, “species” could be
341 defined via biomass (in which case the value of α in Eq. (8) may reflect metabolic
342 scaling as in Dewar & Porté 2008); or via other individual traits (t) on which
343 metabolic requirement depends, $\epsilon(t)$.

344 In either case, VEG opens up ways to understanding the link between functional traits and CSDs
345 that are simply not available to METE.

346 One of METE's great strengths is that it only requires three parameters S , N and E for all
347 of its predictions. How might the above insights be used to improve the predictions of METE
348 without damaging this exceptional parsimony? The answer may lie in the inclusion of a prior
349 distribution for ϵ representing a contribution from the density of species $\rho(\epsilon)$, which plays a role
350 in ecology analogous to the “density of states” in physics describing the distribution in energy-
351 space of available quantum-mechanical particle states. The upper panels of Fig. 3 suggest that

352 this trait distribution should give less weight to the higher values of ϵ , and also to the lower
353 values of ϵ in the Barro Colorado and Hawaiian communities.

354 **REFERENCES**

- 355 Banavar, J. R., Maritan, A., & Volkov, I. (2010). Applications of the principle of maximum
356 entropy: from physics to ecology. *Journal of Physics: Condensed Matter*, 22(6), 063101.
- 357 Bertram, J, and Dewar, R. C. (2015). Combining mechanism and drift in community ecology: a
358 novel statistical mechanics approach, *Theoretical ecology*, 8: 419-35.
- 359 Bertram, J, and Dewar, R. C. (2013). Statistical patterns in tropical tree cover explained by the
360 different water demand of individual trees and grasses, *Ecology*, 94: 2138-44.
- 361 Díaz, S., Purvis, A., Cornelissen, J. H., Mace, G. M., Donoghue, M. J., Ewers, R. M., Jordano, P.
362 & Pearse, W. D. (2013). Functional traits, the phylogeny of function, and ecosystem
363 service vulnerability. *Ecology and evolution*, 3(9), 2958-2975.
- 364 Dewar, R. C., and Porté A. (2008). Statistical mechanics unifies different ecological patterns,
365 *Journal of theoretical biology*, 251: 389-403.
- 366 Dewar, R. C. (2009). Maximum entropy production as an inference algorithm that translates
367 physical assumptions into macroscopic predictions: Don't shoot the messenger. *Entropy*,
368 11(4), 931-944.
- 369 Favretti, M. (2017). Remarks on the maximum entropy principle with application to the
370 maximum entropy theory of ecology. *Entropy*, 20(1), 11.
- 371 Gruner, D.S. (2007) Geological age, ecosystem development, and local resource constraints on
372 arthropod community structure in the Hawaiian Islands. *Biological Journal of the*
373 *Linnean Society*, **90**, 551-570.
- 374 Harte, J., (2011) Maximum entropy and ecology: a theory of abundance, distribution, and
375 energetics. Oxford University Press

- 376 Harte, J. & Newman, E.A. (2014) Maximum entropy as a foundation for ecological theory.
377 *Trends in Ecology and Evolution*, **29**, 384-389.
- 378 Harte, J, T Zillio, E Conlisk, and AB Smith. (2008). Maximum entropy and the state-variable
379 approach to macroecology, *Ecology*, 89: 2700-11.
- 380 Harte, J., Smith, A.B. and Storch, D., 2009. Biodiversity scales from plots to biomes with a
381 universal species–area curve. *Ecology letters*, 12(8), pp.789-797.
- 382 Harte, J., Newman, E. A., & Rominger, A. J. (2017). Metabolic partitioning across individuals in
383 ecological communities. *Global Ecology and Biogeography*, 26(9), 993-997.
- 384 He, F. (2010). Maximum entropy, logistic regression, and species abundance. *Oikos*, 119(4),
385 578-582.
- 386 Hubbell, S., Condit, R., & Foster, R. (2005). Forest census plot on Barro Colorado Island.
387 Retrieved from <http://ctfs.arnarb.harvard.edu/webatlas/datasets/bci>
- 388 McGill, B. (2003). Strong and weak tests of macroecological theory. *Oikos*, 102, 679–685.
- 389 McGill, B. J., Enquist, B. J., Weiher, E., & Westoby, M. (2006). Rebuilding community ecology
390 from functional traits. *Trends in ecology & evolution*, 21(4), 178-185.
- 391 McGill, Brian J, Rampal S Etienne, John S Gray, David Alonso, Marti J Anderson, Habtamu
392 Kassa Benecha, Maria Dornelas, Brian J Enquist, Jessica L Green, and Fangliang He.
393 (2007). Species abundance distributions: moving beyond single prediction theories to
394 integration within an ecological framework, *Ecology letters*, 10: 995-1015.
- 395 McGill, B. J., & Nekola, J. C. (2010). Mechanisms in macroecology: AWOL or purloined letter?
396 Towards a pragmatic view of mechanism. *Oikos*, 119(4), 591-603
- 397 Muller-Landau, H.C., Condit, R.S., Harms, K.E., Marks, C.O., Thomas, S.C., Bunyavejchewin,
398 S., Chuyong, G., Co, L., Davies, S., Foster, R. & Gunatilleke, S. (2006). Comparing

- 399 tropical forest tree size distributions with the predictions of metabolic ecology and
400 equilibrium models. *Ecology Letters*, **9**, 589-602.
- 401 Newman, E.A., Harte, M.E., Lowell, N., Wilber, M. & Harte, J. (2014) Empirical tests of within-
402 and across-species energetics in a diverse plant community. *Ecology*, **95**, 2815-2825.
- 403 Newman, E.A., Wilber, M.Q., Kopper, K.E., Moritz, M.A., Falk, D.A., McKenzie, D. and Harte,
404 J., 2018. Disturbance macroecology: integrating disturbance ecology and macroecology
405 in different-age post-fire stands of a closed-cone pine forest as a case study. bioRxiv,
406 p.309419.
- 407 O'Dwyer, J.P., Rominger, A. & Xiao, X. (2017). Who Constrains the Constraints? Reinterpreting
408 maximum entropy in ecology: a null hypothesis constrained by ecological
409 mechanism. *Ecology letters*, **20**, 832-841.
- 410 Pueyo, S., He, F., & Zillio, T. (2007). The maximum entropy formalism and the idiosyncratic
411 theory of biodiversity. *Ecology Letters*, 10(11), 1017-1028.
- 412 Rominger, A.J. and Merow, C. (2016) meteR: an R package for testing the maximum entropy
413 theory of ecology. *Methods in Ecology and Evolution*, **8**, 241-247.
- 414 Shipley, B., Vile, D., & Garnier, É. (2006). From plant traits to plant communities: a statistical
415 mechanistic approach to biodiversity. *Science*, 314(5800), 812-814.
- 416 Shipley, B. (2010). Community assembly, natural selection and maximum entropy models.
417 *Oikos*, 119(4), 604-609.
- 418 Supp, S.R., Xiao, X., Ernest, S.K.M. and White, E.P., 2012. An experimental test of the response
419 of macroecological patterns to altered species interactions. *Ecology*, 93(12), pp.2505-
420 2511.

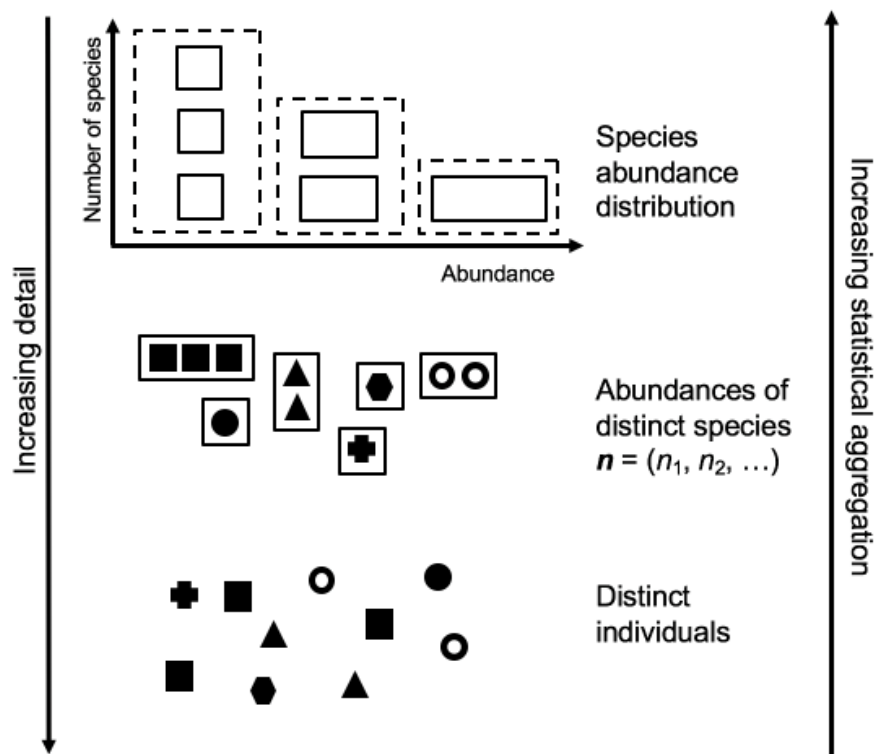
- 421 Supp, S.R. and Ernest, S.M., 2014. Species-level and community-level responses to disturbance:
422 a cross-community analysis. *Ecology*, 95(7), pp.1717-1723.
- 423 Tokeshi, M. (1993). Species abundance patterns and community structure. *Advances in*
424 *ecological research* **24**, 111-186.
- 425 Ulrich, W., Ollik, M., & Ugland, K. I. (2010). A meta-analysis of species–abundance
426 distributions. *Oikos*, 119(7), 1149-1155.
- 427 Vellend, M. (2010). Conceptual synthesis in community ecology. *The Quarterly review of*
428 *biology*, 85(2), 183-206
- 429 Xiao, X., McGlenn, D. & White, E. (2015) A strong test of the Maximum Entropy Theory of
430 Ecology. *American Naturalist*, **185**, E70-E80.
- 431
- 432

433 **FIGURES**

434

435 **Figure 1.** Three levels of detail commonly used for describing ecological communities. At the
436 greatest level of detail (bottom), the distinct identities of individuals and their spatial locations
437 are known. At the intermediate levels of detail found in many well-mixed models such as Lotka-
438 Volterra models (middle), the abundance of each distinct species is known. At the lowest level of
439 detail and highest level of statistical aggregation (top), species identities are lost, and the SAD
440 provides the only description of species diversity.

441

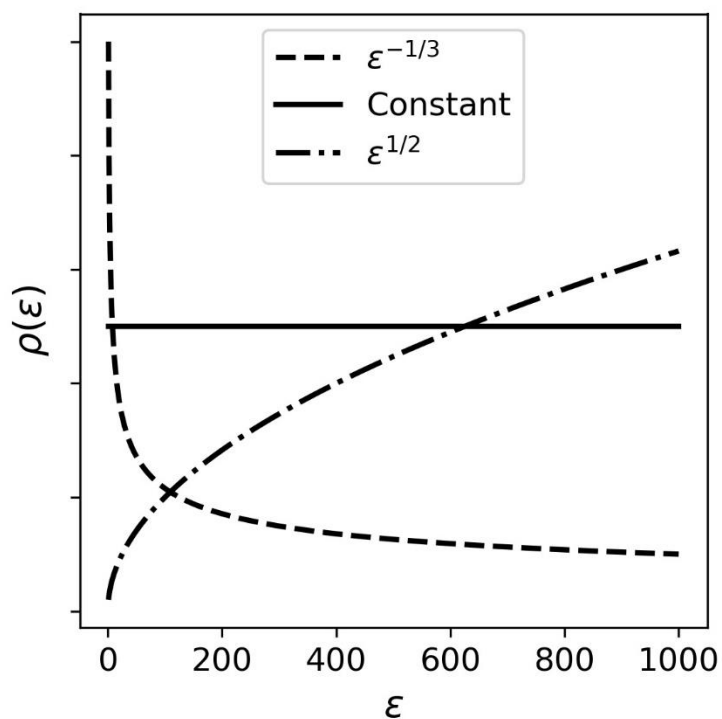


442

443

444

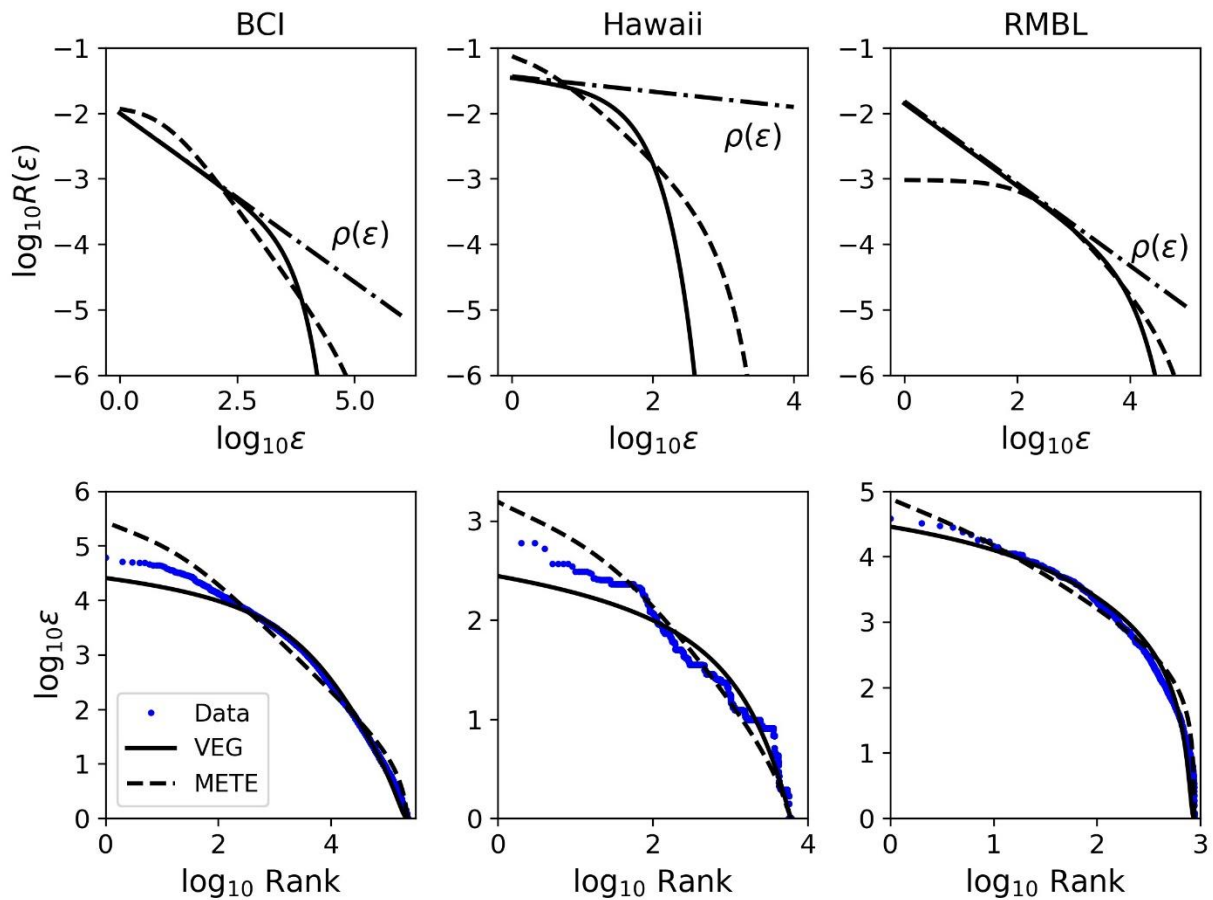
445 **Figure 2.** The function $\rho(\epsilon)$ counts the local density of metabolic rates in the assumed spectrum
446 of possible rates $\epsilon_1 \leq \epsilon_2 \leq \dots$ in VEG. It represents the relatively probability that a randomly
447 selected VEG “species” has metabolic rate ϵ when sampled from all possible “species”.
448



449

450

451 **Figure 3.** Comparison of METE and VEG rank-metabolism relationships in multiple
452 communities: Barro Colorado Island trees (BCI; $\alpha = -0.52, \rho(1) = 10^4$), Hawaiian island
453 arthropods (Hawaii; $\alpha = -0.12, \rho(1) = 10^2$), and Rocky Mountain subalpine meadow plants
454 (RMBL; $\alpha = -0.63, \rho(1) = 10^2$) communities. Each upper/lower panel pair shows the $R(\epsilon)$
455 and rank- ϵ curves for the same METE and VEG ecosystem structure functions.



456

457

458

459 **APPENDIX A: NOTE ON THE DEFINITION OF THE ECOSYSTEM STRUCTURE**

460 **FUNCTION**

461 Our definition of the ecosystem structure function differs slightly from that given in Harte et al.
462 (2008), which reads “[The ecosystem structure function] is the probability that if a species is
463 picked at random [...], then it has abundance n and if an individual is picked at random *from that*
464 *species*, then its metabolic requirement is in the range $\epsilon, \epsilon + d\epsilon$ ” (our italics). The Harte et al.
465 (2008) definition suggests that METE keeps track of species identity. In fact, $R_M(n, \epsilon)$ depends
466 only on n and ϵ , and not species identity. Thus, $R_M(\epsilon|n) = R_M(n, \epsilon) / \sum_n R_M(n, \epsilon)$ is the
467 probability of picking an individual with metabolic requirement ϵ conditional on it coming from
468 a species with abundance n . There is no way within METE to distinguish between different
469 species with the same abundance n , and therefore there is no reason to specify which species the
470 individual is selected from in the definition of the ecosystem structure function (Favretti, 2017).

471

472 **APPENDIX B: DERIVING THE ECOSYSTEM STRUCTURE FUNCTION FROM**
 473 **DISTINGUISHABLE SPECIES**

474 In section 2.1, the METE structure function R_M was inferred directly from community-level
 475 constraints. Here we derive an analogous VEG structure function $R_V(n, \epsilon)$.

476 We start by defining the probability distribution $P(n, \epsilon, i, \mathbf{n})$ as follows: $P(n, \epsilon, i, \mathbf{n})d\epsilon$ is
 477 the joint probability that the community has species abundances \mathbf{n} , that a species picked at
 478 random from the community has species label i , that this chosen species has abundance $n_i = n$,
 479 and that an individual from this chosen species has metabolic requirement in the interval $(\epsilon, \epsilon +$
 480 $d\epsilon)$. $R_V(n, \epsilon)$ is then obtained by marginalizing with respect to i and \mathbf{n} , i.e. $R_V(n, \epsilon) =$
 481 $\sum_i \sum_{\mathbf{n}} P(n, \epsilon, i, \mathbf{n})$ where $n \geq 1$.

482 To marginalize P , we first write it as a product of conditional distributions

$$483 \quad P(n, \epsilon, i, \mathbf{n}) = P(\epsilon|n, i, \mathbf{n})P(n|i, \mathbf{n})P(i|\mathbf{n})P(\mathbf{n}). \quad (\text{B1})$$

484 Here $P(\mathbf{n}) = p(\mathbf{n})$ is the probability that the community has abundance vector \mathbf{n} , $P(i|\mathbf{n}) =$
 485 $(1 - \delta_{n_i}^0)/S(\mathbf{n})$ is the probability that a species picked from a community with abundances \mathbf{n}
 486 has species label i (i.e. 0 if species i is absent, $1/S(\mathbf{n})$ if present), $P(n|i, \mathbf{n}) = \delta_{n_i}^n$ is the
 487 probability that species i has abundance n given the species abundances are \mathbf{n} , and $P(\epsilon|n, i, \mathbf{n})d\epsilon$
 488 is the probability that an individual picked from species i has metabolic requirement in the
 489 interval $(\epsilon, \epsilon + d\epsilon)$ given species i has abundance n and the community abundances are \mathbf{n} . We
 490 thus obtain

$$492 \quad R_V(n, \epsilon) = \sum_i \sum_{\mathbf{n}} P(\epsilon|n, i, \mathbf{n})\delta_{n_i}^n(1 - \delta_{n_i}^0)p(\mathbf{n})/S(\mathbf{n})$$

$$493 \quad = \sum_i \sum_{\mathbf{n}} P(\epsilon|n, i, \mathbf{n})\delta_{n_i}^n p(\mathbf{n})/S(\mathbf{n}) \quad (\text{B2})$$

491 where we have used the fact that $\delta_{n_i}^n(1 - \delta_{n_i}^0) = \delta_{n_i}^n$ for $n \geq 1$.

494 In the case of VEG, all individuals in species i have the same metabolic requirement ϵ_i ,
495 and so $P(\epsilon|n, i, \mathbf{n}) = \delta(\epsilon - \epsilon_i)$ where δ is the Dirac delta function (i.e. the probability that a
496 randomly selected individual from species i has metabolic requirement ϵ is 1 in the immediate
497 vicinity of ϵ_i , and is 0 otherwise). Thus, from (B2) we have

$$503 \quad R_V(n, \epsilon) = \sum_i \delta(\epsilon - \epsilon_i) \sum_{\mathbf{n}} \frac{\delta_{n_i}^n p(\mathbf{n})}{S(\mathbf{n})} \quad (\text{B3})$$

498 An ecologically important special case of (B3) occurs when the variation in the number of
499 species present from one snapshot to the next is small relative to the expected number of species,
500 such that $S(\mathbf{n})$ is approximately constant with value given by $S = \sum_{\mathbf{n}} S(\mathbf{n}) p(\mathbf{n})$. This occurs,
501 for example, if most of the species present have large expected abundances. Eq. (B3) then
502 simplifies to

$$504 \quad R_V(n, \epsilon) = \frac{1}{S} \sum_{\{i|n_i \geq 0\}} \delta(\epsilon - \epsilon_i) P(n_i = n) \quad (\text{B4})$$

505 where $P(n_i = n) = \sum_{\mathbf{n}} \delta_{n_i}^n p(\mathbf{n})$ is the probability that species i has abundance n .

506 In VEG, we have from Eq. (6) (Bertram and Dewar, 2015)

$$507 \quad P(n_i = n) = (1 - e^{-(\mu_1 + \mu_2 \epsilon_i)}) e^{-(\mu_1 + \mu_2 \epsilon_i)n}.$$

508 To make it explicit that this probability depends on the metabolic requirement of species i , we
509 use the notation $p(n|\epsilon_i) \equiv P(n_i = n)$ (that is, $P(n_i = n)$ in VEG is the probability that a species
510 has abundance n given that its metabolic requirement is ϵ_i).

511 Since the ecosystem structure function is a probability density in the continuous variable
512 ϵ , we can introduce a spectral density $\rho(\epsilon)d\epsilon$ that counts the number of metabolic requirement
513 levels ϵ_i in each interval $(\epsilon, \epsilon + d\epsilon)$. From Eq. (B4), $R_V(n, \epsilon)$ can then be written in the form

$$514 \quad R_V(n, \epsilon) = \frac{\rho(\epsilon)p(n|\epsilon)}{S}$$

515 **APPENDIX C: THE VEG SPECIES ABUNDANCE DISTRIBUTION**

516 Assuming a power-law density of states $\rho(\epsilon) \propto \epsilon^\alpha$, we have from Eq. (7)

517
$$R(n) = \int_0^\infty R(n, \epsilon) d\epsilon \propto \int_0^\infty \epsilon^\alpha (1 - e^{-(\mu_1 + \mu_2 \epsilon)}) e^{-(\mu_1 + \mu_2 \epsilon)n} d\epsilon$$

518 By making the substitution $x = \mu_2 \epsilon n$, the integral is found to be:

519
$$\int_0^\infty \epsilon^\alpha (1 - e^{-(\mu_1 + \mu_2 \epsilon)}) e^{-(\mu_1 + \mu_2 \epsilon)n} d\epsilon = \Gamma(\alpha + 1) \frac{e^{-\mu_1 n}}{(\mu_2 n)^{\alpha+1}} \left[1 - e^{-\mu_1} \left(\frac{1}{1 + 1/n} \right)^{\alpha+1} \right]$$

520 where

521
$$\Gamma(\alpha + 1) = \int_0^\infty x^\alpha e^{-x} dx$$

522 is the gamma function.

523 For large n we have $1/(1 + 1/n) \approx 1$ so that:

524
$$R(n) \propto \frac{e^{-\mu_1 n}}{(\mu_2 n)^{\alpha+1}}$$

525

526