1   **Title: General models of ecological diversification. II. Simulations and empirical**

2   **applications**

3

4   Author: Philip M. Novack-Gottshall

5

7

8   RRH: SIMULATING ECOLOGICAL DIVERSIFICATION MODELS

9   LRH: PHILIP M. NOVACK-GOTTSHALL

10

14

15   *Abstract.*—Models of functional ecospace diversification within life-habit frameworks

16   (functional-trait spaces) are increasingly used across community ecology, functional ecology,

17   and paleoecology. In general, these models can be represented by four basic processes, three that

18   have driven causes and one that occurs through a passive process. The driven models include

19   redundancy (caused by forms of functional canalization), partitioning (specialization), and

20   expansion (divergent novelty), but they also share important dynamical similarities with the

21   passive neutral model. In this second of two companion articles, Monte Carlo simulations of

22   these models are used to illustrate their basic statistical dynamics across a range of data

23   structures and implementations. Ecospace frameworks with greater numbers of characters

24     (functional traits) and ordered (multi-state) character types provide more distinct dynamics and

25     greater ability to distinguish the models, but the general dynamics tend to be congruent across all

26     implementations. Classification tree methods are proposed as a powerful means to select among

27     multiple candidate models when using multivariate data sets. Well-preserved Late Ordovician

28     (type Cincinnatian) samples from the Kope and Waynesville Formations are used to illustrate

29     how these models can be inferred in empirical applications. Initial simulations overestimate the

30     ecological disparity of actual assemblages, confirming that actual life habits are highly

31     constrained. Modifications incorporating more realistic assumptions (such as weighting potential

32     life habits according to actual frequencies and adding a parameter controlling the strength of

33     each model's rules) provide better correspondence to actual assemblages. Samples from both

34     formations are best fit by partitioning (and to lesser extent redundancy) models, consistent with a

35     role for local processes. When aggregated as an entire formation, the Kope Formation pool

36     remains best fit by the partitioning model, whereas the entire Waynesville pool is better fit by the

37     redundancy model, implying greater beta diversity within this unit. The 'ecospace' package is

38     provided to implement the simulations and to calculate their dynamics using the R statistical

39     language.

40

41     *Philip M. Novack-Gottshall. Department of Biological Sciences, Benedictine University, Lisle, IL*

42     *60532 . E-mail: pnovack-gottshall@ben.edu*

43

44

**Introduction**

45      Understanding the causes of ecological diversification remains an important goal in

46  paleontology, ecology, and evolutionary biology. Andrew Bush and I recently introduced (2012)

47  several models of diversification in ecospace (functional trait-space) that are useful for

48  conceptually representing a broad range of diversifications, whether at the scale of ecological

49  assembly of communities or whether shaping entire biotas over evolutionary timescales. The

50  models share many similarities with hypotheses used by community and functional ecologists

51  (Villéger et al. 2010, 2011, Mouillot et al. 2013, Vogt et al. 2013, Gerisch 2014), as well as those

52  used by paleontologists studying morphological, and increasingly ecological (functional),

53  disparity (Dineen et al. 2014, Miller et al. 2014, Mitchell and Makovicky 2014, Dick and

54  Maxwell 2015, Dineen et al. 2015, Knope et al. 2015). Most of these ideas invoke the processes

55  of ecological canalization, specialization, and/or divergence during diversification. In a

56  companion article (Novack-Gottshall 2016), I described four of these models in more detail and

57  discussed their usefulness as general models of ecological diversification. The redundancy model

58  occurs when successive species occupy similar regions in ecospace as previously existing

59  species, and can be considered a form of ecological canalization. The partitioning model occurs

60  when successive species progressively subdivide ecospace, such as occurs in niche partitioning

61  and specialization. The expansion model occurs when successive species progressively explore

62  novel portions of ecospace, such as occurs during niche divergence. Despite their mechanistic

63  differences, these three models are all driven (*sensu* McShea 1994) by particular causes. In

64  contrast, the fourth neutral model occurs as a passive model, in which ecospace is inhabited at

65  random.

Note: The line numbers in the margin are 45–66. Let me recount.

**Introduction**

45      Understanding the causes of ecological diversification remains an important goal in

46  paleontology, ecology, and evolutionary biology. Andrew Bush and I recently introduced (2012)

47  several models of diversification in ecospace (functional trait-space) that are useful for

48  conceptually representing a broad range of diversifications, whether at the scale of ecological

49  assembly of communities or whether shaping entire biotas over evolutionary timescales. The

50  models share many similarities with hypotheses used by community and functional ecologists

51  (Villéger et al. 2010, 2011, Mouillot et al. 2013, Vogt et al. 2013, Gerisch 2014), as well as those

52  used by paleontologists studying morphological, and increasingly ecological (functional),

53  disparity (Dineen et al. 2014, Miller et al. 2014, Mitchell and Makovicky 2014, Dick and

54  Maxwell 2015, Dineen et al. 2015, Knope et al. 2015). Most of these ideas invoke the processes

55  of ecological canalization, specialization, and/or divergence during diversification. In a

56  companion article (Novack-Gottshall 2016), I described four of these models in more detail and

57  discussed their usefulness as general models of ecological diversification. The redundancy model

58  occurs when successive species occupy similar regions in ecospace as previously existing

59  species, and can be considered a form of ecological canalization. The partitioning model occurs

60  when successive species progressively subdivide ecospace, such as occurs in niche partitioning

61  and specialization. The expansion model occurs when successive species progressively explore

62  novel portions of ecospace, such as occurs during niche divergence. Despite their mechanistic

63  differences, these three models are all driven (*sensu* McShea 1994) by particular causes. In

64  contrast, the fourth neutral model occurs as a passive model, in which ecospace is inhabited at

65  random.

67    The ecospace that becomes structured during each model's implementation can be

68    envisioned as a landscape (or multivariate ordination) defined by the life habits (functional trait

69    combinations) of constituent organisms (Fig. 1). Because each model ecospace is associated with

70    distinct distributions of life habits (ecological structure) that can be quantified using various

71    disparity statistics, the models also suggest promise as the basis for quantitative tests if

72    implemented through computer simulations. In the companion article, I summarized the expected

73    dynamics for these models, using a range of metrics from the morphological disparity and

74    functional diversity literature. Sensitivity analyses (Ciampaglio et al. 2001, Villier and Eble

75    2009, Mouchet et al. 2010) have demonstrated inherent strengths and weaknesses to each metric,

76    and it remains to be demonstrated how sensitive these dynamics are to different implementations

77    of the models and to different data structures (ecospace frameworks). More important, it remains

78    unclear how to select among alternative models when using multiple, interdependent statistics

79    (Burnham and Anderson 2002, Johnson and Omland 2004, Grueber et al. 2011).

80    In this second of two companion articles, I demonstrate how these models can be

81    implemented as stochastic simulations. Sensitivity analyses are conducted to evaluate how

82    differences in number of characters (functional traits), character types, and various

83    implementations of the simulations affect resulting statistical dynamics. Because the life habits

84    theoretically allowed by any ecospace framework will always be greater than those that exist in

85    reality (either by logical, biological or other constraints), the simulation results suggest several

86    ways they can be made more realistic, i.e., so that they better approximate actual assemblages.

87    The use of classification trees (Breiman et al. 1993) trained on multivariate Monte Carlo data

88    sets (simulations) is proposed as a novel method for selecting among multiple models. As a case

89    study, these methods are applied to well-preserved fossil assemblages from the type Cincinnatian

90    (Upper Ordovician of Ohio, Indiana, and Kentucky). Because the samples hierarchically span

91    multiple localities and stratigraphic levels, they offer useful tests of the generality of the models.

92    In addition to its application for fossil samples, the methods proposed here have value to those

93    studying functional diversity in modern communities, where debate continues regarding the best

94    manner to test important hypotheses on ecological structure.

95

96                    **Implementing the Models using Monte Carlo Simulations**

97        Statistical dynamics for the four models can be obtained using Monte Carlo simulations,

98    which are provided here for a range of ecospace frameworks. The simulations were programmed

99    using R (R Development Core Team 2015), and a package—'ecospace: Simulating Community

100    Assembly and Ecological Diversification Using Ecospace Frameworks' (Novack-Gottshall

101    2015)—is provided to allow others to implement them for their own purposes. In the discussion

102    below, single quotation marks (' ') refer to control parameters or variables called within the

103    package functions.

104

105    Theoretical Ecospace Framework

106        Prior to each simulation, a theoretical ecospace framework is defined (with the function

107    'create_ecospace'), which specifies the realm of possible life habits that all species can inhabit

108    during the simulation. This step is flexible, allowing any number of characters and character

109    types used by functional ecologists (e.g., Mouchet et al. 2010, Villéger et al. 2011),

110    paleoecologists (e.g., Bambach et al. 2007, Bush et al. 2007, Novack-Gottshall 2007), and those

111    studying morphological disparity (e.g., Foote 1994, Ciampaglio et al. 2001). Examples include

112    factors (e.g., diet can be autotrophic, carnivorous, herbivorous, or microbivorous), ordered

113  factors (factors with a specified order, such as for mobility: habitual > intermittent > facultative

114  > passive > sedentary), ordered numerics (discrete or continuous numeric values, such as body

115  size), and binary/numeric character types. Binary states can be treated individually

116  (present=1/absent=0) or can be treated as multiple states within a single character. For example,

117  the character "reproduction" can be treated as including two states [sexual, asexual] with

118  exclusively sexual habits coded as [0,1], exclusively asexual as [1,0], and hermaphrodites as

119  [1,1]. This coding strategy is preferable over the use of discrete factors (sexual, asexual,

120  hermaphroditic) in allowing flexibility for possible (and biologically frequently occurring) state

121  combinations, while maintaining logical functional distances (in a Euclidean or other metric

122  sense), such that hermaphrodites would be closer to sexual and asexual habits that either

123  exclusive mode is to the other. (This example is trivial as it could easily be coded as an ordered

124  factor, but such binary coding schemes can be extended to additional character states where such

125  factorial options are no longer practical.) See Novack-Gottshall (2007) and Supplementary

126  Appendix 1 for additional discussion and examples of the value of multiple binary states within a

127  single character (functional category). When such multiple binary characters are included,

128  subsequent simulations (and the ecospace framework that controls them) specify that "all-

129  absences" are impossible (i.e., an organism cannot be neither sexual nor asexual, [0,0]). This

130  behavior can be controlled with a 'constraint' parameter, which can allow any combination (e.g.,

131  [1,1,1]) or specify a maximum number of "multi-presences" (e.g., setting 'constraint=2' allows

132  [1,0,0], [0,1,0], [0,0,1], [1,1,0], [1,0,1], and [0,1,1] as state combinations, but excludes [1,1,1];

133  setting 'constraint=1' only allows the first three of these combinations).

134      Another feature when defining the ecospace framework is the ability to weight character

135  states, such that certain states occur more often than others. This constraint can be done using

136   customized weights, or calculated automatically using empirical samples (such as a regional

137   species pool) to assign weights based on their frequency of occurrence. A weight of zero

138   excludes that character state from inclusion in the ecospace framework, and thus from

139   simulations based on it. Although not explored here, these weights could be used, with some

140   care, in linking covariation among states across different characters, and offer potential to extend

141   these simulations to studies of morphological disparity, where non-applicable character states—

142   e.g., glabella shape for a mollusk—are frequently encountered.

143

144   Simulations and Implementation as Algorithms

145       All simulations are implemented as Monte Carlo processes in which species are added

146   iteratively to assemblages, with all added species having their character states specified by the

147   model rules below (Fig. 1). Aside from these model rules (and whatever inherent constraints are

148   imposed by the ecospace framework), there are no deterministic components to the stochastic

149   simulations. Simulations begin with the seeding of a specified number of species, chosen at

150   random (with replacement) from either the species pool (if provided) or following the neutral-

151   rule algorithm (if not). Once seeded, the simulations proceed iteratively (character-by-character,

152   species-by-species) by following the appropriate algorithm until terminated at a pre-specified

153   species-richness value (sample size). A 'strength' parameter can be used to specify the probability

154   that a simulation follows each rule, with possible values ranging from 1.0, in which the model

155   rule is always followed, to 0.0, in which the model rule is never followed (and during which the

156   simulation switches to the default neutral model rule.) Although not explored here, it might be

157   profitable in future analyses to modify the model-selection analyses into an optimization routine

158 so that the 'strength' parameter is treated as a free parameter to be estimated by the data, rather

159 than chosen as an arbitrary value.

160 *Neutral rule*.—Iteratively add species independently, assigning characters character-by-

161 character as random multinomial draws (with probabilities assigned by character-state weights, if

162 assigned) from the theoretical ecospace framework. If weighting was assigned to the ecospace

163 framework, character states are assigned in proportion to those weights, on average.

164 This "rule" is not technically a rule, but the absence of other rules: all further simulations

165 default to this neutral algorithm when, for example, a simulation is implemented with a strength

166 parameter of zero, and converge on this model when implemented at strengths less than one. It is

167 important to note that the neutral rule is not a simple permutation (randomization) test

168 (Kowalewski and Novack-Gottshall 2010), drawing at random from the species pool (except for

169 the seed species, if a species pool is provided). Most tests in the functional ecology literature use

170 such simple permutation tests, and such a test can be enacted in the neutral model by setting the

171 number of seed species to a sufficiently large value, such as the size of the species pool. The life

172 habit of each new species is built character-by-character from the realm of theoretically possible

173 states allowed by the ecospace framework. Thus, new species can occupy combinations of

174 character states that did not occur in the species pool (if provided). This is an important feature

175 of the simulations, allowing the entire theoretical ecospace to be explored by the neutral model if

176 given sufficient numbers of species. If it turns out that the actual species pool is structurally

177 different from a similarly parameterized neutral model, it can be concluded that there is an

178 important mechanism controlling how the actual sample was composed.

179 *Redundancy rule*.—Pick one existing species at random and create a new species using that

180 species' characters as a template. A character is modified (using a random multinomial draw

181  from the theoretical ecospace, as in the neutral rule) according to the strength parameter. When

182  the strength parameter is set to unity, all species will be functionally identical to the original seed

183  species; when set to zero, the simulation is identical to the neutral rule. Because new character

184  states can be any allowed by the ecospace framework, there is the possibility of obtaining

185  redundancy greater than that specified by strength parameters less than 1.0 (if, for example, the

186  new randomly chosen character states are identical to those of the template species).

187      *Partitioning rule*.—Calculate distances between all pairs of species, using Euclidean distance

188  if all characters are numeric/binary or ordered numeric types or using Gower distance if any

189  character type is an ordered or unordered factor. There are two algorithms to choose from for the

190  next step, a 'strict' (and default) partitioning implementation and a 'relaxed' implementation. In

191  the strict (or "minimum distant-neighbor") implementation, identify the maximum distances

192  between all pairs of species (the "most-distant neighbors"); the new species will have character

193  states based on the pair of species that is the minimum of these distances. This implementation

194  progressively fills in the largest parts of the ecospace that are least occupied between

195  neighboring species, and on average partitions the ecospace in straight-line gradients between

196  seed species. In the relaxed (or "maximum nearest-neighbor") implementation, identify nearest-

197  neighbor distances between all pairs of species; the new species will have character states based

198  on the pair that is the maximum of these distances. This implementation places new species in

199  the most unoccupied portion of the ecospace that is within the cluster of pre-existing species,

200  often the centroid. (See Fig. 1 for visual portrayal of these alternate implementations.) In both

201  cases, the life habit for each new species is built as a resampled combination of the character

202  states of the chosen neighbor pair, with probabilities controlled by the strength parameter.

203  Ordered, multistate character partitioning (whether factor or numeric) can include any state equal

204    to or between the observed states of existing species. Each newly assigned character is compared

205    with the ecospace framework to confirm that it is an allowed state combination before

206    proceeding to the next character; if the newly built character is disallowed from the ecospace

207    framework (i.e., because it has "dual absences" [0,0], has been excluded based on the species

208    pool, or is not allowed by the ecospace 'constraint' parameter), then the character-selection

209    algorithm is repeated until an allowable character is selected. When simulations proceed to very

210    large sample sizes (>100), this confirmatory process can require long computational times, and

211    produce "new" intermediate species that are functionally identical to pre-existing species. This

212    can occur, for example, when no life habits, or perhaps only one, exist that forms an allowable

213    intermediate between the selected neighbors. This behavior mimics the "pathologically" tight

214    packing of species that has been demonstrated in large sample-size simulations of niche

215    partitioning (Kinzig et al. 1999). This behavior also causes dynamics at such sample sizes to

216    share important similarities with weakened implementations of the redundancy model (see

217    additional discussion below).

218         *Expansion rule*.—Calculate distances between all pairs of species, as done in the partitioning

219    algorithm, and identify the species pair that is maximally distant. The newly added species has

220    character states chosen at random that are equal to or more extreme (if ordered or numeric

221    character types, or that are the same or different for unordered factors) than those in this species

222    pair, with the strength parameter controlling the probability of following this rule. As with the

223    partitioning rule, each newly assigned character is confirmed to be an allowed character-state

224    combination before proceeding with remaining characters and species. Like the partitioning rule,

225    the algorithm can take long computational times to run to completion for large sample sizes, and

226   shares the similar property that functionally identical life habits may occur by virtue of saturation

227   of ecospace-allowable life habits.

228   *Other simulation assumptions and limitations*.—The simulations explicitly assume that

229   dispersal is guaranteed to all species, provided that new species have appropriate character states

230   as proscribed by the ecospace framework and the model rules. This is an important distinction

231   from Hubbell's neutral model (2001) and many other spatially structured models in community

232   ecology, but consistent with how many species pools are defined in ecological studies (Cornell

233   and Harrison 2014). Extinction, speciation, nor emigration is allowed during the course of a

234   simulation (although they can play important roles in the definition of the original species pool).

235   All species that immigrate to the assemblage remain there, with their character states retained,

236   throughout the course of the simulation. Such ecological and evolutionary processes (character

237   displacement, competitive exclusion, habitat filtering, etc.) are only present as implemented by

238   the model rules governing addition of new species to the assemblage.

239   The framework and simulation algorithms currently do not incorporate relative

240   abundance (cf., Bush et al. 2007, Deline 2009, Deline et al. 2012), life-history characteristics,

241   population demographics, dispersal and local extinction, spatial structure, speciation and

242   evolution, or other constraints on the regional pool. The incorporation of relative abundance

243   could be especially worthwhile for ecological studies, as it would allow the simulations to better

244   mimic the assembly of ecological communities. For example, abundance plays an important role

245   in immigration from regional species pools (Hubbell 2001, Patzkowsky and Holland 2007,

246   Cornell and Harrison 2014), and abundant taxa can have larger influences on the composition of

247   functional traits within communities and they can have a greater impact on the functional identity

248   of taxa that settle later (Fargione et al. 2003, Guillemot et al. 2011). Functional diversity studies

249    increasingly incorporate relative abundance ((Villéger et al. 2008, Laliberté and Legendre 2010,

250    Fontana et al. 2015), but the models used here are limited to presence/absence data, which some

251    functional ecology studies (e.g., Ackerly and Cornwell 2007) have demonstrated produce similar

252    results as those incorporating abundance data.

253        The simulations also currently do not incorporate phylogenetic structure in an explicit

254    sense. The three driven models are Markovian in that pre-existing taxa affect which functional

255    traits can enter at later stages of the simulation, but there is no explicit assumption of heritability

256    or phylogenetic relatedness among simulated species. In contrast, the neutral model is non-

257    Markovian in that all taxa are added independently of one another, without regard to their

258    functional traits. The lack of such phylogenetic structure perhaps hinders the application of these

259    models to diversification within individual, evolving lineages, but the models are intended as

260    reasonable approximations of large-scale evolutionary diversifications involving many disparate

261    lineages, and where the focus is on the ecospace structure of the overall biota. The existing

262    simulations could be modified to include such features, and would prove worthwhile for future

263    studies.

264

265    Calculating Functional Diversity and Ecological Disparity Statistics

266        As each simulation proceeds, the theoretical ecospace (functional-trait space) becomes

267    progressively filled with new species having combinations of character states (life habits)

268    selected by the model rules (Fig. 1). Because the mechanisms controlling community assembly

269    vary among models, it is expected that the statistical dynamics will likewise vary. These model-

270    specific dynamics can be calculated as a function of species richness. Here, eight widely used

271  statistics are calculated, drawn from the morphological disparity and functional-diversity

272  literature (Ciampaglio et al. 2001, Wills 2001, Villéger et al. 2008, Mouchet et al. 2010).

273      The four morphological (here ecological) disparity measures include the following (Foote

274  1993, Ciampaglio et al. 2001, Wills 2001). Unique trait-combination (life-habit) richness (H)

275  measures the number of unique life habits within each sample. Mean distance (D) and maximum

276  range (M) measure the mean and maximal distance, respectively, in functional-trait space among

277  species, using Euclidean distance if all characters are numeric/binary or ordered numeric types,

278  or using Gower distance (1971) if any character type is an ordered or unordered factor. Total

279  variance (V) measures the sum of variances for each character state across species (Van Valen

280  1974); when using factor character types, this statistic cannot be measured and is excluded.

281      The four functional diversity measures include the following (Mason et al. 2005,

282  Anderson et al. 2006, Villéger et al. 2008, Laliberté and Legendre 2010, Mouchet et al. 2010,

283  Mouillot et al. 2013). Functional richness (FRic) measures the minimal convex-hull volume in

284  multidimensional principal coordinates analysis (PCoA) trait-space ordination. Functional

285  evenness (FEve) measures the evenness of minimum-spanning-tree lengths between species in

286  PCoA trait-space. Functional divergence (FDiv) measures the mean distance of species from the

287  PCoA trait-space centroid. Functional dispersion (FDis) measures the total deviance of taxa from

288  the circle with radius equal to mean distance from PCoA trait-space centroid. FRic and FDiv use

289  a subset (here set to 3 by default) of the PCoA axes for their calculation (because their convex-

290  hull-volume calculation requires more species than functional traits), whereas the other two use

291  the entire PCoA space.

292      These statistics are calculated with the function 'calc_metrics', wrapping around function

293  'dbFD' in the 'FD' package (Laliberté and Shipley 2014) to calculate functional-diversity

294    statistics. The four functional diversity statistics are not calculated for samples having less than

295    four unique life habits, for the same reason as just explained. The calculation of FRic can be

296    computationally demanding, and has proven truculent to many users when dealing with

297    idiosyncratic data structures (see help file in Laliberté and Shipley 2014). Some of these issues

298    are resolved in the implementation here. Specifically, the package defaults to three PCoA axes

299    ('m=3') as a generally useful and computationally tractable number of axes, defaults to Lingoes

300    correction (Legendre and Anderson 1999) when non-Eucliden PCoA axes are calculated, has

301    corrected (starting with 'FD' v. 1.0–12) a previously unnoticed rounding error in the calculation

302    of PCoA scores which caused computation problems with highly redundant datasets, and has

303    been optimized to prevent overwriting errors (of temporarily stored vertices files) when used in a

304    parallel-processing computer environment. Although including more PCoA axes allows greater

305    statistical power (cf., Villéger et al. 2011, Maire et al. 2015), the use of three here provides

306    satisfactory resolution for the often functionally redundant data sets in these simulations. This

307    choice is also warranted because such statistics are calculated directly from multivariate PCoA

308    trait-spaces (instead of dendograms) and because the ecospace frameworks used here have few

309    factor character types (which are especially sensitive to few axes). Three PCoA axes also provide

310    the largest computationally tractable number to allow calculation of these statistics across all

311    sample sizes and models considered.

312        Sensitivity analyses of these statistics (Ciampaglio et al. 2001, Mouchet et al. 2010,

313    Novack-Gottshall 2010, Maire et al. 2015) have concluded that, despite much correlation among

314    them, each statistic has its strengths for particular uses and contexts. H, M, V, FRic, and FDis are

315    all useful measures of disparity (or dispersion of species within functional-trait space). D and

316    FDiv are useful for characterizing internal structure (i.e., clumping or inhomogeneities within the

317    trait-space). FEve is useful for characterizing the extent of spacing among species within the

318    trait-space. Such sensitivity studies consistently recommend that all are useful for particular

319    purposes, with selection left to the researcher, a solution that is decidedly indecisive. Below, I

320    demonstrate that classification trees (trained on all statistics from simulated datasets) offer a

321    productive solution to this subjectivity.

322

323    Model Selection

324        A difficulty in implementing model selection with these simulations is that every

325    simulation iteration (i.e., each addition of a species) yields eight statistics. This lack of statistical

326    independence is typically addressed through model-fitting procedures (such as likelihood-ratio

327    tests or the Akaike information criterion [AIC], Anderson et al. 2000, Burnham and Anderson

328    2002, Johnson and Omland 2004, Grueber et al. 2011) that incorporate dependency relationships

329    explicitly (such as for a multivariate normal model). However, these methods become

330    prohibitively complicated when the statistics involved are not approximately normally

331    distributed. In this case, H is a discrete binomial distribution. Depending on the model, several

332    statistics are highly skewed, with D best fit by gamma distributions and M and V best fit by

333    Weibull distributions (confirmed using maximum likelihood goodness-of-fit tests); in contrast,

334    the four functional diversity statistics are all reasonably well fit by normal distributions. Most

335    functional ecology studies thus use one (or a few) of these statistics for a particular test,

336    subjecting them to a simple permutation test to reject a particular null hypothesis. Because all

337    functional diversity /disparity metrics contribute potentially useful information on functional

338    (ecospace) structure (Ciampaglio et al. 2001, Villier and Eble 2009, Mouchet et al. 2010), there

339    is value in retaining all suitable metrics when possible. The obstacle is doing so while

340 simultaneously considering multiple candidate models (Johnson and Omland 2004, Grueber et

341 al. 2011).

342     Here I propose a novel and simple method to conduct model-fitting in such cases.

343 Classification trees are a powerful and flexible tool for classifying complex datasets with non-

344 linear, localized, and other complicated relationships among variables (Breiman et al. 1984,

345 Cutler et al. 2007), and they are being increasingly used by ecologists (De'ath and Fabricius

346 2000, De'ath 2002, Sullivan et al. 2006, Cutler et al. 2007, Boyer 2010, Durst and Roth 2012)

347 and paleontologists (Finnegan et al. 2012, 2015). The basic algorithm is to recursively subdivide

348 heterogeneous data sets (e.g., composed of different classes/models) into subsets that are as

349 homogeneous as possible, using whatever values of variable(s) are most powerful to do so. The

350 result, often portrayed visually as a decision tree, is a statistical model in which values of

351 predictor variables are used to classify the original class/model identities. Studies demonstrate

352 that the method (and its continuous-variable counterpart regression trees) performs as well as

353 non-linear (GAM and GLM) regressions and constrained ordination, and is exceptionally well

354 suited to scenarios involving complex (heterogeneous) classifications (Breiman et al. 1993,

355 Prasad et al. 2006, Cutler et al. 2007). Performance can be improved using "random forest"

356 variants, in which replicate trees are made, each time leaving out random subsets of variables

357 and data, until a stable model solution is achieved (Prasad et al. 2006, Cutler et al. 2007). This

358 variant also prevents overfitting of the model, and allows ranking of predictor variables by their

359 overall importance. (Use of alternate tree-building algorithms below confirmed that random

360 forests performed best for the current analyses.) When applied to new data sets (either a

361 validation/test set to test the tree's classification performance or empirical samples), each

362 resulting tree in the forest casts a vote for which class (here, model) each sample should be

363    classified. The proportion of votes for particular models, which scale from zero to one, can be

364    used as a measure of support for particular candidate models, analogous to how Akaike

365    weighting allows candidate models to be compared by their relative support (Anderson et al.

366    2000, Johnson and Omland 2004, Grueber et al. 2011).

367        The typical work flow to implement classification trees as a form of model-fitting

368    includes the following steps. See Supplementary Appendix 2 for technical details on how these

369    steps were implemented in analyses below.

370        1. Simulate training data sets: Run many simulations of the candidate models (here, the

371            four ecological diversification models) to be considered, calculating relevant

372            summary statistics (here, the eight disparity and functional diversity statistics) for

373            each iteration. This large Monte-Carlo data set serves as a training data set to train the

374            classification tree.

375        2. Build classification tree: Use random-forest classification to identify which

376            combination of variables and their statistics are most predictive of each original

377            model: Model ~ S + H + D + M + V + FRic + FEve + FDiv + FDis

378        3. Test tree with validation data set: Test the tree's performance using validation data

379            sets, additional simulation data sets of known model that were not included when

380            training the tree. A general (non-overfit) classification tree should perform

381            approximately as well on a validation data set as it did on its training data set.

382        4. Classify empirical samples using the random-forest tree: Submit empirical samples

383            (with calculated disparity /diversity statistics but unknown model identity) to the tree

384            for classification, treating proportion of votes cast for each model as a relative

385            measure of model support.

386    Methodological analyses (Supplementary Fig. 9) demonstrate that classification trees

387    perform well as model-fitting methods. When conducted across a variety of ecospace

388    frameworks (spanning different numbers of characters, seed species, and character types, and

389    incorporating weighting by empirical species pools), trees were able to classify 75–97% of their

390    training data sets correctly (Supplementary Appendix 2); validation tests showed similar

391    performance, with ~2% fewer samples being classified correctly. When applied to more complex

392    cases, such as the ability to distinguish nine models of subtly distinct dynamics (neutral model,

393    four [including both strict and relaxed partitioning versions] process-driven models with 100%

394    rule-following, and the four process-driven models with 50% rule-following), the training trees

395    still classified 90% of training data correctly (84% of validation data), even when the dynamics

396    of the models were quite similar. This performance extends across all sample sizes, with the

397    classification tree still sufficiently powerful to distinguish these models >33% of the time

398    correctly at sample sizes as small as 6 (the minimum sample size when the models could be

399    classified, as 5 species were seeded without any assembly rules). In nearly all cases, the correct

400    model also received the largest proportion of classification votes.

401    The trees were also successful in identifying samples produced from simulations quite

402    foreign from those used to train the sample (Supplementary Fig. 9). When trained on the nine

403    50% and 100% rule-following models, the resulting tree generally classified validation data sets

404    produced at 95% and 90% rule-following as members of the 100% models, with some votes for

405    the alternate 50% rules. Validation samples produced at 75% generally were classified with

406    mixed support for both the 100% and 50% models. (This generality declines at sample sizes

407    greater than 30–40 as dynamics reach stable asymptotic values; at larger sample sizes, it is

408    worthwhile to consult more complex classification trees to identify the strength of these models.)

409     These results demonstrate that classification tree votes can provide a measure of model-

410     selection support, and this support can be extended, at least in this case, beyond the particular

411     training models. Whether trained on 9 training sets (as in Supplementary Fig. 9), or on larger

412     numbers (13- or 21-model training sets, if adding 75%- or 75%-, 90%-, and 95%-strength

413     models), the correct model consistently achieves support levels greater than 0.20, from which it

414     is possible to recommend this support value as a benchmark of "substantial" support for the

415     correct model, and values above 0.40 tend to be associated with "unambiguously strong" support

416     for the correct model. (In other words, support values of 0.2 and 0.4 here can serve analogous

417     roles as the 0.1 and 0.9 Akaike-weight benchmarks used in AIC-based model selection.)

418     Although 0.4 may seem low compared to 0.9, the difference can be explained by the large

419     number of candidate models [9, 13, or 21] being considered here, and the fact that the correct

420     model in these tests often achieves support values much greater than 0.4 whereas poorly

421     supported models tend to receive far fewer votes. It is recommended that future analyses using

422     different model implementations run sensitivity analyses to validate the generality of these

423     benchmarks, as they likely depend on the number and dynamical similarity of models being

424     considered.

425     The classification-tree approach to model selection shares several similarities with the

426     methods recently termed approximate Bayesian computation (ABC) in that model selection can

427     proceed without a known likelihood function or an understanding of the dependence between

428     variables (Beaumont 2010, Slater et al. 2012). However, ABC is primarily used when estimating

429     one or a few parameters, and often within structurally similar or hierarchically nested models;

430     they are not well suited for choosing among models with distinct parameters (such as is the case

431     here). ABC also performs poorly when considering many summary statistics simultaneously, the

432  "curse of dimensionality" (Beaumont 2010) (and also the case here, where eight interdependent

433  statistics are used). In contrast, classification trees are powerful and flexible methods that were

434  developed to classify such heterogeneous data sets using many predictive variables (Breiman et

435  al. 1993). Although they have not been used before for multivariate model selection, they

436  perform remarkably well and appear to be well suited to this goal.

437

438                    **Simulation Dynamics and Sensitivity to Different Ecospace Frameworks**

439          The four models can be distinguished by examining the dynamics of the functional

440  diversity/disparity statistics as a function of species richness. These dynamics are examined here

441  across a range of ecospace-framework simulations in order to understand their general dynamics,

442  their sensitivity to different ways of creating ecospace frameworks, and their statistical power in

443  model inference when models are known to operate in simulations.

444          *Popular ecospace frameworks*.—Figure 2 illustrates the dynamics for two influential

445  ecospace frameworks, that of Bush and Bambach (Bambach et al. 2007, Bush et al. 2007, Bush

446  and Bambach 2011, Bush et al. 2011) and Novack-Gottshall (2007), which have been adapted

447  for a variety of studies (e.g., Xiao and Laflamme 2009, Bush and Novack-Gottshall 2012, Ros et

448  al. 2012, Bush and Pruss 2013, Laflamme et al. 2013, Dineen et al. 2014, Aberhan and Kiessling

449  2015, Dick and Maxwell 2015, Dineen et al. 2015, Mondal and Harries 2015, O'Brien and Caron

450  2015). Bambach's original ecospace framework (1983, 1985) consisted of three factor characters

451  (diet, activity/motility, and tiering, with the last ordered) and 11 character states, allowing 48

452  possible unique life habits. The Bush and Bambach ecospace framework retained these three

453  characters, but added additional character states (totaling 19) and reformulated motility as an

454  ordered factor. With the addition of osmotrophy as a diet category (Laflamme et al. 2009, Bush

455     and Bambach 2011, Bush et al. 2011), this yields 252 possible unique life habits. The Novack-

456     Gottshall framework more finely subdivides life-habit dimensions according to substrate

457     relationships, microhabitat, motility, diet, and foraging habits, and is modified here to include 18

458     characters (6 ordered numeric and 12 multi-state binary) and 37 character states, yielding nearly

459     3.6 trillion possible life habits. (See Supplementary Appendix 1 for list of characters and states.)

460     In all three frameworks, many of these theoretically possible life habits are unlikely to ever be

461     realized given logical or functional constraints (cf., Bambach et al. 2007).

462             Model dynamics for both ecospace frameworks are congruent, despite their rather

463     different structures. In both cases, life-habit richness (H) increases as a function of species

464     richness (S) for all models except redundancy, with the values for strict partitioning (and for

465     Bush and Bambach, also relaxed partitioning) reduced compared to the other models. The neutral

466     and expansion model slope for Bush and Bambach is also reduced slightly compared to that of

467     Novack-Gottshall, as the smaller realized ecospace becomes saturated more quickly. Trends for

468     mean distance (D) are also similar, with neutral remaining stable, expansion increasing slightly

469     (more so for Novack-Gottshall), and the others declining asymptotically to varying extents. Total

470     variance (V) is highly correlated with D for Novack-Gottshall, although the presence of factor

471     character types in Bush and Bambach precludes its measurement. Maximum range (M) is quite

472     distinct across frameworks, with all Bush and Bambach simulations quickly filling the available

473     ecospace. The Novack-Gottshall framework, in contrast, does not reach asymptotic limits within

474     50 species, with the rank order of model dynamics generally paralleling that for the other

475     disparity metrics (D and V): expansion has the greatest range values, followed closely by neutral,

476     significantly greater than both partitioning models, and with redundancy remaining constant at

477     low range. Functional richness (FRic) increases for all models, except redundancy for Novack-

478    Gottshall which decreases, with the greatest values for expansion and neutral models,

479    intermediate for strict partitioning, and least for redundancy. For Bush and Bambach, relaxed

480    partitioning produces the largest FRic values of all models, whereas it remains intermediate for

481    Novack-Gottshall. Functional evenness (FEve) remains nearly constant for neutral and expansion

482    and declines asymptotically for the other models, and the rank order is similar in both cases:

483    expansion and neutral at top, followed by both partitioning models, and redundancy at bottom.

484    Functional divergence (FDiv) has generally non-linear dynamics. Its dynamics increase in both

485    frameworks for redundancy whereas they decrease for other models, with relaxed partitioning

486    generally decreasing at the fastest rate (although eventually crossing neutral and expansion

487    models into an intermediate value at high sample sizes for Novack-Gottshall). Functional

488    dispersion (FDis) has similar dynamics to D (and V), although the dynamics for the expansion

489    and neutral models increase asymptotically instead of remaining constant.

490        That these two quite different frameworks share generally similar dynamics is good news

491    for those seeking to apply these statistics for model inference. The most important differences

492    include nearly continuous overlap between the neutral and expansion models in the Bush and

493    Bambach framework, whereas the expansion model generally has greater values and more

494    distinct dynamics for Novack-Gottshall. This is an important consideration if the goal (cf., Bush

495    et al. 2007) is to distinguish an active process (expansion) from a passive process (neutral).

496    Trend dynamics generally also have smaller error margins for the Novack-Gottshall framework,

497    which enables more powerful model selection.

498        This statistical power can be evaluated using classification tree methods, by training a

499    random forest tree on the simulation statistics and counting the proportion of an independently

500    simulated (validation) data set that is classified correctly. The Bush and Bambach framework

501    was able to correctly identify 73% of validation samples correctly. (See Supplementary

502    Appendix 2 for details and classification rates.) Greatest variable importance (in order from most

503    to moderate importance) was awarded to H, FEve, FDis, D, and FRic. This ranking is consistent

504    with Figure 2A, where these dynamics generally overlap least with those of other models.

505    Confusion matrices, which illustrate both the frequency that the correct model was selected as

506    well as which models they were misclassified, reiterate that the tree had difficulty distinguishing

507    the expansion from neutral dynamics in the Bush and Bambach framework, with ~50% of

508    expansion samples misclassified as neutral (and vice versa). In contrast, the two partitioning

509    models were correctly classified more than 80% of the time (and usually confused with the

510    alternate partitioning parametrization), and the redundancy model was nearly always properly

511    classified.

512        The Novack-Gottshall framework provided greater power to distinguish the models, with

513    89% classified correctly using a validation data set. Variable importance rankings rated FEve the

514    most important, followed in order by D, H, FDiv, and FDis, four of which were also considered

515    most important in the Bush and Bambach tree. The confusion matrix also demonstrates some

516    difficulty distinguishing the expansion and neutral models, but to a much lesser extent: ~75% of

517    expansion models were correctly distinguished from the neutral models (and vice versa), with

518    98% of partitioning and exactly 100% of redundancy models correctly classified. Thus, the more

519    multidimensional Novack-Gottshall framework, albeit more complicated and less intuitively

520    appealing for general summary applications, provides greater statistical power for inferring

521    important ecological and evolutionary models.

522        These results can be generalized by examining other ecospace frameworks, offering

523    recommendations for what forms of ecospace framework choices are most likely to provide

524 informative results for future analyses. The most important differences between these two

525 ecospace frameworks concern differences in number of characters and character types, thus it is

526 worth examining the effect of these variables. In addition, analyses also examine the effect that

527 number of seed species has on dynamics and the ability to distinguish resulting models.

528      *Number of ecospace characters.*—Ecospace frameworks with more characters and states

529 ought to have greater opportunities for the models to be distinguished, and thus be more

530 powerful frameworks. Simulations— implemented using ecospace frameworks with 5, 15, and

531 25 characters of mixed character types—bear this out, but only to a modest extent

532 (Supplementary Fig. 6). The dynamics across simulations are generally similar, in much the

533 same way as for the frameworks discussed above (Fig. 2). In particular, the smaller 5-character

534 simulation (Supplementary Fig. 6A) strongly resembles the dynamics of the Bush and Bambach

535 framework (Fig. 2A) whereas the larger 15-character simulation strongly resembles the dynamics

536 of the Novack-Gottshall framework (Fig. 2B), attesting that many of the differences result from

537 differences in the numbers of characters in each framework. However, the classification rates

538 improve only modestly as additional characters are added (83%, 85%, and 86% of models,

539 respectively), suggesting that number of characters is not the primary driver of performance

540 between the two frameworks. This conclusion reiterates one made in prior sensitivity analyses of

541 functional diversity (Maire et al. 2015), which likewise found that number of characters

542 accounted for limited improvement in performance. Examination of confusion matrices

543 (Supplementary Appendix 2) demonstrates that most of the improvement in the larger

544 simulations was the result of improved performance at distinguishing the partitioning models,

545 both from each other and from the other models, made possible because the larger number of

546 characters allows for greater opportunities for intermediate life habits.

547        *Character types.*—Another important difference between the two frameworks considered

548    above is their character types, with the Bush and Bambach framework using ordered and

549    unordered factors and the Novack-Gottshall framework using ordered and unordered numeric

550    (binary) characters. Maire et al. (2015), using sensitivity analyses involving simulated data sets,

551    concluded that character type carried the second-most importance in the performance of

552    functional-trait spaces (ecospace frameworks), second only to whether statistics were calculated

553    via a multidimensional ordination versus a functional dendogram, a comparison not considered

554    in the current study (where all functional-diversity metrics are calculated using the more

555    powerful ordination method). They demonstrated that continuous and ordinal (ordered numeric)

556    character types performed best, and discrete categorical (unordered factor) types performed

557    worst. This conclusion is partially borne out here (Supplementary Fig. 7) using simulated

558    ecospace frameworks built from 15 characters of varying types: factors, ordered factors, ordered

559    numeric, and binaries (unordered numeric). Again, dynamics are generally similar across

560    simulations, with the greatest difference occurring with ordered factors (Supplementary Fig. 7B),

561    which has rather distinct dynamics from the other frameworks for D, M, FDiv, and FDis, and to

562    a lesser extent for FRic. The dynamics for factors and ordered and unordered numeric (binaries)

563    are all quite similar. Thus, there is negligible difference caused by distance metric used, with

564    Gower distance used for both factor types and Euclidean distance for both numeric types.

565    Classification rates on ordered factors were also substantially improved (94% of trained models

566    classified correctly) compared to the other character types (78% for unordered factors, 79% for

567    ordered numerics, and 81% for binaries), with most of the difference due to improved ability of

568    the ordered-factor tree to distinguish the neutral and expansion models. (See Supplementary

569    Appendix 2 for additional details, including confusion matrices.) Overall, ordered character

570    types, and especially ordered factors, perform better than the other character types, although the

571    improvement is modest.

572         *Number of seed species in simulations.*—Actual ecological communities do not normally

573    start with pre-specified numbers of initial colonizing "seed" species; thus, selecting how many to

574    choose in a simulation should reflect knowledge of the ecological and evolutionary scenarios one

575    seeks to mimic. It is worthwhile to evaluate the effect that this arbitrary choice might have on

576    resulting simulations. Supplementary Figure 8 demonstrates that the effect can be substantial for

577    some models, although the dynamics for the neutral and expansion models are essentially

578    invariant across implementations tested here. The greatest difference occurs for FDiv, with the

579    partitioning dynamics dramatically changing behavior across the three scenarios: when seeded

580    with three species, the dynamics are relatively constant at low values; when seeded with five

581    species, the strict version decreases and the relaxed version remains constant after an initial

582    decrease; and when seeded with ten species, both trends decrease at different rates. In all cases,

583    however, the strict version has values greater than the relaxed one. The statistic M and to a lesser

584    extent FRic also vary with number of seed species, with the redundancy and partitioning values

585    substantially suppressed when seeded with three (and to a lesser extent five) species and

586    substantially increased when seeded with ten species. Although some dynamics vary

587    dramatically due to differences in seed species, the classification rates are less variable. The

588    simulation seeded with three species classified 85% of training data correctly (or 82% when

589    excluding the redundancy model, in which many statistics were absent because the functional-

590    diversity statistics require a minimum of five unique life habits), 85% for the simulation seeded

591    with five species, and 75% for the simulation seeded with ten species (see Supplementary

592    Appendix 2 for details, including confusion matrices). Confusion matrices demonstrate that the

593    neutral and expansion models are most difficult to distinguish in all three simulations, with the

594    ten-seed-species simulation also showing greater difficulty distinguishing the other models.

595    Thus, it is worthwhile—assuming the ecological and evolutionary context justifies it—using

596    moderately few numbers of species to seed the simulations because using larger numbers

597    impedes the simulations from following the model rules. Although not entirely justified by the

598    overall classification tree results, the dynamics suggest that using too few species to seed the

599    simulations can suppress the dynamics of the partitioning and redundancy models, which are

600    explicitly constrained to exploring the part of the ecospace seeded by the initially occurring

601    species.

602          *Which statistics are most valuable?*—Another benefit of using random-forest

603    classification trees as a form of model selection is that they provide an explicit mechanism to

604    evaluate which statistics are most valuable in classifying the models. As discussed earlier, this is

605    accomplished by excluding some variables at random from the tree-training algorithm

606    (somewhat analogous to bootstrap support in cladistics), which prevents overfitting of the model

607    and aids evaluating the relative variable contribution to the final tree model (Prasad et al. 2006,

608    Cutler et al. 2007, Strobl et al. 2007). Variable importance rankings for the classification trees

609    used with these simulations considered FEve the most valuable statistic, on average, for

610    distinguishing the models, followed by D, H, and FDis. In models where V could be included, it

611    was considered the fourth most important variable. FDiv and FRic consistently had intermediate

612    (but worthwhile) value, and M and S consistently had low value. The poor ranking of S is

613    reassuring, as it is the independent variable in all models, and only useful in combination with

614    other statistics. This ranking is intuitive, as the high-value statistics tended to have the least

615    overlap among models in the simulations presented above (Fig. 2, Supplementary Figs. 6–8).

616 Although functional ecologists focus almost exclusively on their newly invented statistics, those

617 traditionally used in morphological disparity are often just as powerful (or more so), and also

618 computationally more straightforward to measure. Thus, they ought to be considered more

619 frequently in functional diversity studies. However, a benefit of using classification trees as a

620 basis for model selection is that one ultimately does not have to choose which statistics to

621 include (cf., Mouchet et al. 2010). The tree algorithm is sufficiently powerful and flexible to

622 "learn" which variables to rely on in each circumstance: it always sees the forest through their

623 trees.

624

### Model Inference of Late Ordovician Samples and Effect of Geographic and

625

### Temporal Scale

626

627 *Empirical samples*.—These methods are applied to a case study of well-preserved fossil

628 biotas from the type Cincinnatian (Late Ordovician) Kope and Waynesville Formations of

629 southwest Ohio, northern Kentucky, and southeast Indiana. Lithologically, these units are

630 composed of shales interbedded with thin limestone tempestites, representing offshore, soft-

631 substrate conditions below storm-wave base in the Cincinnati Arch epeiric sea (Holland 1993).

632 Taphonomic conditions are excellent, with preservation by obrution deposits, and most

633 assemblages include evidence of fossils preserved in life position, articulated specimens, and

634 autochthonous and parautochthonous communities (Frey 1987a, b, Schumacher and Shrake

635 1997, Hughes and Cooper 1999, Aucoin et al. 2015). The shelly and trace-fossil biotas are also

636 extensively well studied, with much attention given to their life habits (Pojeta 1971, Richards

637 1972, Frey 1987a, b, 1989, Brandt et al. 1995, Lescinsky 1995, Brandt 1996, Feldmann 1996,

638 Sandy 1996, Schumacher and Shrake 1997, Leighton 1998, Gaines et al. 1999, Meyer et al.

639    2002, Morris and Felton 2003, Novack-Gottshall and Miller 2003, English and Babcock 2007,

640    Freeman et al. 2013). These two formations were selected for analysis as a case study because

641    they represent exemplary preservation within a single depositional environment, but the results

642    are likely to apply to other type Cincinnatian units.

643          218 collections from the Kope and Waynesville Formations (totaling 2322 occurrences;

644    min=2, mean=8.2, median=8, max=23 taxa) were downloaded from the Paleobiology Database

645    (www.paleobiodb.org), only including those that included the entire biota. Source references

646    include Dalvé (1948), Browne (1964), Frey (1987a, 1988, 1989), Novack-Gottshall and Miller

647    (2003), and Holland and Patzkowsky (2007). The size of each collection was confirmed to

648    represent a typical field sample (i.e., a hand sample, slab, or small bedding plane). All

649    stratigraphic members were included, as they are defined more by biofacies than by major

650    differences in depositional environment (Holland et al. 2001, Holland and Patzkowsky 2007).

651          92 Kope samples were collected from eight stratigraphic sections (roughly equivalent to

652    outcrops) spanning the Fulton, Economy, Southgate, and McMicken Members; the 126

653    Waynesville samples represented fifteen sections spanning the Fort Ancient, Clarksville, and

654    Blanchester Members and their informal "trilobite shale" and "*Treptoceras duseri* shale" facies.

655    Analyses were conducted at multiple scales: individual samples (representing autochthonous

656    local communities), stratigraphic sections, members, and each formation in aggregate. Because

657    the analyses hierarchically span multiple localities and stratigraphic levels, they offer useful tests

658    of the generality of the models. If the ecological models presented above provide good fits for

659    the functional structure within individual samples, then we might expect different structures (and

660    models) to be recorded at the spatially and temporally mixed scales of the larger aggregate levels

661    where regional and evolutionary processes are more likely to play out.

662      *Ecospace framework and ecological disparity.*—Life habits (functional traits) of 237

663    unique Kope and Waynesville taxa were coded into a modified ecospace framework from

664    Novack-Gottshall (2007) that included 18 characters (6 ordered numeric and 12 multi-state

665    binary) and 37 character states. Given the sensitivity analyses above, the large number of

666    characters and the presence of these character types ought to be sufficient for distinguishing the

667    relevant models. Life habits were coded according to inferences of functional morphology, body

668    size, ichnology, *in situ* preservation, biotic associations recording direct interactions, and

669    interpretation of geographic and depositional environment patterns, in consultation with 67 peer-

670    reviewed, published references, many studying this fauna directly. Body-size measurements

671    were made using the methods of Novack-Gottshall (2008b), primarily from Feldmann (1996),

672    supplemented with the *United States Geological Survey Professional Paper* volume 1066 series

673    and the *Treatise on Invertebrate Paleontology*. Ordered numeric states (e.g., body volume,

674    mobility, and distance from seafloor) were rescaled as five discrete, equidistant bins (seven for

675    body volume) so that they ranged from zero to unity. See Supplementary Appendix 1 for list of

676    characters and states, and an example of how life-habit (functional-trait) inferences were coded,

677    using the Ohio state fossil, the trilobite *Isotelus maximus*.

678        Of the 237 taxa in the species pool (Supplementary Appendix 3), 101 were identified to

679    (and coded at) species-level and 136 were coded at genus-level, with just 7 of these taxa having

680    uncoded (unreliably inferred) states. The relative completeness of this coding is important, as

681    sensitivity analyses have demonstrated that missing data (traits or species) can lead to inaccurate

682    functional-diversity measurements (Pakeman 2014); this impact is also mitigated here by having

683    replicate samples and the use of a model-selection routine that explicitly incorporates variation

684    among simulation trials. Indeterminate taxa (e.g., trepostome bryozoan indet. or *Platystrophia*

685   sp.) that occurred within individual samples (but were excluded from the aggregate 237-taxon

686   species pool unless their occurrence was the sole member of that taxon, see below) were coded

687   for a particular state only when all other members of that taxon within the Kope-Waynesville

688   species pool unanimously shared that common state; otherwise, the state was listed as NA

689   (missing). Although there are differences in taxa and functional traits between the Kope and

690   Waynesville Formation species pools (especially resulting from the Richmondian invasion,

691   Holland and Patzkowsky 2007, 2009), a single aggregate pool is used here for simplicity. Use of

692   formation-specific species-pool models do not alter the current results.

693        Eight ecological disparity and functional diversity statistics were calculated for the 218

694   Kope and Waynesville samples, and for section, member, and formation-level aggregates of

695   these samples. The routine used to produce temporal and geographical aggregates only retained

696   indeterminate taxa if no other members of the same taxon were found in the aggregate; this logic

697   also applied to genus-level occurrences, excluding them when a named species of the genus was

698   present. This left 135 unique taxa (representing 83 unique life habits) in the Kope Formation

699   aggregate and 174 unique taxa (93 unique life habits) in the Waynesville Formation aggregate.

700   Eight samples were excluded from model-fitting protocols because they had fewer than four life

701   habits, the minimum needed to calculate the functional-diversity statistics.

702        *Simulations, constraints, and model selection.*—Two criteria must be met when

703   considering candidate models for these samples (Burnham and Anderson 2002, Johnson and

704   Omland 2004). The first concerns adequacy: do the models predict values of similar magnitude

705   to those found in the samples? The second concerns model selection: among adequate alternative

706   candidates, which model best represents the empirical samples? "Best" in this context can be

707   defined in multiple ways, such as model simplicity, maximum likelihood, or Chi-square or

708    Kolmogorov-Smirnoff statistics (Burnham and Anderson 2002), but classification trees trained

709    (and validated) on simulated data sets are used here because of the non-independent,

710    multidimensional nature of the model statistics.

711         Five seed species were used in the simulations, as this approximated (the $12^{th}$ percentile)

712    the sample size of the smallest samples, and this value allows calculation of all functional-

713    diversity statistics and provides powerful model discrimination using classification trees. To

714    obtain adequate models, simulations were implemented iteratively, altering constraints of the

715    ecospace framework to produce models with statistics similar to those found in empirical

716    samples and aggregates (Fig. 3). Note that simulations proceeded in identical manner in each

717    simulation; only constraints delimiting the nature of the theoretically possible ecospace were

718    modified. As discussed below, the alterations needed to achieve realistic statistics are themselves

719    informative as to the types of constraints that may exist in defining the realized

720    Kope/Waynesville ecospace.

721         The simulations used when building the models in Figure 2B used an ecospace

722    framework that allowed at most "two presences" for each binary character ('constraint=2'; see

723    "Theoretical Ecospace Framework" above). These models all fail the first test of adequacy.

724    Except for H, FEve and FDiv, the model simulations predict disparity and functional diversity

725    statistics much larger than the values found in the actual Kope and Waynesville samples

726    (compare with points in Fig. 3A). It is unwise to select the best model when all perform poorly.

727    The inadequacy is attributable to the relatively unconstrained framework, which allows for 930

728    billion unique life habits to occur in simulations, many of which are unlikely or impossible to

729    occur in nature (such as animals that feed using any two foraging combinations among ambient,

730    filter, attachment, mass, and raptorial habits). The highly multidimensional nature of this

731    framework provides so many opportunities for unrealistic life habits that empirical samples are

732    bound to be depauperate by comparison.

733         Simulations in Figure 3A further constrain the framework so that life habits can inhabit

734    only a single binary character state ('constraint=1'; e.g., only infaunal or epifaunal are allowed

735    but not semi-infaunal). This additional constraint diminishes all statistics by a negligible amount,

736    but the overall dynamics are largely unchanged, with the largest differences occurring in the

737    partitioning simulations. Empirical samples continue to have statistics substantially below those

738    predicted by these models, implying that even this more constrained framework—allowing only

739    1.13 billion unique life habits—is too permissive to represent life habits that actually occur in

740    nature. For example, it "permits" animals the size of humans to crawl between grains of

741    sediment, a logically impossible life habit. However, many of the life habits created by this

742    framework are biologically possible, even if they require some creativity to envision an organism

743    inhabiting such a life habit. As an example, the framework could happen upon a "buried"

744    microscopic, fluid-feeding carnivore that feeds facultatively attached to its prey on hard

745    substrates while living far above the sea floor. Although seemingly implausible, this is a typical

746    life habit for numerous parasites (cf., Huntley and Scarponi 2012). The exposure of such unlikely

747    (or at least rarely fossilized) life habits (see Novack-Gottshall 2007 for others) is not a flaw or

748    bias in the simulations or analyses, and can reveal important constraints that structure actual

749    biotas (Raup 1966, Seilacher 1970, Thomas and Reif 1993).

750         The generally unchanged dynamics in Figures 3A and 2B, despite their ability to

751    represent orders-of-magnitude different numbers of life habits, also demonstrate that the

752    inclusion of large numbers of potentially illogical or biologically unreasonable life habits has

753    little effect on resulting dynamics. The imposed "single-presence" constraint is also biologically

754    unrealistic, as 52% of the Kope/Waynesville taxa are coded using "two-presence" characters

755    (e.g., corals reproduce both sexually and asexually, semi-infaunal bivalves are common, and

756    many crinoids, brachiopods, and bryozoans are known to attach to lithic and biotic substrates).

757    Different kinds of constraints must be imposed to allow the models to approximate actual

758    Ordovician samples.

759         The next simulations relax the "one-presence" constraint to allow for "two-presence"

760    characters, but use the Kope/Waynesville aggregate pool to weight the inhabitation of character

761    states. For example, because 87% of the pool taxa are coded as epifaunal, 9% as infaunal, and

762    4% as semi-infaunal, simulated life habits within the neutral model (and seed species for all

763    models) will mimic these proportions, on average. Similarly, because no taxa are coded as

764    simultaneously autotrophic and carnivorous, that combination is removed as a possible character

765    state across all simulations, as are other such non-occurring ones. Imposing such empirical

766    weighting still allows for 57 billion unique life habits, and allows the ecospace framework to

767    mimic important features of the empirical species pool, albeit indirectly. For example, although

768    the simulation rules do not specify actual predator:prey or producer:consumer trophic-group

769    ratios (Van Valkenburgh 1988, Dunne et al. 2008, Mitchell et al. 2012, Hatton et al. 2015), size-

770    diet relationships (Van Valkenburgh and Molnar 2002, Codron et al. 2013), or other regulating

771    factors, the weighting allows the simulation to indirectly mirror these empirical ecological rules.

772    As a result of this weighting, most of the simulated life habits also correspond to biologically

773    reasonable life habits commonly inhabited by animals, living or extinct. Neutral simulations

774    yield life habits differing by ~4 states (out of 37 in the framework) from ones actually occurring

775    in the Late Ordovician species pool (after allowing the simulation to proceed through completion

776    at 50 species per sample). The dynamically restricted models produce life habits differing by less

777    than 1 state from actual Ordovician life habits for redundancy (with 95% rule-following) and

778    strict partitioning simulations, and ~2 states for the relaxed partitioning simulations. The

779    expansion model, as intended, explores the most disparate (and slightly unreasonable) life habits,

780    but still produces life habits that differ by ~10 states from actual Ordovician life habits. Thus, the

781    life habits produced by the simulations, when weighted by empirical species pools, are

782    biologically quite reasonable approximations of reality. (It should be emphasized, however, that

783    the model-selection method only considers the disparity structure of each assemblage, and not

784    the identity of simulated life habits.) Further analysis of the nature of these constraints is beyond

785    the scope of this study, but it seems reasonable that the implementation of empirical weighting is

786    a practical solution to restrict the theoretical ecospace to those approximate regions locally

787    allowed by whatever combination of historical, biological (adaptational), and structural

788    constraints shaped the species pool (Raup 1966, Seilacher 1970, Thomas and Reif 1993).

789          Implementing these weightings also result in simulations (Fig. 3B) that reflect adequately

790    the statistics of empirical samples across all eight statistics, with dynamics substantially

791    diminished compared to the unweighted simulations. Statistical dynamics are also quite altered,

792    although the rank order of models remains generally unchanged: expansion (also the least

793    altered, with asymptotic values negligibly diminished from prior unweighted implementations)

794    has the largest disparity/diversity values, followed by neutral (now dynamically distinct from

795    expansion), the two partitioning implementations, and with redundancy at lowest values. The

796    distinct model dynamics also result in improved performance for the classification tree, with

797    96% of the training data sets classified correctly (compared to 91% for the one- and two-

798    constraint simulations), primarily resulting from improved ability to distinguish the expansion

799    and neutral models (see Supplementary Appendix 2 for details).

800    Although the empirical samples adequately overlay model dynamics, they follow the

801    mean trends only rarely, occurring more frequently in the region between the passive neutral

802    model and the other active models. This could indicate that the best models to represent the

803    samples are not 100% implementations of the models, but instead weakened versions of the

804    models. This is illustrated in Figure 3C where eight additional "weaker" simulations are added,

805    in which the model rules are followed on average 90% and 50% of the time. In other words, for

806    the 50%-strength implementation, a computational "coin flip" was used in the assignment for

807    each taxon's character states, whether to assign it at random (weighted by the species pool) or

808    whether to use the model rules to assign it (but still only allowing states present in the species

809    pool). The neutral model represents a 0% simulation, in which no active model rules are

810    followed, and model dynamics coalesce toward the neutral model as they become weaker and

811    weaker. Most empirical samples overlay the weakened model dynamics. Although many of the

812    nine models share similar dynamics (more so given the often overlapping variation around each

813    mean trend line), the classification tree displays a remarkable ability to distinguish them, with

814    90% of training data sets classified correctly, and 73% correct when tested against independently

815    created validation data sets.

816    The classification tree used below in model selection of samples is simpler, using just the

817    100% and 50% rule-following models in Figure 2C as training data sets; it classified 92% of

818    training data sets and 78% of validation data sets correctly (Supplementary Appendix 2).

819    Although the training data sets here represent just two strength levels (50% and 100%),

820    validation samples produced at different strengths (75%, 90%, and 95%) indicate the simple

821    "two-strength" classification tree is able to provide generalizable approximations for other

822    model-strength implementations (Supplementary Appendix 2; Supplementary Fig. 9).

823   Performance remains strong at sample sizes as low as six, the minimum sample size when the

824   models could be classified (i.e., because five were seeded without any assembly rules). (It is

825   reassuring that if the simulated assemblages with five or less species are submitted to model

826   selection, they are overwhelmingly, and correctly, classified as representing neutral models.)

827   More complicated trees (i.e., those trained on 75%, 90%, and 95%-strength models) are used to

828   confirm the strength of models, especially for larger samples, where these more complex trees

829   are better suited for distinguishing their less dissimilar dynamics.

830          The majority of Kope and Waynesville samples are classified (Fig. 4, Supplementary

831   Appendix 4) as representing partitioning models (73% and 74% respectively), with most

832   classified as 100%-strength versions of the model (51% and 43%, respectively). Some samples

833   are classified as weak (50%-strength) redundancy models (14% and 16% respectively), and

834   negligible numbers are classified as neutral or expansion models. A Chi-squared test confirms

835   that the model classifications are statistically indistinguishable between formations ($X^2$=48,

836   df=36, $p$-value=0.087). Average model support (measured in terms of proportion of tree-

837   classification votes cast) for the "best" model is 0.50, above the threshold of 0.20 in which

838   support can be considered substantial and 0.40 when it can be considered unambiguously strong

839   (see Supplementary Fig. 9). Only one sample had winning support below this threshold, a small

840   seven-taxon Waynesville sample deemed neutral with support of 0.18. The relatively low

841   numbers of samples classified as weakened partitioning models (less than 5% of samples) also

842   suggests that most Kope and Waynesville samples are best represented by strong (>75%, and

843   likely >90%) partitioning models, based on sensitivity analyses in Supplementary Figure 9. This

844   is confirmed using more complicated classification trees: when additional strength simulations

845   are included (i.e., when adding 90% as in Fig. 3C, or adding 75%, 90%, and 95% to the training

846    data sets), most samples remain classified as strong (90, 95% and 100%) versions of strict

847    partitioning, 100%-strength relaxed partitioning, and 75% and 90%-weakened redundancy

848    models (Supplementary Appendix 4). Classification trees and simulations that used separate

849    species pools for each formation (i.e., creating a Kope tree trained only on Kope species pool

850    simulations) also produced similar model classifications, as expected given that 42% of genera

851    and 69% of life habits are identical in each separate pool.

852          Representative samples classified by the four driven models are illustrated in Figure 5 on

853    two-dimensional non-metric multidimensional scaling plots. As expected, samples classified as

854    partitioning models demonstrate continuous gradations, either in a linear fashion for the strict

855    partitioning version of the model or in a central cluster in the relaxed version. Samples classified

856    as weakened redundancy demonstrate discrete clusters of life habits, often with multiple taxa

857    sharing similar or identical life habits. At large sample sizes, both the partitioning and

858    redundancy models yield many instances of taxa with nearly overlapping life habits, although

859    those produced by partitioning tend to have more continuous gradients between these clusters

860    (compare Figs. 1 and 5). Samples classified as expansion tend toward relatively large distances

861    between life habits, often have the centroid empty, and encompass a broad range of the ecospace,

862    even at small sample sizes.

863          If the model classifications indicate evidence for ecological structure within these

864    generally autochthonous samples, then one might expect different models to emerge at larger

865    temporal and spatial scales, where samples are aggregated into individual outcrops (sections) and

866    stratigraphic members and formations (Patzkowsky and Holland 2003, Holland 2010,

867    Tomašových and Kidwell 2010, Hautmann 2014). For example, if beta diversity is substantial,

868    one might expect a preponderance of redundancy models, as samples with similar structure

869    duplicate each other in aggregation. Alternatively, one might expect greater frequency of neutral

870    models, as patterns at the larger scales increasingly reflect the regional species pool.

871         Such a transition is partially borne out here. The Kope and Waynesville stratigraphic

872    sections have model classifications similar to those of their individual samples (Fig. 4),

873    dominated by 100%-strength partitioning (albeit only the strict implementation) and the

874    remainder largely weak redundancy. The sole section deemed neutral has the smallest sample

875    size and lowest (but still substantial) model support of 0.34. These classifications are generally

876    supported when using more complex classification trees, with all eight Kope sections best

877    classified as 90–100% strict partitioning models, but a greater number of Waynesville sections

878    classified as 50–90% redundancy models (Supplementary Appendix 4).

879         At the larger sample sizes of stratigraphic members and formations, it is possible to use

880    more complex classification trees (i.e., adding 90%-strength training data sets or adding 75%,

881    90%, and 95- strength implementations) to conduct model selection. This is especially

882    worthwhile because the dynamics of redundancy and both partitioning models are similar if even

883    slightly weakened (i.e. 90–99%). This is evident by their overlapping dynamics in Figure 3C and

884    when examining model classifications and confusion matrices on validation data sets

885    (Supplementary Fig. 9, Supplementary Appendix 2). In particular, slightly weakened (90%-95%)

886    redundancy models are typically misclassified by the simple tree as strong strict partitioning

887    models at sample sizes greater than 30 (Supplementary Fig. 9), and moderately weak (75–90%)

888    relaxed partitioning models are typically misclassified as strong relaxed partitioning and weak

889    redundancy models (Supplementary Fig. 9). Note that these misclassifications tend to only occur

890    at sample sizes greater than 30, where the more asymptotic dynamics make it easier for the

891    simple classification tree to reject implementations that are different from those in the training

892    data set. One might argue that these classifications point to indistinguishable dynamics for these

893    models, regardless of whether named "redundancy" or "partitioning." Yet, the simple tree is able

894    to correctly distinguish these models more than 33% of the time (Supplementary Fig. 9). In

895    contrast, the more complex three-strength (50%, 75%, and 100%) and five-strength (50%, 75%,

896    90%, 95%, and 100%) trees classified 90% of training data sets correctly, and displays improved

897    ability to distinguish all weakened implementations of each model.

898         When using the five-strength tree, all four Kope members are classified as strong (90–

899    95%) strict partitioning models, consistent with model classifications at smaller scales (Fig. 4).

900    The Waynesville members have divided support, with three best classified as 90–95% strict

901    partitioning and two as 90% redundant models (Supplementary Appendix 4). Classifications

902    using the three-strength tree are identical, but with the 95% votes classified instead as 90%-

903    strength implementations.

904         The entire Kope Formation species pool, using the five-strength tree (Fig. 4), is classified

905    as 90% strict partitioning (0.52 support), with substantial support for 90% redundancy model

906    (0.34), whereas the Waynesville formation pool is classified as 90% redundancy (0.78 support)

907    with the remainder for 95% strict partitioning. The three-strength tree provides similar

908    classifications (Supplementary Appendix 4). This classification is intuitively evident when

909    considering the ratio of unique life habits to taxa (H/S) in each formation: the Waynesville pool

910    has substantially more true redundancy (81 of 174 taxa, or 47%, are functionally identical)

911    compared to the Kope (with 52 of 135 taxa, or 37% functionally identical). The structure of each

912    formation is visualized in Figure 5. Although they look quite similar (the most obvious

913    difference being the presence of corals in the Waynesville, one of many new immigrants during

914    the Richmondian invasion, Holland and Patzkowsky 2007, 2009), the classification trees are

915    sufficiently sensitive to their structural differences to classify them accordingly.

916        When comparing classifications across spatially and temporally mixed scales, the Kope

917    Formation remains consistently classified by strong (>90%) strict partitioning models. There is

918    no evidence of change in ecospace structure across scales. In contrast, the Waynesville is

919    typically classified as representing partitioning models (both relaxed and strict versions) at the

920    scale of samples, switching to greater support for ~90% redundant models at larger scales,

921    especially when considering the entire formation species pool. The consistent lack of support for

922    neutral and expansion models at any scale is evidence that the ecospace structure of these Late

923    Ordovician samples is decidedly nonrandom and restricted. Even when weighted by the life-

924    habit traits known to occur in the Late Ordovician regional species pool, the ecospace actually

925    inhabited by this biota remains substantially constrained.

926        Although the simulations do not incorporate phylogenetic structure (newly added species

927    are added solely on account of their functional traits), it is promising that model selection is not a

928    simple artifact of which taxa are present within a sample (Fig. 5). If it were, one might expect

929    that samples classified as redundancy or partitioning be dominated by one or a few taxonomic

930    groups, whereas samples classified as expansion might have greater taxonomic diversity. That is

931    not the case here, where a wide range of major taxonomic groups is present at all scales, and

932    where clusters of functionally similar species likewise tend to be taxonomically diverse. The

933    same ecological structure reoccurs when taxonomically quite different biotas are analyzed.

934        Ecological and evolutionary theory underlying these models (see companion article

935    Novack-Gottshall 2016) allows speculation as to the nature of processes structuring these biotas.

936    At the scale of individual samples, partitioning dynamics are consistent with ecological niche-

937     partitioning: co-existing taxa have similar (but not identical) life habits. Holland and Patzkowsky

938     (2007, Patzkowsky and Holland 2007) demonstrated similar tight packing of species along

939     onshore-offshore gradients in these intervals, even more so in the Waynesville samples. This

940     partitioning could be manifested at the scale of local communities or through evolutionary

941     processes at the regional scale, in which speciation produces multiple variants with slightly

942     different specializations. Consideration of phylogenetic relatedness could offer a useful test of

943     these claims (Gerhold et al. 2015), in particular whether newly evolving taxa within the regional

944     species pool were more likely to have life habits different from those previously existing.

945     Patzkowsky and Holland (2003) found no evidence for saturation within Cincinnatian samples,

946     suggesting regional processes as the more important factor. This conclusion is tentatively

947     supported here, especially for the Kope Formation, where the same model support across all

948     scales could suggest a shared cause. Within the Waynesville sequence, Patzkowsky and Holland

949     (2007) identified greater beta diversity, driven in large part by faunal incursions during the

950     Richmondian invasion. This could explain the redundancy dynamics found within the

951     Waynesville Formation and its members, in which the accumulation of many different life habits

952     within individual samples results in functional duplication when aggregated.

953          These results, however, do not rule out an important role for local processes. Although

954     not illustrated here, individual Kope and Waynesville samples, even those classified by the same

955     model, do not share the same life-habits (or taxa). In other words, when plotted on a common

956     ordination (such as for the partitioning samples in Fig. 5), taxonomic and functional composition

957     varies substantially among samples. Model selection results suggest simply that most samples

958     share a similar level of ecological disparity (Fig. 4), one best represented by the partitioning

959     model. This consistency of structure suggests that local processes could still play an important

960    role in regulating these communities. Among the many possible life habits available within the

961    regional species pool, individual communities were preferentially composed of groups of

962    generally similar but non-identical life habits, although the particular life habits present in any

963    setting could vary a great deal.

964

965                    **Relevance to broader Phanerozoic trends in ecospace utilization**

966            The general correspondence between empirical samples and the dynamics of the model

967    simulations suggests that the simulations provide reasonable null models for understanding how

968    ecospace is structured in extant and fossil biotas, at many spatio-temporal scales. Support for the

969    partitioning model, and to a lesser extent redundancy, indicates that these Late Ordovician biotas

970    had relatively constrained ecospace structures. Broader discussion of this model support is

971    warranted, especially within the context of how these results might generalize to understanding

972    Phanerozoic-wide trends in ecospace utilization. In particular, one might want to know whether

973    these results are consistent with an expansion in ecospace utilization in later Phanerozoic biotas,

974    as has been widely claimed (Bambach 1983, Vermeij 1987, Knoll and Bambach 2000, Bush et

975    al. 2007, Novack-Gottshall 2007, Bush and Bambach 2011, Vermeij 2011, 2013). A conclusion

976    cannot be made at this point, but speculation might be made based on the interactions of three

977    factors: changes in the species pool used (which changes allowable ecospace traits and resulting

978    simulation dynamics), measures of functional diversity/disparity in later biotas, and changes in

979    taxonomic richness.

980            The most important of these factors is the species pool one chooses to use (cf., Cornell

981    1999, de Bello 2012). The analyses above used a narrowly defined one, the aggregate pool of

982    species known from a single habitat in one small (tri-state) region during a short (~8 m.y.) time

983    interval. Analyses using a Kope-only versus Waynesville-only pool had negligible effect, but

984    future analyses should test the sensitivity of simulation dynamics to more distinct species pools.

985    Use of a much larger pool would allow for a greater range of life-habits in the framework, which

986    could increase dynamical values in simulations (although this remains to be determined). The

987    effect of this change might be to increase classification-tree support for the partitioning and

988    redundancy models, and especially so if using a Paleozoic-wide or Phanerozoic-wide species

989    pool, given the greater range of epifaunal and infaunal tiering (Ausich and Bottjer 1982, Bottjer

990    and Ausich 1986) and body sizes (Novack-Gottshall 2008a, Heim et al. 2015, Zhang et al. 2015)

991    evolved by later biotas. Because the analyses above rescaled these ordered numeric characters,

992    incorporating broader ranges would depress the empirical disparity statistics calculated here to

993    some extent. However, the effect of this is likely minor, as post-Ordovician ranges would only

994    add a single bin (given the logarithmic binning scale used in the framework) and the largest,

995    most deeply infaunal, and tallest tiering animals were likely proportionally minor components of

996    these biotas (which diminishes their effect for these extreme character states because of the

997    empirical weighting used). It remains to be determined what pool is best suited to such analyses

998    given the major evolutionary changes throughout the Phanerozoic, but the use of more inclusive

999    pools would likely demonstrate that these Ordovician samples are even more functionally

1000    restricted (i.e., less ecologically disparate) compared to later Phanerozoic samples.

1001        Species richness has also increased during the Phanerozoic, both globally (Sepkoski

1002    1981, Alroy et al. 2008) and within individual assemblages (Bambach 1977, Powell and

1003    Kowalewski 2002, Bush and Bambach 2004, Kowalewski et al. 2006). The effect of this

1004    increase, by itself, is also likely to be limited here because the dynamics tend to reach asymptotic

1005    values at moderate (>20) species richness values, which means that simply increasing sample

1006    sizes, all things equal, will yield negligible changes to statistical values. The greater effect will

1007    be related to how these larger biotas utilize ecospace. Novack-Gottshall (2007) demonstrated that

1008    extant marine biotas had greater ecological disparity (measured as D, after accounting for

1009    differences in sample sizes) than Early–Middle Paleozoic ones, despite insignificantly greater

1010    number of life habits. Increasing values for statistics D and H through time could provide greater

1011    support for the expansion model for such biotas, although it remains uncertain whether this

1012    increase would be attenuated by using a larger species pool. It is worth reiterating that one would

1013    expect these statistics to increase with increasing species richness, even if ecospace was occupied

1014    by a passive process (i.e., that expected by the neutral model). The discussion above includes

1015    many subjective hypotheticals, and so this author is unwilling to make a wager on the outcome.

1016    Without conducting formal simulations and statistical analyses, it simply remains too early to

1017    predict whether the Phanerozoic trend would be better supported by the driven expansion model,

1018    the passive neutral one, or perhaps another. It seems a worthwhile goal to conduct such analyses.

1019

1020                                    **Conclusions**

1021        Despite documentation of synoptic paleoecological trends across the Phanerozoic and

1022    speculation about their causes, we lack, in many—perhaps most—cases, specific quantitative

1023    claims on ecospace inhabitation that can be tested analytically. The models and analytical

1024    methods proposed here, none of which are entirely novel, offer potentially fruitful avenues for

1025    such tests, whether applied to individual assemblages or to the entire biosphere throughout the

1026    Geozoic. In particular, the analyses presented here, and in an accompanying article, lead to the

1027    following conclusions.

1028    1.  Simulations of the redundancy, partitioning, expansion, and neutral models demonstrate

1029        dynamical consistency in functional diversity/disparity statistics across different data

1030        structures (ecospace frameworks), number and type of functional traits (characters), and

1031        implementations of the simulations. However, ecospace frameworks with greater numbers of

1032        functional traits, use of ordered factors, and modest numbers of seed species tend to be

1033        statistically more powerful for differentiating the models, especially the dynamically similar

1034        neutral and expansion models and the often-similar partitioning and redundancy models. The

1035        'ecospace' package provides R functions to conduct simulations for any ecospace framework

1036        and to calculate a wide range of functional diversity/disparity statistics.

1037    2.  Classification trees are a powerful method for rigorously classifying these models in a multi-

1038        model inference framework (Breiman et al. 1993, Cutler et al. 2007), with relative support

1039        allocated to candidate models according to proportion of votes from the classification tree.

1040        Classification trees are successful in identifying models correctly, even when the statistical

1041        dynamics are similar and when tested with foreign data sets unlike those used to train the

1042        tree. The trees are also able to identify which statistics are most valuable. This method

1043        identifies number of unique life-habits (functional trait combinations, H) as the most

1044        important statistical discriminant, followed by functional evenness (FEve) and functional

1045        dispersion (FDis). However, because all metrics retain useful information on ecospace

1046        structure and tree algorithms perform well using large number of predictive variables, it is

1047        recommended that analyses of ecological disparity/functional diversity use all statistics when

1048        conducting model-selection analyses.

1049    3.  Comparison of stochastic simulation dynamics to those of Ordovician empirical samples

1050        demonstrates that actual fossil assemblages are substantially constrained in their inhabitation

1051    of life habits, compared to what is possible in the theoretical ecospace framework. Although

1052    the identity of constraints are not analyzed here specific, they likely reflect a combination of

1053    logically impossible trait combinations, maladaptive strategies, inherent covariation among

1054    functional traits, ecologically meaningful restrictions to the regional species pool, and other

1055    factors. Incorporating constraints to the ecospace framework (such as limiting allowed life-

1056    habit combinations and weighting functional traits by those occurring in the empirical species

1057    pool) causes simulation dynamics to converge on more realistic values, allowing simple

1058    approximations for many of these constraints. However, doing so still demonstrates that these

1059    Late Ordovician biotas had substantially constrained ecospace structures.

1060    4. Empirical application of the classification trees demonstrate that Late Ordovician (type

1061    Cincinnatian) samples from the Kope and Waynesville Formations are primarily best fit by

1062    partitioning models. When larger stratigraphic and temporal aggregates are analyzed, the

1063    entire Kope Formation pool remains best fit by the partitioning model, but the aggregate

1064    Waynesville pool is better fit by the redundancy model. This structural transition in the

1065    Waynesville Formation can be biologically interpreted by greater beta diversity in this unit,

1066    likely related to faunal incursions caused by the Richmondian invasion. However, the

1067    consistency of support for the partitioning model at small scales suggests an important role

1068    for local processes.

1069    5. Most hypotheses regarding patterns in ecospace utilization across the Phanerozoic are

1070    superficially consistent with multiple models of ecological diversification, despite being

1071    caused by distinct processes. Most documented trends are equally consistent with passive

1072    processes. Statistical tests that consider alternative stochastic models must be conducted

1073    before we can confidently claim that ecological patterns across the Geozoic history of life

1074 had driven causes. Similar concerns are shared with ongoing research in community ecology

1075 and functional ecology.

1076

1091

1092           **Literature Cited**

1093

1094 Aberhan, M., and W. Kiessling. 2015. Persistent ecological shifts in marine molluscan

1095   assemblages across the end-Cretaceous mass extinction. Proceedings of the National

1096   Academy of Sciences 112(23):7207-7212.

1097 Ackerly, D. D., and W. K. Cornwell. 2007. A trait-based approach to community assembly:
1098     partitioning of species trait values into within- and among-community components.
1099     Ecology Letters 10(2):135-145.
1100 Alroy, J., M. Aberhan, D. J. Bottjer, M. Foote, F. T. Fursich, P. J. Harries, A. J. W. Hendy, S. M.
1101     Holland, L. C. Ivany, W. Kiessling, M. A. Kosnik, C. R. Marshall, A. J. McGowan, A. I.
1102     Miller, T. D. Olszewski, M. E. Patzkowsky, S. E. Peters, L. Villier, P. J. Wagner, N.
1103     Bonuso, P. S. Borkow, B. Brenneis, M. E. Clapham, L. M. Fall, C. A. Ferguson, V. L.
1104     Hanson, A. Z. Krug, K. M. Layou, E. H. Leckey, S. Nürnberg, C. M. Powers, J. A. Sessa,
1105     C. Simpson, A. Tomašových, and C. C. Visaggi. 2008. Phanerozoic trends in the global
1106     diversity of marine invertebrates. Science 321(5885):97-100.
1107 Anderson, D. R., K. P. Burnham, and W. L. Thompson. 2000. Null hypothesis testing: problems,
1108     prevalence, and an alternative. Journal of Wildlife Management 64(4):912-923.
1109 Anderson, M. J., K. E. Ellingsen, and B. H. McArdle. 2006. Multivariate dispersion as a measure
1110     of beta diversity. Ecology Letters 9(6):683-693.
1111 Aucoin, C. D., B. F. Dattilo, C. E. Brett, and D. L. Cooper. 2015. Preliminary report on the
1112     Oldenburg "butter shale" in the Upper Ordovician (Katian; Richmondian) Waynesville
1113     Formation, USA. Estonian Journal of Earth Sciences 64:3-7.
1114 Ausich, W. I., and D. J. Bottjer. 1982. Tiering in suspension feeding communities on soft
1115     substrata throughout the Phanerozoic. Science 216:173-174.
1116 Bambach, R. K. 1977. Species richness in marine benthic habitats through the Phanerozoic.
1117     Paleobiology 3(2):152-167.
1118 ---. 1983. Ecospace utilization and guilds in marine communities through the Phanerozoic. Pp.
1119     719–746. *In* M. J. S. Tevesz, and P. L. McCall, eds. Biotic Interactions in Recent and
1120     Fossil Benthic Communities. Plenum, New York.
1121 ---. 1985. Classes and adaptive variety: the ecology of diversification in marine faunas through
1122     the Phanerozoic. Pp. 191–253. *In* J. W. Valentine, ed. Phanerozoic Diversity Patterns:
1123     Profiles in Macroevolution. Princeton University Press, Princeton, NJ.
1124 Bambach, R. K., A. M. Bush, and D. H. Erwin. 2007. Autecology and the filling of ecospace:
1125     key metazoan radiations. Palaeontology 50(1):1-22.
1126 Beaumont, M. A. 2010. Approximate Bayesian computation in evolution and ecology. Annual
1127     review of ecology, evolution, and systematics 41(1):379-406.

1128      Bottjer, D. J., and W. I. Ausich. 1986. Phanerozoic development of tiering in soft substrata
1129          suspension-feeding communities. Paleobiology 12:400-420.

1130      Boyer, A. G. 2010. Consistent ecological selectivity through time in Pacific Island avian
1131          extinctions. Conservation Biology 24(2):511-519.

1132      Brandt, D. S. 1996. Epizoans on *Flexicalymene* (Trilobita) and implications for trilobite
1133          paleoecology. Journal of Paleontology 70:442-449.

1134      Brandt, D. S., D. L. Meyer, and P. B. Lask. 1995. *Isotelus* (Trilobita) "hunting burrow" from
1135          Upper Ordovician strata, Ohio. Journal of Paleontology 69:1079-1083.

1136      Breiman, L. 2006. randomForest: Breiman and Cutler's random forests for classification and
1137          regression, Version 4.6-10. cran.r-project.org/web/packages/randomForest.

1138      Breiman, L., J. H. Friedman, R. A. Olshen, and C. J. Stone. 1984. Classification and Regression
1139          Trees. Wadsworth & Brooks.

1140      ---. 1993. Classification and Regression Trees. Chapman and Hall/CRC press.

1141      Browne, R. G. 1964. The coral horizons and stratigraphy of the Upper Richmond Group in
1142          Kentucky West of the Cincinnati Arch. Journal of Paleontology 38(2):385-392.

1143      Burnham, K. P., and D. R. Anderson. 2002. Model Selection and Multi-Model Inference: A
1144          Practical Information-Theoretic Approach. Springer, New York.

1145      Bush, A. M., and R. K. Bambach. 2004. Did alpha diversity increase during the Phanerozoic?
1146          Lifting the veils of taphonomic, latitudinal, and environmental biases. The Journal of
1147          geology 112(6):625-642.

1148      ---. 2011. Paleoecologic megatrends in marine Metazoa. Annual Review of Earth and Planetary
1149          Sciences 39:241-269.

1150      Bush, A. M., R. K. Bambach, and G. M. Daley. 2007. Changes in theoretical ecospace utilization
1151          in marine fossil assemblages between the mid-Paleozoic and late Cenozoic. Paleobiology
1152          33(1):76-97.

1153      Bush, A. M., R. K. Bambach, and D. H. Erwin. 2011. Ecospace utilization during the Ediacaran
1154          radiation and the Cambrian eco-explosion. Pp. 111-134. *In* M. Laflamme, J. D.
1155          Schiffbauer, and S. Q. Dornbos, eds. Quantifying the Evolution of Early Life: Numerical
1156          Approaches to the Evaluation of Fossils and Ancient Ecosystems. Springer, New York.

1157 Bush, A. M., and P. M. Novack-Gottshall. 2012. Modelling the ecological-functional
1158      diversification of marine Metazoa on geological time scales. Biology Letters 8(1):151-
1159      155.

1160 Bush, A. M., and S. B. Pruss. 2013. Theoretical ecospace for ecosystem paleobiology: energy,
1161      nutrients, biominerals, and macroevolution. Pp. X-XX. *In* A. M. Bush, S. B. Pruss, and J.
1162      L. Payne, eds. Ecosystem Paleobiology and Geobiology. Short Courses in Paleontology
1163      19. Paleontological Society and Paleontological Research Institute, Ithaca, NY.

1164 Ciampaglio, C. N., M. Kemp, and D. W. McShea. 2001. Detecting changes in morphospace
1165      occupation patterns in the fossil record: characterization and analysis of measures of
1166      disparity. Paleobiology 27(4):695-715.

1167 Codron, D., C. Carbone, and M. Clauss. 2013. Ecological interactions in dinosaur communities:
1168      influences of small offspring and complex ontogenetic life histories. PLoS One
1169      8(10):e77110.

1170 Cornell, H. V. 1999. Unsaturation and regional influences on species richness in ecological
1171      communities: a review of the evidence. Ecoscience 6:303-315.

1172 Cornell, H. V., and S. P. Harrison. 2014. What are species pools and when are they important?
1173      Annual review of ecology, evolution, and systematics 45(1):45-67.

1174 Cutler, D. R., T. C. Edwards, K. H. Beard, A. Cutler, K. T. Hess, J. Gibson, and J. J. Lawler.
1175      2007. Random forests for classification in ecology. Ecology 88(11):2783-2792.

1176 Dalvé, E. 1948. The fossil fauna of the Ordovician in the Cincinnati region. University Museum,
1177      Department of Geology and Geography, University of Cincinnati.

1178 De'ath, G. 2002. Multivariate regression trees: a new technique for modeling species-
1179      environment relationships. Ecology 83(4):1105-1117.

1180 De'ath, G., and K. E. Fabricius. 2000. Classification and regression trees: a powerful yet simple
1181      technique for ecological data analysis. Ecology 81(11):3178-3192.

1182 de Bello, F. 2012. The quest for trait convergence and divergence in community assembly: are
1183      null-models the magic wand? Global Ecology and Biogeography 21(3):312-317.

1184 Deline, B. 2009. The effects of rarity and abundance distributions on measurements of local
1185      morphological disparity. Paleobiology 35(2):175-189.

1186   Deline, B., W. I. Ausich, and C. E. Brett. 2012. Comparing taxonomic and geographic scales in
1187         the morphologic disparity of Ordovician through Early Silurian Laurentian crinoids.
1188         Paleobiology 38(4):538-553.

1189   Dick, D. G., and E. E. Maxwell. 2015. The evolution and extinction of the ichthyosaurs from the
1190         perspective of quantitative ecospace modelling. Biology Letters 11(7).

1191   Dineen, A. A., M. L. Fraiser, and P. M. Sheehan. 2014. Quantifying functional diversity in pre-
1192         and post-extinction paleocommunities: A test of ecological restructuring after the end-
1193         Permian mass extinction. Earth-Science Reviews 136:339-349.

1194   Dineen, A. A., M. L. Fraiser, and J. Tong. 2015. Low functional evenness in a post-extinction
1195         Anisian (Middle Triassic) paleocommunity: A case study of the Leidapo Member
1196         (Qingyan Formation), south China. Global and Planetary Change 133:79-86.

1197   Dunne, J. A., R. J. Williams, N. D. Martinez, R. A. Wood, and D. H. Erwin. 2008. Compilation
1198         and network analyses of Cambrian food webs. PLoS Biol 6(4):e102.

1199   Durst, P. A. P., and V. L. Roth. 2012. Classification tree methods provide a multifactorial
1200         approach to predicting insular body size evolution in rodents. American Naturalist
1201         179(4):545-553.

1202   English, A. M., and L. E. Babcock. 2007. Feeding behaviour of two Ordovician trilobites
1203         inferred from trace fossils and non-biomineralised anatomy, Ohio and Kentucky, USA.
1204         Memoirs of the Association of Australasian Palaeontologists (34):537.

1205   Fargione, J., C. S. Brown, and D. Tilman. 2003. Community assembly and invasion: An
1206         experimental test of neutral versus niche processes. Proceedings of the National Academy
1207         of Sciences 100(15):8916-8920.

1208   Feldmann, R. M. 1996. Fossils of Ohio.

1209   Finnegan, S., S. C. Anderson, P. G. Harnik, C. Simpson, D. P. Tittensor, J. E. Byrnes, Z. V.
1210         Finkel, D. R. Lindberg, L. H. Liow, R. Lockwood, H. K. Lotze, C. R. McClain, J. L.
1211         McGuire, A. O'Dea, and J. M. Pandolfi. 2015. Paleontological baselines for evaluating
1212         extinction risk in the modern oceans. Science 348(6234):567-570.

1213   Finnegan, S., N. A. Heim, S. E. Peters, and W. W. Fischer. 2012. Climate change and the
1214         selective signature of the Late Ordovician mass extinction. Proceedings of the National
1215         Academy of Sciences.

1216   Fontana, S., O. L. Petchey, and F. Pomati. 2015. Individual-level trait diversity concepts and
1217           indices to comprehensively describe community change in multidimensional trait space.
1218           Functional Ecology 123(11):1391-1399.

1219   Foote, M. 1993. Discordance and concordance between morphological and taxonomic diversity.
1220           Paleobiology 19:185-204.

1221   ---. 1994. Morphological disparity in Ordovician-Devonian crinoids and the early saturation of
1222           morphological space. Paleobiology 20:320-344.

1223   Fordyce, D., and T. W. Cronin. 1993. Trilobite vision: a comparison of schizochroal and
1224           holochroal eyes with the compound eyes of modern arthropods. Paleobiology 19(3):288-
1225           303.

1226   Fortey, R. A., and R. M. Owens. 1999. The trilobite exoskeleton. Pp. 537-562. *In* E. Savazzi, ed.
1227           Functional Morphology of the Invertebrate Skeleton. John Wiley and Sons, Ltd., New
1228           York.

1229   Freeman, R. L., B. F. Dattilo, A. Morse, M. Blair, S. Felton, and J. Pojeta. 2013. The "curse of
1230           Rafinesquina": Negative taphonomic feedback exerted by strophomenid shells on storm-
1231           buried lingulids in the Cincinnatian Series (Katian, Ordovician) of Ohio. PALAIOS
1232           28(6):359-372.

1233   Frey, R. C. 1987a. The occurrence of pelecypods in Early Paleozoic epeiric-sea environments:
1234           Late Ordovician of the Cincinnati, Ohio area. PALAIOS 2:3-23.

1235   ---. 1987b. The paleoecology of a Late Ordovician shale unit from southwest Ohio and
1236           southeastern Indiana. Journal of Paleontology 61:242-267.

1237   ---. 1988. Paleoecology of *Treptoceras duseri* (Michelinoceratida, Proteoceratidae) from Late
1238           Ordovician of southwestern Ohio. New Mexico Bureau of Mines and Mineral Resources
1239           Memoir 44:79-101.

1240   ---. 1989. Paleoecology of a well-preserved nautiloid assemblage from a Late Ordovician shale
1241           unit, southwestern Ohio. Journal of Paleontology 63:604-620.

1242   Gaines, R. R., M. L. Droser, and N. C. Hughes. 1999. The ichnological record in Ordovician
1243           mudstones: examples from the Cincinnatian strata of Ohio and Kentucky (USA). Acta
1244           Universitatis Carolinae, Geologica (1/2):163-166.

1245     Gerhold, P., J. F. Cahill, M. Winter, I. V. Bartish, and A. Prinzing. 2015. Phylogenetic patterns
1246          are not proxies of community assembly mechanisms (they are far better). Functional
1247          Ecology 29(5):600-614.

1248     Gerisch, M. 2014. Non-random patterns of functional redundancy revealed in ground beetle
1249          communities facing an extreme flood event. Functional Ecology 28(6):1504-1512.

1250     Gower, J. C. 1971. A general coefficient of similarity and some of its properties. Biometrics
1251          27(4):857-871.

1252     Grueber, C. E., S. Nakagawa, R. J. Laws, and I. G. Jamieson. 2011. Multimodel inference in
1253          ecology and evolution: challenges and solutions. Journal of Evolutionary Biology
1254          24(4):699-711.

1255     Guillemot, N., M. Kulbicki, P. Chabanet, and L. Vigliola. 2011. Functional redundancy patterns
1256          reveal non-random assembly rules in a species-rich marine assemblage. PLoS One
1257          6(10):e26735.

1258     Hatton, I. A., K. S. McCann, J. M. Fryxell, T. J. Davies, M. Smerlak, A. R. E. Sinclair, and M.
1259          Loreau. 2015. The predator-prey power law: Biomass scaling across terrestrial and
1260          aquatic biomes. Science 349(6252).

1261     Hautmann, M. 2014. Diversification and diversity partitioning. Paleobiology:162-176.

1262     Hegna, T. A. 2010. The function of forks: *Isotelus*-type hypostomes and trilobite feeding.
1263          Lethaia 43(3):411-419.

1264     Heim, N. A., M. L. Knope, E. K. Schaal, S. C. Wang, and J. L. Payne. 2015. Cope's rule in the
1265          evolution of marine animals. Science 347(6224):867-870.

1266     Holland, S. M. 1993. Sequence stratigraphy of a carbonate-clastic ramp: the Cincinnatian Series
1267          (Upper Ordovician) in its type area. Geological Society of America Bulletin 105:306-
1268          322.

1269     ---. 2010. Additive diversity partitioning in palaeobiology: revisiting Sepkoski's question.
1270          Palaeontology 53(6):1237-1254.

1271     Holland, S. M., A. I. Miller, B. F. Dattillo, and D. L. Meyer. 2001. The detection and importance
1272          of subtle biofacies in lithologically uniform strata: the Upper Ordovician Kope Formation
1273          of the Cincinnati, Ohio region. PALAIOS 16:205-217.

Holland, S. M., and M. E. Patzkowsky. 2007. Gradient ecology of a biotic invasion: biofacies of the type Cincinnatian series (Upper Ordovician), Cincinnati, Ohio region, USA. PALAIOS 22(4):392-407.

---. 2009. The Richmondian invasion: understanding the faunal response to climate change through stratigraphic paleobiology. Type Cincinnatian (Upper Ordovician) outcrops, northern Kentucky, southwestern Ohio, and southeastern Indiana. Pp. 1-67. The Richmondian Invasion in the type Cincinnatian Series. Fieldtrip guidebook, Ninth North American Paleontological Convention. Cincinnati, OH.

Hubbell, S. P. 2001. The Unified Theory of Biodiversity and Biogeography. Princeton University Press, Princeton, NJ.

Hughes, N. C., and D. L. Cooper. 1999. Paleobiologic and taphonomic aspects of the "*Granulosa*" trilobite cluster, Kope Formation (Upper Ordovician, Cincinnati region). Journal of Paleontology 73(2):306-319.

Huntley, J. W., and D. Scarponi. 2012. Evolutionary and ecological implications of trematode parasitism of modern and fossil northern Adriatic bivalves. Paleobiology 38(1):40-51.

Johnson, J. B., and K. S. Omland. 2004. Model selection in ecology and evolution. Trends in Ecology & Evolution 19(2):101-108.

Kinzig, A. P., S. A. Levin, J. Dushoff, and S. Pacala. 1999. Limiting similarity, species packing, and system stability for hierarchical competition-colonization models. American Naturalist 153:371-383.

Knoll, A. H., and R. K. Bambach. 2000. Directionality in the history of life: diffusion from the left wall or repeated scaling of the right? Paleobiology (Supplement) 26(sp4):1-14.

Knope, M. L., N. A. Heim, L. O. Frishkoff, and J. L. Payne. 2015. Limited role of functional differentiation in early diversification of animals. Nature Communications 6.

Kowalewski, M., W. Kiessling, M. Aberhan, F. T. Fürsich, D. Scarponi, S. L. Barbour Wood, and A. P. Hoffmeister. 2006. Ecological, taxonomic, and taphonomic components of the post-Paleozoic increase in sample-level species diversity of marine benthos. Paleobiology 32(4):533-561.

Kowalewski, M., and P. Novack-Gottshall. 2010. Resampling methods in paleontology. Pp. 19-54. *In* J. Alroy, and G. Hunt, eds. Quantitative Methods in Paleobiology. Short Courses in

1304        Paleontology 16. Paleontological Society and Paleontological Research Institute, Ithaca,
1305        NY.

1306  Laflamme, M., S. A. F. Darroch, S. M. Tweedt, K. J. Peterson, and D. H. Erwin. 2013. The end
1307        of the Ediacara biota: Extinction, biotic replacement, or Cheshire Cat? Gondwana
1308        Research 23(2):558-573.

1309  Laflamme, M., S. Xiao, and M. Kowalewski. 2009. Osmotrophy in modular Ediacara organisms.
1310        Proceedings of the National Academy of Sciences (U.S.A.) 106(34):14438-14443.

1311  Laliberté, E., and P. Legendre. 2010. A distance-based framework for measuring functional
1312        diversity from multiple traits. Ecology 91(1):299-305.

1313  Laliberté, E., and B. Shipley. 2014. FD: Measuring functional diversity from multiple traits, and
1314        other tools for functional ecology, Version 1.0-12.

1315  Legendre, P., and M. J. Anderson. 1999. Distance-based redundancy analysis: testing
1316        multispecies responses in multifactorial ecological experiments. Ecological Monographs
1317        69(1):1-24.

1318  Leighton, L. R. 1998. Constraining functional hypotheses: controls on the morphology of the
1319        concavo-convex brachiopod *Rafinesquina*. Lethaia 31:293-307.

1320  Lescinsky, H. L. 1995. The life orientation of concavo-convex brachiopods: overturning the
1321        paradigm. Paleobiology 21:520-551.

1322  Maire, E., G. Grenouillet, S. Brosse, and S. Villéger. 2015. How many dimensions are needed to
1323        accurately assess functional diversity? A pragmatic approach for assessing the quality of
1324        functional spaces. Global Ecology and Biogeography 24(6):728-740.

1325  Mason, N. W. H., D. Mouillot, W. G. Lee, and J. B. Wilson. 2005. Functional richness,
1326        functional evenness and functional divergence: the primary components of functional
1327        diversity. Oikos 111(1):112-118.

1328  McShea, D. W. 1994. Mechanisms of large-scale evolutionary trends. Evolution 48(6):1747-
1329        1763.

1330  Meyer, D. L., A. I. Miller, S. M. Holland, and B. F. Dattilo. 2002. Crinoid distribution and
1331        feeding morphology through a depositional sequence: Kope and Fairview Formations,
1332        Upper Ordovician, Cincinnati Arch region. Journal of Paleontology 76(4):725-732.

1333  Miller, J. H., A. K. Behrensmeyer, A. Du, S. K. Lyons, D. Patterson, A. Tóth, A. Villaseñor, E.
1334        Kanga, and D. Reed. 2014. Ecological fidelity of functional traits based on species

presence-absence in a modern mammalian bone assemblage (Amboseli, Kenya). Paleobiology 40(4):560-583.

Mitchell, J. S., and P. J. Makovicky. 2014. Low ecological disparity in Early Cretaceous birds. Proceedings of the Royal Society B: Biological Sciences 281(1787).

Mitchell, J. S., P. D. Roopnarine, and K. D. Angielczyk. 2012. Late Cretaceous restructuring of terrestrial communities facilitated the end-Cretaceous mass extinction in North America. Proceedings of the National Academy of Sciences 109(46):18857-18861.

Mondal, S., and P. J. Harries. 2015. Phanerozoic trends in ecospace utilization: The bivalve perspective. Earth-Science Reviews.

Morris, R. W., and S. H. Felton. 2003. Paleoecologic associations and secondary tiering of *Cornulites* on crinoids and bivalves in the Upper Ordovician (Cincinnatian) of southwestern Ohio, southeastern Indiana, and northern Kentucky. PALAIOS 18(6):546-558.

Mouchet, M. A., S. Villéger, N. W. H. Mason, and D. Mouillot. 2010. Functional diversity measures: an overview of their redundancy and their ability to discriminate community assembly rules. Functional Ecology 24(4):867-876.

Mouillot, D., N. A. J. Graham, S. Villéger, N. W. H. Mason, and D. R. Bellwood. 2013. A functional approach reveals community responses to disturbances. Trends in Ecology and Evolution 28(3):167-177.

Novack-Gottshall, P. 2015. ecospace: Simulating Community Assembly and Ecological Diversification Using Ecospace Frameworks, Version 1.0.1. cran.r-project.org/package=ecospace.

Novack-Gottshall, P. M. 2007. Using a theoretical ecospace to quantify the ecological diversity of Paleozoic and modern marine biotas. Paleobiology 33(2):273-294.

---. 2008a. Ecosystem-wide body size trends in Cambrian-Devonian marine invertebrate lineages. Paleobiology 34(2).

---. 2008b. Using simple body-size metrics to estimate fossil body volume: empirical validation using diverse Paleozoic invertebrates. PALAIOS 23(3):163-173.

---. 2010. Performance of functional diversity metrics applied as measures of disparity. GSA Abstracts with Programs 42(5):140.

---. 2016. General models of ecological diversification. I. Conceptual synthesis. Paleobiology.

1366    Novack-Gottshall, P. M., and A. I. Miller. 2003. Comparative taxonomic richness and abundance
1367         of Late Ordovician gastropods and bivalves in mollusc-rich strata of the Cincinnati Arch.
1368         PALAIOS 18(6):559-571.

1369    O'Brien, L. J., and J.-B. Caron. 2015. Paleocommunity analysis of the Burgess Shale Tulip Beds,
1370         Mount Stephen, British Columbia: comparison with the Walcott Quarry and implications
1371         for community variation in the Burgess Shale. Paleobiology FirstView:1-27.

1372    Pakeman, R. J. 2014. Functional trait metrics are sensitive to the completeness of the species'
1373         trait data? Methods in Ecology and Evolution 5(1):9-15.

1374    Patzkowsky, M. E., and S. M. Holland. 2003. Lack of community saturation at the beginning of
1375         the Paleozoic plateau: the dominance of regional over local processes. Paleobiology
1376         29:545-560.

1377    ---. 2007. Diversity partitioning of a Late Ordovician marine biotic invasion: controls on
1378         diversity in regional ecosystems. Paleobiology 33(2):295-309.

1379    Pojeta, J., Jr. 1971. Review of Ordovician pelecypods. United States Geological Survey
1380         Professional Paper 695.

1381    Powell, M. G., and M. Kowalewski. 2002. Increase in evenness and sampled alpha diversity
1382         through the Phanerozoic: comparison of early Paleozoic and Cenozoic marine fossil
1383         assemblages. Geology 30(4):331-334.

1384    Prasad, A., L. Iverson, and A. Liaw. 2006. Newer classification and regression tree techniques:
1385         bagging and random forests for ecological prediction. Ecosystems 9(2):181-199.

1386    R Development Core Team. 2015. R: A language and environment for statistical computing,
1387         Version 3.2.0. R Foundation for Statistical Computing, Vienna, Austria.

1388    Raup, D. M. 1966. Geometric analysis of shell coiling: general problems. Journal of
1389         Paleontology 40:1178-1190.

1390    Richards, R. P. 1972. Autecology of Richmondian brachiopods (Late Ordovician of Indiana and
1391         Ohio). Journal of Paleontology 46:386-405.

1392    Ros, S., M. D. Renzi, S. E. Damborenea, and A. Marquez-Aliaga. 2012. Early Triassic–Early
1393         Jurassic bivalve diversity dynamics Pp. 1-19. Treatise Online. University of Kansas,
1394         Lawrence, KS.

Rudkin, D. M., G. A. Young, R. J. Elias, and E. P. Dobrzanski. 2003. The world's biggest trilobite—*Isotelus rex* new species from the Upper Ordovician of Northern Manitoba, Canada. Journal of Paleontology 77(1):99-112.

Sandy, M. R. 1996. Oldest record of peduncular attachment of brachiopods to crinoid stems, Upper Ordovician, Ohio, U.S.A. Journal of Paleontology 70:532-534.

Schumacher, G. A., and D. L. Shrake. 1997. Paleoecology and comparative taphonomy of an *Isotelus* (Trilobita) fossil lagerstätten from the Waynesville Formation (Upper Ordovician, Cincinnatian Series) of southwestern Ohio. Pp. 131-161. *In* C. E. Brett, and G. C. Baird, eds. Paleontological Events: Stratigraphic, Ecological, and Evolutionary Implications. Columbia University Press, New York.

Seilacher, A. 1970. Arbeitskonzept zur Konstruktions-Morphologie. Lethaia 3:393-396.

Sepkoski, J. J., Jr. 1981. A factor analytic description of the Phanerozoic marine fossil record. Paleobiology 7(1):36-53.

Slater, G. J., L. J. Harmon, D. Wegmann, P. Joyce, L. J. Revell, and M. E. Alfaro. 2012. Fitting models of continuous trait evolution to incompletely sampled comparative data using approximate Bayesian computation. Evolution 66(3):752-762.

Strobl, C., A.-L. Boulesteix, A. Zeileis, and T. Hothorn. 2007. Bias in random forest variable importance measures: Illustrations, sources and a solution. BMC Bioinformatics 8(1):25.

Sullivan, M., M. Jones, D. Lee, S. Marsden, A. Fielding, and E. Young. 2006. A comparison of predictive methods in extinction risk studies: contrasts and decision trees. Biodiversity & Conservation 15(6):1977-1991.

Thomas, R. D. K., and W. E. Reif. 1993. The skeleton space: a finite set of organic designs. Evolution 47:341-360.

Tomašových, A., and S. M. Kidwell. 2010. Predicting the effects of increasing temporal scale on species composition, diversity, and rank-abundance distributions. Paleobiology 36(4):672-695.

Van Valen, L. 1974. Multivariate structural statistics in natural history. Journal of Theoretical Biology 45:235-247.

Van Valkenburgh, B. 1988. Trophic diversity in past and present guilds of large predatory mammals. Paleobiology 14:155-173.

1425 Van Valkenburgh, B., and R. E. Molnar. 2002. Dinosaurian and mammalian predators compared.
1426    Paleobiology 28:527-543.

1427 Vermeij, G. J. 1987. Evolution and Escalation: An Ecological History of Life. Princeton
1428    University Press.

1429 ---. 2011. The energetics of modernization: The last one hundred million years of biotic
1430    evolution. Paleontological Research 15(2):54-61.

1431 ---. 2013. On escalation. Annual Review of Earth and Planetary Sciences 41(1):1-19.

1432 Villéger, S., N. W. H. Mason, and D. Mouillot. 2008. New multidimensional functional diversity
1433    indices for a multifaceted framework in functional ecology. Ecology 89(8):2290-2301.

1434 Villéger, S., J. R. Miranda, D. F. Hernández, and D. Mouillot. 2010. Contrasting changes in
1435    taxonomic vs. functional diversity of tropical fish communities after habitat degradation.
1436    Ecological Applications 20(6):1512-1522.

1437 Villéger, S., P. M. Novack-Gottshall, and D. Mouillot. 2011. The multidimensionality of the
1438    niche reveals functional diversity changes in benthic marine biotas across geological
1439    time. Ecology Letters 14(6):561-568.

1440 Villier, L., and G. J. Eble. 2009. Assessing the robustness of disparity estimates: the impact of
1441    morphometric scheme, temporal scale, and taxonomic level in spatangoid echinoids.

1442 Vogt, R. J., P. R. Peres-Neto, and B. E. Beisner. 2013. Using functional traits to investigate the
1443    determinants of crustacean zooplankton community structure. Oikos 122(12):1700-1709.

1444 Wills, M. A. 2001. Morphological disparity: a primer. Pp. 55-143. *In* J. M. Adrain, G. D.
1445    Edgecombe, and B. S. Lieberman, eds. Fossils, phylogeny, and form: an analytical
1446    approach. Kluwer Academic/Plenum Publishers, New York.

1447 Xiao, S., and M. Laflamme. 2009. On the eve of animal radiation: phylogeny, ecology and
1448    evolution of the Ediacara biota. Trends in Ecology & Evolution 24(1):31-40.

1449 Zhang, Z., M. Augustin, and J. L. Payne. 2015. Phanerozoic trends in brachiopod body size from
1450    synoptic data. Paleobiology 41(03):491-501.

1451

Figure 1. Typical examples of simulated 50-species assemblages produced by the four model

rules. Assemblages are plotted on common non-metric multi-dimensional scaling ordination

space of functional traits to allow comparative evaluation of model behavior. Five gray diamonds

represent common "seed" species whose life habits were assigned stochastically using an 18-

character (functional trait space) ecospace framework (modified from Novack-Gottshall 2007),

imposing a realistic constraint that each life habit could have at most two character states within

a given character. Numbers illustrate the addition of five species to each assemblage (after seed

species), with remaining 40 species as hollow circles. All model rules, except redundancy, were

1461    enacted at 100% rule-following for each simulation; redundancy rules were weakened such that

1462    all successive species have habits 95% similar to pre-existing ones; at 100% enactment, later life

1463    habits are limited to the "seed" species. In the neutral model, functional traits of all species are

1464    chosen independently at random, and the entire ecospace becomes inhabited through passive

1465    processes. In the redundancy model, new species have life habits similar to pre-existing ones,

1466    producing an ecospace with distinct clusters. In the partitioning model, new species inhabit life

1467    habits intermediate to pre-existing ones. This model can be enacted in a strict form (larger

1468    image) in which new species are restricted to gradients between preexisting species (typically

1469    leaving the center empty) and a relaxed form (inset) in which new species progressively fill in

1470    empty regions of the space originally defined by the seed species. In the expansion model, new

1471    species progressively inhabit novel life habits, producing an ecospace that expands its breadth

1472    over the simulation, while leaving the original region uninhabited. Figure 2B shows the average

1473    dynamics of eight structural statistics when such simulations were repeated 1000 times.

1474

1475

Figure 2. Model dynamics for eight functional-diversity statistics for the ecospace frameworks of (A) Bush and Bambach (Bambach et al. 2007, Bush et al. 2007, Bush and Bambach 2011, Bush et al. 2011) and (B) Novack-Gottshall (2007). In both cases, all models are enacted at 100% rule-following with five seed species. The second case (part B) also stipulates that each life habit can have at most two character states within a given life-habit character (see text for explanation of

1482    'constraint' parameter). Dynamics are reported as a function of increasing species richness, up to

1483    50 species (i.e., there is a common abscissa in all graphs); the ordinate has been truncated to

1484    focus on the trend lines for each statistic. Statistics include life-habit richness (H), mean distance

1485    (D), maximum range (M), total variance (V), functional richness (FRic), evenness (FEve),

1486    divergence (FDiv), and dispersion (FDis); see text for description of each statistic. Trend lines

1487    are the average of 1000 simulations. Vertical bars near top right of each graph illustrate the

1488    average standard deviation around each set of trends; the standard deviation for mean distance in

1489    the Bush and Bambach framework equals 0.0117 and extends beyond the borders of the graph.

1490    PartS and PartR represent the strict and relaxed versions of the partitioning model. The variance

1491    trend in part A is omitted because the characters in the Bush and Bambach framework are all

1492    factorial, preventing its calculation. The uppermost H trend line in part B is an overlap between

1493    the neutral, relaxed partitioning, and expansion trend lines. Despite their rather different

1494    structures, the model dynamics for each ecospace framework are quite similar. Notable

1495    differences include total overlap between neutral and expansion models in the Bush and

1496    Bambach framework, and generally smaller error margins for the Novack-Gottshall framework,

1497    which allows more powerful model selection using classification-tree methods (89% of

1498    validation models classified correctly versus 73% with the Bush and Bambach framework).

1499

1500

Figure 3. Comparison of Late Ordovician samples to dynamics of different simulation

implementations: (A) 100% rule-following implementation with one-state constraint, (B) 100%

rule-following implementation with two-state constraint and empirical weighting of states, (C)

100%, 90%, and 50% rule-following with two-state constraint and empirical weighting. The

legend and graphical interpretation is the same as Figure 2, and variability around the mean trend

lines are of similar magnitude. Statistics for Late Ordovician (type Cincinnatian) Kope and

1508   Waynesville Formation samples are represented by circles and triangles, respectively. In part A,

1509   simulations were enacted at 100% rule-following and each life habit was constrained to have just

1510   one character state within a given life-habit character; for example, a taxon could be either sexual

1511   or asexual, but not both (hermaphroditic). The dynamics are similar to those for the two-state

1512   implementation in Figure 2B, although values are slightly diminished, especially for the

1513   partitioning models. In all cases, the empirical sample statistics lay well below the mean trend

1514   lines, implying a poor fit between these simulations and reality. In part B, the simulations

1515   included the two-state constraint, but were modified so that the combined Kope and Waynesville

1516   species pool was used to weight how seed species were chosen (i.e., they were chosen at random

1517   from the species pool) and the inhabitation of character states of all successive species were

1518   weighted by their frequency of occurrence in the species pool. The fit between the empirical

1519   samples and the model dynamics is much improved (that is, there is substantially more overlap).

1520   In part C, eight additional "weaker" simulations are added, in which the model rules are followed

1521   on average 90% and 50% of the time. The average trend lines for these weaker simulations are

1522   slightly thinner (but with the same colors and line types) than the 100% implementations to make

1523   it easier to distinguish them (although the 90% expansion model trend mostly overlays the 100%

1524   implementation). A 0% simulation is always present, the neutral model, in which no model rules

1525   are followed. The empirical samples frequently overlay weaker versions of the models, and all

1526   models coalesce toward the neutral model as they become weaker and weaker.

1527

Figure 4. Relative model support for Kope (left column) and Waynesville (right column) samples at scales of individual samples, aggregated stratigraphic sections, members, and entire formation species pool. Model support for samples and sections was calculated using a classification tree trained on the 100% and 50% simulations used in Figure 3C; model support for members and formations used the classification tree trained on 50%, 75%, 90%, 95%, and 100% simulations, which is a more powerful classifier at larger sample sizes. Models are only listed (middle column) when they have support in a particular scale. Support for all but the formation-scale is

1536     the number of samples/sections/members that were classified for each model. Support for the

1537     formations (bottom row) is the proportion of votes for each model. The aggregate Kope

1538     Formation is overwhelmingly (0.52) classified as the 90% strict partitioning model, although

1539     there is also substantial support (0.34) for the 90% redundancy model. The aggregate

1540     Waynesville Formation has overwhelming support (0.78) for the 90% redundancy model, with

1541     less but substantial support (0.22) for the 95% strict partitioning model.

1542

1543
1544 Figure 5. Ordinations of four representative Kope and Waynesville samples best fit by various

1545 models, and the corresponding ordination of each formation species pool. Ordinations were

1546 conducted as two-dimensional non-metric multidimensional scaling on a common

1547 Kope/Waynesville species pool to allow comparison across graphs. Text above each graph notes

1548 which model was best fit (i.e., had the most votes in the classification tree) for each sample, as

1549   well as the relative support for that "vote-winning" model. Also listed are the Paleobiology

1550   Database sample collection number, its taxonomic richness and number of unique life habits (H),

1551   and symbols representing major taxonomic groups. Note that jittering (moving overlapping

1552   points by tiny amounts) was not used in this figure; points that overlap incompletely represent

1553   distinct life habits. Aside from at the formation-scale, there is relatively little absolute

1554   redundancy among life habits.

1555

1556

Supplementary Figure 6. Comparing statistical dynamics for different ecospace framework

structures: varying number of characters, (A) 5 characters, (B) 15 characters, and (C) 25

characters. Each framework had mixed character types, in identical proportions (40% binary,

20% three-state factor, 20% five-state factor, and, 20% five-state ordered numeric character

types). 5 "seed" species were chosen at random to begin each simulation. Other simulation

details and graphical interpretation are the same as is Figure 2. Trends in total variance were

excluded because the inclusion of factors prevented their calculation. The dynamics are generally

similar, although larger frameworks allow modestly more powerful model selection using

classification-tree methods (83%, 85%, and 86% of training models, respectively, classified

correctly using classification-tree methods). See Supplementary Appendix 2 for additional

details.

1568

1569
1570 Supplementary Figure 7. Comparing statistical dynamics for different ecospace framework

1571 structures: varying character types, (A) factor, (B) ordered factor, (C) ordered numeric, and (D)

1572 binary. Each framework had 15 characters, four states per character (except for binary, which

1573 had two binary states per character), and five seed species. Trends in total variance were

1574 excluded in parts A and B because the inclusion of factors prevented their calculation. Other

1575 simulation details and graphical interpretation are the same as in Supplementary Figure 6.

1576 Dynamics are generally similar, but frameworks built with ordered factors performed

1577 substantially better (94% of trained models classified correctly) than the others (78% for

1578 unordered factors, 79% for ordered numerics, and 81% for binaries). See Supplementary

1579 Appendix 2 for additional details.

1580

1581
1582 Supplementary Figure 8. Comparing statistical dynamics for different ecospace framework

1583 structures: varying number of "seed" species chosen at random at start of simulation, (A) 3

1584 species, (B) 5 species, and (C) 10 species. Each framework had 15 mixed character types, such

1585 that part B is identical to Supplementary Figure 6B. Trends are only plotted starting at 5 species

1586 for comparative purposes, which explains idiosyncratic behaviors at low sample sizes (i.e.,

1587 missing trend lines in part A and overlapping models in part C). Other simulation details and

1588 graphical interpretation are the same as is Supplementary Figure 6. Note that the functional-

1589 diversity statistics could not be calculated for the redundancy model in part A because their

1590    calculation requires a minimum of four unique life habits; however, all statistics (except V) were

1591    included as potential predictor variables in the classification tree algorithm. Dynamics are

1592    generally similar across models, but simulations with fewer seed species provide the most

1593    powerful model selection using classification-tree methods (85% of models classified correctly

1594    for both 3-seed and 5-seed simulations). Starting with larger numbers of seed species impedes

1595    enacting distinct model rules, and results in 75% of models classified correctly. See

1596    Supplementary Appendix 2 for additional details.

1597

1599    Supplementary Figure 9. Performance of classification tree used on validation samples as a

1600    function of sample size. This is the tree used to classify empirical samples in analyses for Figures

1601    4 and 5. The classification tree was trained on 5,265,500 simulated samples spanning nine

1602    models (identified in figure legend, with mean trends for 100% and 50% training data sets

1603    visualized in Figure 3C): the neutral model; the four models that implement the redundancy, two

1604    versions of partitioning, and expansion rules; and four "weakened" versions of these four

1605    models, in which rules were followed at random 50% of the time. Numbers in model names are

1606    abbreviated, with "-1" referring to the 100%-rule-following implementation and "-0.50" referring

1607    to the 50% implementation. Performance of the tree was evaluated against 2,375,000 validation

1608    samples (i.e., new samples not used when training the tree), that included simulated data not only

1609    from the 50% and 100% rule-following simulations, but also those from 75%, 90%, and 95%

1610    implementations. In other words, there were 1000 simulations per model-sample size * 95

1611    sample sizes [5–100 species per sample] * 25 models [with independent samples from the neutral

1612    model included in each set]).

1613    Each row corresponds to a validation data set, with 100% at the top, 95% below that, continuing

1614    to 50% at the bottom. Each column corresponds to a model, with expansion at left, followed

1615    rightward by neutral, relaxed implementation partitioning, strict partitioning, and redundancy.

1616    Trend lines illustrate the proportion of votes for each model in the classification tree, as a

1617    function of sample size. For example, the top-left graph (illustrating how "known" 100%-

1618    expansion validation samples were classified), shows that nearly all such validation samples

1619    were correctly classified as simulated from the 100%-expansion model, even at low sample

1620    sizes, and with essentially all such samples classified perfectly at sample sizes as low as 15. That

1621    nearly 40% of samples with just 6 species are correctly classified is especially promising because

1622   this is the smallest sample size when expansion rules were implemented (recall that that five

1623   species were assigned at random to seed the simulation). Similar performance occurs for most of

1624   the 100% implementations. Classification performance for the 50% implementations is less

1625   accurate, although the correct model always receives majority support, nearly always at or above

1626   support values of 0.40.

1627   Patterns of support for the intermediate (and novel) validation samples, those produced by

1628   implementing model rules 75%, 90%, and 95% of the time, demonstrate that the use of

1629   classification trees produces powerful and intuitively understandable performance as a model-

1630   selection method. In general, model support at the 95% level mimics closely that at the 100%

1631   level, with somewhat weaker (but still majority) support for the more similar model. This makes

1632   intuitive sense; the 95% samples are produced by the 100% expansion model 95% of the

1633   simulation, on average. As one looks down a column, relative support transitions across

1634   classified models as expected: support for the 100% model diminishes and gives way to support

1635   for the 50% model. The switch from one model to another occurs generally at the half-way point

1636   (75%), and this point is also where most confusion in model classification occurs. This pattern

1637   further suggests that support for different weights of the same model (e.g., the 75% expansion

1638   samples are classified as both 100% and 50% expansion models), can provide evidence for a

1639   weakened model (i.e., intermediate between the 100% and 50% implementations), especially for

1640   smaller sample sizes (below 30–40 taxa). This mutual support for alternative models is lessened

1641   at larger sample sizes, where the simulations used to train the classification tree generally have

1642   reduced variability around the asymptotes at these sample sizes, and which leads to increased

1643   classification confusion as the tree is more powerful at rejecting the trained implementations.

1644   Trends down the "neutral" column are largely identical, reflecting the fact that the neutral model

1645    is always the absence of any assembly rule (either always 100% neutral or 0% an alternative

1646    model), and it is reassuring that the independent "neutral" simulations are consistently classified

1647    across iterations.

1648    Other patterns are worth noting when evaluating classification trees as a method of model

1649    selection. Across sample sizes, the correct model consistently achieves minimum support levels

1650    greater than 0.2, from which it is possible to conclude that this support value be recommended as

1651    evidence of "substantial" support for the correct model, and values above 0.40 tend to be

1652    associated with "unambiguously strong" support for the correct model. These recommended

1653    support levels are also found for more complicated trees (i.e., those trained on 50%, 75%, 90%,

1654    95%, and 100%-strength implementations), especially for sample sizes greater than 40 species.

1655    (It is worth noting that these "critical" 0.2 and 0.4 support values are the result of validation tests

1656    of these models, and they may not be generalizable to other implementations; however, their

1657    consistency when evaluated across 9-, 13-, and 21-training set implementations is promising.)

1658    Model support for simulations created with weakened (<90% implementation) models tend to be

1659    confused with other models (i.e., they have greater numbers of misclassifications). Weakened

1660    implementations of the partitioning and redundancy models (redundancy-0.5, partitioningR-0.5,

1661    and partitioningS-0.5), in particular, get confused for one another; see Supplementary Appendix

1662    2 for additional confirmation of this claim. This is not a weakness of the classification tree but a

1663    confirmation of its power; in their weakened form, these three models are truly dynamically

1664    similar. The classification tree is still sufficiently powerful to distinguish these models at least

1665    40% of the time, even at small sample sizes.

1666 **Supplementary Appendix 1. Example of how life-habit character states were inferred and**

1667 **coded. See Supplementary Table 1 for raw character-state codings.**

1668 Mature asaphid trilobite *Isotelus maximus* specimens have a shell volume of ca. 270 cm3, using caliper

1669 measurements of major axes (Feldmann 1996) and volumetric allometries from Novack-Gottshall

1670 (2008b). It is coded as a carnivorous predator given its large, forked hypostome with a rasp-like texture,

1671 likely supplemented with limbs to dispatch prey (Fortey and Owens 1999, Hegna 2010). *Rusophycus*

1672 trace fossils further support this foraging habit, suggesting evidence of raptorial attacks on *Palaeophycus*-

1673 producing shallow-burrowing worms using large limbs (Brandt et al. 1995, Rudkin et al. 2003, English

1674 and Babcock 2007). Raised holochroal eyes well suited to deep water (Fordyce and Cronin 1993) and

1675 absence of other trace fossils suggests a hunting strategy where the trilobite swam habitually above the

1676 seafloor searching for shallow-infaunal prey at some distance. Swimming height is less well constrained,

1677 but was most likely nektobenthic given flattened exoskeleton profile and previously discussed trace fossil

1678 evidence. Given there are no asexual arthropods, it likely reproduced gonochoristically.

1679 **Supplementary Table 1.** 18 characters are in bold, with 37 character states numbered. First six characters

1680 are coded as ordered numeric factors and remainder are coded as binary states. (Such binary codings are

1681 used instead of discrete factors to allow for logically possible state combinations. For example, a

1682 hermaphrodite would be coded as reproductive code [1,1] for being both asexual and sexual, and eating

1683 food from on top of and within the seafloor would be coded as primary feeding microhabitat code [1,1].

1684 See Novack-Gottshall (2007) for definitions and discussion of characters and states.

1) **Skeletal body volume**: 1 (100–1000 cm3)
2) **Mobility**: 1 (Habitually mobile)
3) **Primary stratification**: 0.5 (1-10 cm from seafloor)
4) **Immediate stratification**: 0.5 (1-10 cm from seafloor)
5) **Primary food stratification**: 0.25 (0.1-1.0 cm from seafloor)
6) **Immediate food stratification**: 0.75 (10-100 cm from seafloor)

**Reproduction**:
7. Asexual: 0 (No)
8. Sexual: 1 (Yes)

**Substrate/medium composition:**
9. Biotic: 0 (No)
10. Lithic: 1 (Yes)
11. Fluidic: 0 (No)

**Substrate consistency:**
12. Hard: 0 (No)
13. Soft: 1 (Yes)
14. Insubstantial: 0 (No)

**Substrate relationship:**
15. Attached: 0 (No)
16. Free-living: 1 (Yes)

**Primary microhabitat:**
17. Above primary substrate: 1 (Yes)
18. Within primary substrate: 0 (No)

**Immediate microhabitat:**
19. Above immediate substrate: 1 (Yes)
20. Within immediate substrate: 0 (No)

**Support:**
21. Supported: 0 (No)
22. Self-supported: 1 (Yes)

**Primary feeding microhabitat:**
23. Above primary substrate: 1 (Yes)
24. Within primary substrate: 1 (Yes)

**Immediate feeding microhabitat:**
25. Above immediate substrate: 1 (Yes)
26. Within immediate substrate: 1 (Yes)

**Diet:**
27. Autotroph: 0 (No)
28. Microbivore: 0 (No)
29. Carnivore: 1 (Yes)

**Physical condition of food:**
30. Incorporeal feeder: 0 (No)
31. Particle feeder: 0 (No)
32. Bulk feeder: 1 (Yes)

**Feeding habit:**
33. Ambient feeder: 0 (No)
34. Filter feeder: 0 (No)
35. Attachment feeder: 0 (No)
36. Mass feeder: 0 (No)
37. Raptorial feeder: 1 (Yes)

1685

1686 **Supplementary Appendix 2. Technical details on classification tree methods, and**

1687 **confusion-matrix results compared across various simulations.**

1688 Classification trees are a powerful and flexible tool for classifying complex datasets that have

1689 complex, often non-linear and localized relationships among variables (Breiman et al. 1984,

1690 Cutler et al. 2007) and are being increasingly used by ecologists (De'ath and Fabricius 2000,

1691 De'ath 2002, Sullivan et al. 2006, Cutler et al. 2007, Boyer 2010, Durst and Roth 2012) and

1692 paleontologists (Finnegan et al. 2012, Finnegan et al. 2015). Except where noted in text,

1693 classification-tree analyses were run in the following manner: classification trees were trained on

1694 1000-replicate data sets (per model, each implemented 100% rule-following) with sample sizes

1695 of 6 to 50 species. (In other words, the training set for each tree included 225,000 samples, equal

1696 to 1000 simulations * 45 sample sizes [6–50 species per sample], or 45,000 simulated samples

1697 for each model * 5 models.) Random-forest classification was conducted in the R randomForest

1698 package (Breiman 2006), building ensembles of 50 bootstrapped replicate trees ('ntrees=50')—in

1699 which some samples and variables are excluded at random during tree building to enhance

1700 statistical power and reduce overfitting—to terminal node sizes of 1 ('nodesize=1'). These

1701 algorithmic parameters were justified by sensitivity analyses.

1702 The basic structure was to predict (classify) each model as a function of species richness and

1703 eight functional diversity/disparity statistics:

1704 $$Model \sim S + H + D + M + V + FRic + FEve + FDiv + FDis,$$

1705 with variables defined in text. In the following confusion matrices, the total variance (V) statistic

1706 was excluded for consistency, as this statistic could not be calculated for ecospace -framework

1707 simulations with factors present. (Although the classification tree considers V a relatively

1708 important statistic for model selection, its high correlation with equally important mean distance

1709    (D) yields relatively little effect, less than 1%, when excluding it.) The four functional diversity

1710    statistics were included in all analyses, even in the one case (3 seed species, Supplementary Fig.

1711    8) where they could not be calculated (excluding them had a negligible effect on model

1712    classification).

1713    Classification rates and confusion matrices provide measures of the ability of classification trees

1714    to classify model trends. Validation tests were not run for most of the following trials, so

1715    classification rates may be elevated slightly (<10%) because of potential model overfitting, but

1716    relative rankings for comparative purposes should remain accurate. Validation tests (with data

1717    sets of the same size) were run on several simulations (e.g., those in Fig. 2 and Fig. 3C), with

1718    validation classification rates reported in brackets.

1719    The first column in the confusion matrices is the known ("true") model from which samples were

1720    simulated. The next five columns show the models the tree classified these samples. In an ideal

1721    classification tree (i.e., in which all models are predicted perfectly without error), for example,

1722    all "Expansion-1" samples (those samples simulated following the expansion model 100% of the

1723    time) would be classified into this same model. In other words, the diagonals in the confusion

1724    matrix show correctly classified models. The final column ("class.error") tallies the proportion of

1725    each known model that was incorrectly classified.

1726 **Comparing number of characters in ecospace frameworks (Supplementary Fig. 6 and**
1727 **Supplementary Fig. 7B)**

1728 **5 mixed characters: 83.38% correct**

```
1729  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1730  Expansion-1        29920    13958        751            348            23         0.3351
1731  Neutral            11626    32145        804            402            23         0.2857
1732  PartitioningR-1      767      899      40073           3147           114         0.1095
1733  PartitioningS-1      337      454       3537          40493           176         0.1001
1734  Redundancy-1           1        1          7              9         44982         0.0004
```

1735 **15 mixed characters (and 5 seed species): 84.65% correct**

```
1736  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1737  Expansion-1        28839    15867        140            154             0         0.3591
1738  Neutral            14364    30344        163            129             0         0.3257
1739  PartitioningR-1      117      128      43404           1351             0         0.0355
1740  PartitioningS-1       70       52       2001          42877             0         0.0472
1741  Redundancy-1           0        0          0              0         45000         0.0000
```

1742 **25 mixed characters: 85.94% correct**

```
1743  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1744  Expansion-1        29658    15265         43             34             0         0.3409
1745  Neutral            13955    30985         32             28             0         0.3114
1746  PartitioningR-1       23       18      44045            913             1         0.0212
1747  PartitioningS-1       10       15       1295          43678             2         0.0294
1748  Redundancy-1           0        0          0              0         45000         0.0000
```

1749 **Comparing number of "seed" species used at start of simulations (Supplementary Fig. 8)**

1750 **3 seed species: 84.53% correct**

```
1751  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1752  Expansion-1        29230    15744          6             20             0         0.3504
1753  Neutral            14069    30911          7             13             0         0.3131
1754  PartitioningR-1        4       10      41158           3815            13         0.0854
1755  PartitioningS-1        7       16       1070          43903             4         0.0244
1756  Redundancy-1           0        0          0              7         44993         0.0002
```

1757 **10 seed species: 74.68% correct**

```
1758  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1759  Expansion-1        26956    14223       1487           1277          1057         0.4010
1760  Neutral            14034    27106       1520           1299          1041         0.3976
1761  PartitioningR-1     1415     1442      37161           3882          1100         0.1742
1762  PartitioningS-1     1152     1269       5414          36075          1090         0.1983
1763  Redundancy-1        1022     1065       1069           1109         40735         0.0948
```

1764 **Comparing character types used in ecospace framework (Supplementary Fig. 7)**

1765 **Binary character types: 81.08% correct**

```
1766  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1767  Expansion-1        26723    18042        187             48             0         0.4062
1768  Neutral            17896    26847      ,  198             59             0         0.4034
1769  PartitioningR-1      167      137      42086           2610             0         0.0648
1770  PartitioningS-1       29       23       3172          41774             2         0.0717
1771  Redundancy-1           0        0          0              0         45000         0.0000
```

1772 **Factor character types: 78.15% correct**

```
1773  True model      Expansion-1 Neutral  PartitioningR-1 PartitioningS-1 Redundancy-1  class.error
1774  Expansion-1        24391  , 20520         72             17             0         0.4580
1775  Neutral          , 20608    24318         56             18             0         0.4596
1776  PartitioningR-1       35       24      41589           3351             1         0.0758
1777  PartitioningS-1        2       10       4435          40551             2         0.0989
1778  Redundancy-1           0        0          0              1         44999         0.0000
```

#### 1779 **Ordered factor character types: 94.06% correct**

```
1780  True model   Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1781  Expansion-1        41526     2816              306              352             0       0.0772
1782  Neutral             2501    41301              718              480             0       0.0822
1783  PartitioningR-1      357      911            42126             1606             0       0.0639
1784  PartitioningS-1      273      475             2569            41683             0       0.0737
1785  Redundancy-1           0        0                0                0         45000       0.0000
```

#### 1786 **Ordered numeric character types: 79.23% correct**

```
1787  True model   Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1788  Expansion-1        25708    18686              432              173             1       0.4287
1789  Neutral            18913    25492              422              173             0       0.4335
1790  PartitioningR-1      291      375            41042             3284             8       0.0880
1791  PartitioningS-1       70       94             3790            41024            22       0.0884
1792  Redundancy-1           0        0                0                0         45000       0.0000
```

#### 1793 <u>**Other ecospace frameworks:**</u>

#### 1794 **Bush and Bambach ecospace framework (Fig. 2A): 74.85% correct**

```
1795  True model   Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1796  Expansion-1        24198  , 20127              448              182            45       0.4623
1797  Neutral          , 20157    24170              452              167            54       0.4629
1798  PartitioningR-1      357      267            37250             6702           424       0.1722
1799  PartitioningS-1      170      140             5278            37939          1473       0.1569
1800  Redundancy-1           0        4               24              121         44851       0.0033
```

#### 1801 **[Same, but using a validation data set: 72.47% for validation data set]**

```
1802   True model    Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1803   Expansion-1        21102    23300              376              169            53       0.5311
1804   Neutral          , 20814    23493              471              162            60       0.4779
1805   PartitioningR-1      425      351            36388             7363           473       0.1914
1806   PartitioningS-1      266      242             5206            37361        , 1925       0.1698
1807   Redundancy-1           3        0               32              244         44721       0.0062
```

#### 1808 **Manuscript ecospace framework, limiting binary characters to two mutually "present"**
#### 1809 **states ['constraint=2'] (Fig. 2B): 90.79% correct**

```
1810  True model   Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1811  Expansion-1        34808    10091               81            , 20             0       0.2265
1812  Neutral             8095    36758              113               34             0       0.1832
1813  PartitioningR-1       53       70            44060              817             0       0.0209
1814  PartitioningS-1       10       23             1062            43905             0       0.0243
1815  Redundancy-1           0        0                0                0         45000       0.0000
```

#### 1816 **[Same, but using a validation data set: 89.10% for validation data set]**

```
1817   True model    Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1818   Expansion-1        32566    12335               77               22             0       0.2763
1819   Neutral            10004    34863               96               37             0       0.2253
1820   PartitioningR-1       54       69            44135              742             0       0.0192
1821   PartitioningS-1       14       29             1054            43903             0       0.0244
1822   Redundancy-1           0        0                0                0         45000       0.0000
```

#### 1823 **Manuscript ecospace framework, limiting binary characters to single "present" state**
#### 1824 **['constraint=1'] (Fig. 3A): 91.08% correct**

```
1825  True model   Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
1826  Expansion-1        34898     9891              135               76             0       0.2245
1827  Neutral             7040    37686            , 191               83             0       0.1625
1828  PartitioningR-1      121      140            43638             1101             0       0.0303
1829  PartitioningS-1       47       50             1194            43709             0       0.0287
1830  Redundancy-1           0        0                0                0         45000       0.0000
```

**Manuscript ecospace framework, limiting binary characters to two mutually "present" states and weighing state occupation by frequency within empirical Kope/Waynesville-formation species pool ['constraint=2'] (Fig. 3B): 95.51% correct**

```
True model      Expansion-1  Neutral  PartitioningR-1  PartitioningS-1  Redundancy-1  class.error
Expansion-1         44489      259            63               72           117        0.0114
Neutral               304    42343          1697              502           154        0.0590
PartitioningR-1       129     2849         40394             1491           137        0.1024
PartitioningS-1       112      514          1444            42750           180        0.0500
Redundancy-1           33       11            11               13         44932        0.0015
```

**Manuscript ecospace framework, limiting binary characters to two mutually "present" states ['constraint=2'], weighing state occupation by frequency within empirical Kope/Waynesville-formation species pool, and including weakened (50%-rule-following) models (Fig. 3C in part)**

This trial illustrates the confusion matrix for the classification tree used in the manuscript to classify empirical Kope and Waynesville formation samples. The tree was trained on simulated data spanning nine models: the neutral model; the four models that implement the redundancy, two versions of partitioning, and expansion rules; and four "weakened" versions of these four models, in which rules were followed at random 50% of the time (i.e., the 100% and 50% training data sets visualized in Figure 3C). Model names are abbreviated, with "-1" referring to the 100%-rule-following implementation and "-0.5" referring to the 50% implementation. The classification tree included the total variance (V) statistic (unlike those above).

The confusion matrix was produced from a validation data set, i.e., samples excluded from the training data set used to build the classification tree. The training data set included 5,265,000 samples (9 models * 3000 simulations * 195 sample sizes [5–200 species per sample]) and the validation data set (whose confusion matrix is illustrated below) included 855,000 samples (1000 simulations * 95 sample sizes [5–100 species per sample] = 95,000 samples per model * 9 models). The classification rate on the training data set was 91.7% (79.9% when restricted to 1000 samples, sample sizes below 50, and excluding V from the training data set, for comparative purposes with above confusion matrices) and the overall classification rate on the validation data set was 78.1% (75.2% when restricted to sample sizes below 50 and excluding V, for comparative purposes with above confusion matrices). See Supplementary Figure 9 for an alternate version of the confusion matrix, illustrating performance as a function of sample size and including additional validation samples.

```
Truth       Exp-0.5  Exp-1  Neutral  PartR-0.5  PartR-1  PartS-0.5  PartS-1  Redund-0.5  Redund-1  class.error
Exp-0.5       92049    816      937        285       82        327       80         288       136      0.0311
Exp-1           713  93969       68         35       19         27       12          49       108      0.0109
Neutral        1275    112    65234      12889      581       7598      208        6974       129      0.3133
PartR-0.5       895    102    17747      47997     2468      19126      403        6120       142      0.4948
PartR-1         265     54     1116       2615    83205       3416     1253        2939       137      0.1242
PartS-0.5      1164     83    13221      26207     2442      36869      731       14157       126      0.6119
PartS-1         151     51      158        218     1261        378    91424         425       934      0.0376
Redund-0.5      810    151     9706       5842     2555      12988      725       62071       152      0.3466
Redund-1         27     31        9         16       21         11       30          16     94839      0.0017
```

**Supplementary Appendix 3. Spreadsheet with life-habit/functional-trait codings for the Kope and Waynesville Formation species pool.**

KWTraits.csv is a comma-separated value (.csv) format file listing the aggregate species pool for the Kope and Waynesville Formation used in empirical analyses. (The file is also included as a data file within the ecospace R package.) The first three columns list taxonomic information. The remaining columns list ecospace character states (functional traits). See Supplementary Appendix 1 and Novack-Gottshall (2007) for information on characters and states. See text for explanation of how multistate characters were rescaled.


**Supplementary Appendix 4. Model-selection support data files for Kope and Waynesville Formation samples, stratigraphic section, member, and formation aggregates.**

Files are in comma-separated value (.csv) format. The first five columns describe the Paleobiology Database collection identification number, scale (hand sample, stratigraphic section, etc.) of the sample, and stratigraphic/section names. Columns 6–14 list sample size (S, species richness) and values for eight disparity statistics (with NA designating when a statistic could not be calculated, because there were fewer than four unique life habits in the sample); see text for descriptions and abbreviations of statistics. The last column identifies which model has the best support among those candidates considered. The remaining columns list the classification-tree support each sample has for each candidate model considered.

emp2-modelfits.csv lists model support using the classification tree trained on the 50% and 100%-strength training data sets. emp3-modelfits.csv lists model support for the tree trained on 50%, 90%, and 100% training data. emp5-modelfits.csv lists model support for the tree trained on 50%, 75%, 90%, 95%, and 100% training data.