

Information Theory Broadens the Spectrum of Molecular Ecology and Evolution

Rutgers University has made this article freely available. Please share how this access benefits you.
Your story matters. [\[https://rucore.libraries.rutgers.edu/rutgers-lib/54634/story/\]](https://rucore.libraries.rutgers.edu/rutgers-lib/54634/story/)

This work is an **ACCEPTED MANUSCRIPT (AM)**

This is the author's manuscript for a work that has been accepted for publication. Changes resulting from the publishing process, such as copyediting, final layout, and pagination, may not be reflected in this document. The publisher takes permanent responsibility for the work. Content and layout follow publisher's submission requirements.

Citation for this version and the definitive version are shown below.

Citation to Publisher Sherwin, William B., Chao, Anne, Jost, Lou & Smouse, Peter E. (2017). Information Theory Broadens the Spectrum of Molecular Ecology and Evolution. *Trends in Ecology And Evolution* 32(12), 948-963. <https://doi.org/10.1016/j.tree.2017.09.012>.

Citation to this Version: Sherwin, William B., Chao, Anne, Jost, Lou & Smouse, Peter E. (2017). Information Theory Broadens the Spectrum of Molecular Ecology and Evolution. *Trends in Ecology And Evolution* 32(12), 948-963. Retrieved from [doi:10.7282/T3Z89GD4](https://doi.org/10.7282/T3Z89GD4).



Terms of Use: Copyright for scholarly resources published in RUcore is retained by the copyright holder. By virtue of its appearance in this open access medium, you are free to use this resource, with proper attribution, in educational and other non-commercial settings. Other uses, such as reproduction or republication, may require the permission of the copyright holder.

Article begins on next page

Information Theory Broadens the Spectrum of Molecular Ecology and Evolution

W.B. Sherwin ^{1*}

A. Chao ²

L. Jost ³

P.E. Smouse ⁴

Addresses:

¹ Evolution and Ecology Research Centre, School of Biological Earth and Environmental Science,
University of New South Wales, Sydney NSW 2052 AUSTRALIA.

and

Murdoch University Cetacean Research Unit, Murdoch University, South Road, Murdoch, WA 6150
AUSTRALIA.

² Institute of Statistics, National Tsing Hua University, Hsin-Chu 30043, Taiwan.

³ EcoMinga Foundation, Via a Runtun, Baños, Tungurahua, Ecuador

⁴ Department of Ecology, Evolution and Natural Resources, School of Environmental and Biological
Sciences, Rutgers University, New Brunswick, NJ 08901-8551, USA,

* Correspondence W.Sherwin@unsw.edu.au

Keywords: Shannon; diversity-profile; entropy; selection; linkage-disequilibrium; gene-expression

Word-length for main text (excluding comments such as [insert fig ** here]): 3422

28 **Highlights**

29

30 Diversity of molecules or species is best expressed in terms of a diversity profile

31 Such profiles are useful in studies spanning bioinformatics to physical landscapes

32 Shannon information is a neglected but particularly informative part of the profile

33 Shannon now has robust theoretical background for molecular ecology and evolution

34

35

36 **Abstract**

37 Information or entropy analysis of diversity is used extensively in community ecology, and has
38 recently been exploited for prediction and analysis in molecular ecology and evolution. Information
39 measures belong to a spectrum (or 'q-profile') of measures whose contrasting properties provide a
40 rich summary of diversity, including allelic richness ($q=0$), Shannon information ($q=1$), and
41 heterozygosity ($q=2$). We present the merits of information measures for describing and
42 forecasting molecular variation within and among groups, comparing forecasts with data, and
43 evaluating underlying processes such as dispersal. Importantly, information measures directly link
44 causal processes and divergence outcomes, have straightforward relationship to allele frequency
45 differences (including monotonicity that $q=2$ lacks), and show additivity across hierarchical layers
46 such as ecology, behaviour, cellular processes, and non-genetic inheritance.

47

48 **An Information Theory Base for Evolutionary Genetics and Ecology**

49 [\[Insert glossary sidebar about here\]](#)

50 Evolutionary ecology aims to make and test forecasts about the behaviour of variants, at all levels
51 from molecular through species to landscapes, but until recently this field has paid little attention
52 to one of three main ways for doing this. Shannon **information (or entropy) theory** (see Glossary)
53 is widely-used to predict macro-patterns from micro-parts, by methods such as finding the model
54 that maximises the entropy [1], in fields as diverse as community ecology [2] and energy
55 generation, where these methods compare favourably with alternative predictions [3]. Genes
56 obviously carry information – about evolutionary history, recent demography, and possible future

57 trajectories [4, 5] – but information theory has rarely been used to investigate molecular evolution
58 [6-9]. Shannon’s information index ${}^1H (=H)$ [10], a fundamental component of information theory,
59 is the most commonly used abundance-sensitive measure of species diversity within a community
60 [11]:

61
$${}^1H = -\sum_{i=1}^S p_i \ln p_i$$
 Equation 1

62 where p_i is the proportional abundance of the i^{th} species in a community of S different species
63 (Tutorial Box 1, and Supplement Box S1, whose footnote has definitions of all symbols). 1H also
64 applies to a population containing genetic variants (such as allelic types) [12], and can be thought
65 of as the ability to spell out different messages by rearranging individual alleles (Box 1-l) – with
66 higher diversity, a greater range of messages can be spelt out. More formally, higher 1H means that
67 there is reduced certainty about what type to expect when a single allele is randomly sampled.

68

69 [Insert Box 1 about here]

70

71 Historically, molecular ecology has quantified diversity with two alternative entropies (Boxes 1,S1):
72 ${}^0H = S-1$; and 2H (Boxes 1,S1). 2H is the chance of choosing two different allelic types from the
73 population, called ‘heterozygosity’ (${}^2H = H_e$, or for species in communities, called Gini-Simpson
74 and variants) [4, 5, 13].

75
$${}^2H = 1 - \sum_{i=1}^S p_i^2$$
 Equation 2

76 With higher diversity, the chance of randomly choosing two different alleles increases, so 2H
77 increases.

78

79 To create the ' q -profile' or 'spectrum', one converts (${}^qH = {}^0H, {}^1H, {}^2H, \text{ etc.}$) to a common scale
80 of 'effective numbers' (${}^qD = {}^0D, {}^1D, {}^2D \text{ etc}$), which represent the number of equally-frequent
81 alleles that would be needed to yield the observed qH in the sample, which typically contains
82 alleles at unequal frequencies [14-16] (Boxes 1-II,S1,S2). The qD -measures are sometimes called
83 'true diversity', and the use of the qD profile has long been recommended in ecology, because each
84 q -value provides different insights into the composition of the diversity [17, 18]. For example,
85 higher q -values emphasise the more common variants (Box 1-I,II). As well as these 'alpha'-
86 measures for alleles in a single population, each q -value has measures of diversification among
87 locations, which are the beta-measures (Box 1-III) that are critical in evolution and conservation [19,
88 20]. Finally, the total diversity within and among localities is labelled gamma diversity.

89

90 Partial qD profiles are already used in evolution and ecology, to assess community response to
91 changed conditions [19], as well as possible correlations between species- and gene-diversity [21].
92 Also, methods to infer selection or population-size changes exploit the way these processes have
93 different effects on number of variants S ($q=0$) and heterozygosity ($q=2$) [22-24]). Many recent
94 publications include ($q=1$) (Box S3), yet few authors [25-28] have exploited systematic q -profiles
95 (Box 1-II), to obtain the power described further below.

96

97 The diversity profile is most useful if we can forecast its shape under specified histories of selection,
98 population size, dispersal, etc. We can then either test departures from those predictions and/or

99 estimate forces that underlie the diversity patterns we encounter. Classically, ($q=2$) theory predicts
100 values of within- and among-population measures, including heterozygosity 2H , $Jost-D$, and F_{ST} ,
101 under conditions such as neutrality, selection, subdivision, dispersal, altered population size [4, 5,
102 29]. Below we show that after a slow start [30-34], ($q=1$) predictive theory is now catching up to
103 ($q=2$) theory [12, 35-37], and 1H has been proposed as a primary measure of evolvability [38, 39].

104

105 As a reading guide, the initial sections outline how ($q=1$) molecular information can now predict
106 and measure processes such as adaptation and dispersal (with additional detail in supplements).
107 Those wanting less technical detail might skip directly to the penultimate section evaluating
108 strengths and weaknesses of each element of the q -profile, such as the effect of q on sensitivity to
109 rare alleles (Box 1-I), and the poor performance of some conventional ($q=2$) measures for analysis
110 of among-population (β) differentiation.

111

112 We consider both neutral and adaptive genetic variants, and also discuss haploid, diploid and clonal
113 organisms. We consider genes ('loci') with only two variants ('alleles'), such as typical SNPs (single-
114 nucleotide polymorphisms), as well as multiallelic loci, evolving under either the SMM (stepwise
115 mutation model - some microsatellite loci) or the IAM (infinite alleles model for 'haplotypes' or
116 'haptigs' [40] of variants linked on the same DNA molecule [5]). We also discuss continuous traits
117 determined by variation of multiple genes.

118

119 **Measuring and Predicting Shannon's Information for Neutral Alleles**

120 Within-population (alpha) equations have recently been derived to predict Shannon entropy 1H
121 and diversity 1D , within-populations (alpha), for neutral genetic information (having no effect on
122 adaptation), for equilibrium with constant effective population size N_e and mutation rate μ (Box
123 S1). The new equations show good fit to both simulated and real data sets for loci evolving under
124 several mutation processes - SNP [36], SMM [12, 37] and IAM [12, 37]. Boxes S1 and S6 show how
125 to calculate 1H , to compare data with theoretical predictions. For example, in a rainforest tree
126 *Elaeocarpus sedentarius*, calculated and predicted 1H agreed [12, 41]. There are many recent
127 examples of the empirical and theoretical use of Shannon measures to assess how molecular
128 diversity is affected by factors such as: endangerment and bottlenecks; subdivision; environmental
129 gradients; pedigree inbreeding; and invasions expansions introductions and host jumps (Box S3). In
130 most cases, measures with different q -values yield similar interpretations, but we later discuss
131 informative cases where they differ.

132

133 Beta differentiation of molecular diversity among populations, species, and landscapes can be
134 summarised using measures of each order of q . The ($q=0$) β -measures are based on the proportion
135 of allelic types that that are not shared by two populations (Box S1). *Jost-D* is a $q=2$ β -measure,
136 and there are related measures such as F_{ST} [29, 42] and pairwise comparison algorithms, including
137 STRUCTURE [43] and AMOVA [44] (Box S1). The ($q=1$) measures were explicitly designed to be
138 hierarchical [1, 10], so partitioning of molecular diversity is particularly easy, leading to
139 differentiation measures Mutual Information I , and Shannon Differentiation. These can be
140 derived from a contingency table of differentiation of allele frequency between geographic
141 locations (Boxes 2,S1) [1, 12, 33, 35, 37]. Box 2 shows how dimensions can be added to

142 incorporate diversity within and among different habitats, landscapes, etc. [35, 45-47], because log-
143 linear χ^2 (or I) is completely additive [48].

144 [Insert Box 2 about here]

145 For any pair of populations, lower dispersal, smaller population size, or greater elapsed time since
146 separation will increase molecular differentiation (and hence, mutual information I). At
147 equilibrium, simulation results have been used to derive an inverse relationship between mutual
148 information (I , $q=1$) and effective dispersal rate ($N_e m$), over a very wide range of effective
149 population sizes ($N_e \geq 10$) and dispersal rates ($0.001 \leq m \leq 0.30$ *i.e.* 30% dispersing per
150 generation) [12] (Box S1). This equation can be used to estimate dispersal from genetic data, and
151 outperforms predictions based on ($q=2$) in simulation studies (Figure 1A), as well as in laboratory
152 colonies of *Drosophila* with known $N_e m$ [12]. The waratah *Telopea speciosissima* showed a strong
153 negative correlation between F_{ST} ($q=2$) and $N_e m$ estimated from mutual information I ($q=1$), as
154 expected because F_{ST} and I are each inversely related to $N_e m$ [49]. Dispersal can also be assessed
155 using I in clonal species [50], haploids, and for other ploidies such as a three-species hybrid moss
156 *Sphagnum x falcatulurm* [45].

157 [Insert Fig 1 about here]

158

159 Other equations predicting equilibrium I are based on rigorous theory, rather than simulations,
160 but require knowledge of mutation rate. An equation for bi-allelic SNPs [36] (Box S1) fits closely to
161 simulation outcomes, as well as to data from laboratory colonies of *Drosophila* of known $N_e m$
162 (Figure 1B). Also, equations for IAM and SMM loci [37] (Box S1) agree with simulation outcomes,

163 and with observed genetic divergence for SMM (microsatellite) loci in starlings (*Sturnus vulgaris*)
164 introduced to Australia [37]. There are many other recent examples of the use of molecular mutual
165 information I to investigate hypotheses about temporal or spatial environmental gradients or
166 barriers, in a variety of free-living and parasitic plants, animals, and fungi, plus simulations (Box S3).
167 In most cases, different q -values show similar results, but we discuss divergent cases later.
168 Calculations, sampling and programs, are in Box S6.

169

170 **Using Information Methods to Scan the Genome for Adaptive Innovation**

171 There is a boom in searches for adaptive genomic regions, for evolutionary interest, choosing
172 reintroduction sources for conservation [51] and identifying human disease loci [52, 53]. Strategies
173 to detect these regions include ‘evolve-resequence’ or ‘transplant’ experiments [54], and assessing
174 genomic differentiation over landscapes (‘landscape genomics’) [55]. Selection acts via differential
175 fitness of variants of any heritable characteristic: DNA sequence variation, expression variability
176 due to interaction between environment and multiple genes, or non-genetic inheritance such as
177 epigenetics [56, 57]. Partly because of these many modes, there are many signals of selection, but
178 each has a high rate of false-positives [58], which can be minimised by using multiple alternative
179 approaches, including ($q=1$) methods.

180

181 For detecting selection, Shannon-based metrics are particularly appropriate, in ways such as their
182 greater sensitivity (than $q=2$) to rare alleles that may be important for conservation management or
183 detecting new (potentially adaptive) alleles [59]. For directional selection, which favours a single

184 advantageous allele [5], Shannon information is proposed as a natural measure of evolvability [38,
185 39]. The ultimate result will be a single allele at 100%, so 1H tends to zero, with dynamics analysed
186 by logit transformation that is algebraically related to mutual information between the allele
187 distributions, before and after selection [39, 60]. (Box S4). Unlike directional selection, balancing
188 selection maintains equilibrium proportions of two or more alleles, for example due to high
189 heterozygote fitness [5], and the predicted allele proportions can be used to calculate the
190 equilibrium value of Shannon information 1H (Box S4) [35]. In addition, much adaptive variation is
191 based on quantitative traits controlled by multiple loci. Information approaches to directional
192 selection on multilocus traits (Box S4 [61-64]) indicate that such evolution favours gene duplication
193 [65]. The multilocus analog of balancing selection is called stabilising selection, for which there are
194 also treatments based on information statistics [66].

195

196 **Detecting when Selection Changes Patterns in either Populations or Expression Profiles**

197 For diploids, departure from single-locus Hardy-Weinberg expectations (HWE) due to selection,
198 non-random mating and other forces is summarised by F_{IS} in ($q=2$) form (Box S4) [4, 5].

199 Equivalent for ($q=1$) include the conventional log-linear-chisquare for fit to random-mating
200 expectations (HWE Box S4), as well as extensions to mixed mating systems, such as mixed selfing
201 and outcrossing, in plants and invertebrates (Boxes S1, S3, S6.3) [12, 34, 67].

202

203 Analysis of selection, and other evolutionary processes, requires us to evaluate 'haplotypes' of
204 variants in multiple parts of the genome, including SNPs, insertions, deletions [40]. Discriminating

205 the individual and combined effects of haplotype elements requires analysis of ‘Linkage
206 disequilibrium’ (LD) which is non-random association among such variants, due to physical linkage
207 (‘true LD’ [68]), historical population size and admixture [69], or ‘epistatic’ selection in which
208 environmental conditions favour particular combinations of variants at multiple positions [7, 52, 53,
209 70]. LD can be expressed with Mutual Information I ($q=1$) for association between variants at
210 different positions, which minimises false-discovery of LD (Boxes S5, S6.3) [7, 47, 52, 53, 70-72] and
211 can analyse multi-locus associations of both continuous normal [52] or non-normal [71] traits. This
212 approach has identified the combined effect of SNPs in two protein-coding loci, upon an addictive
213 phenotype [73]. Also, I can be used to assess whether adaptive changes have resulted in different
214 haplotypes in different locations (Box S5) [47, 74, 75]. Physical linkage is not the only way genes
215 interact, and information methods are used to integrate the potentially selective effects of
216 expression networks, DNA sequence frequencies, and linkage disequilibrium [7, 70, 76, 77], for
217 example helping us predict phenotypic outcomes of RNA modification (splicing) in blood-clotting
218 factors [78]. A partial profile ($q = 1,2$) proved best for assessing expression patterns [79].
219 Programs for analyses are in Box S6.3.

220

221 **Choosing Measures – the Relative Merits of Information-based Measures** 222 **and Other Measures.**

223 A full diversity profile exploits the sensitivities, and strengths, of each element ($q = 0, 1, 2$).
224 Reassuringly, they often give similar results (Box S3), but to understand cases where they differ, or
225 to choose measures for particular applications, we must assess the sensitivities, strengths and
226 weaknesses of each profile element ($q=0,1,2$). For example, sensitivity to rare alleles may be an

227 advantage for conservation or adaptation studies (Box 3), but this sensitivity means that missing
228 rare alleles can seriously degrade ($q=0$) estimates unless appropriate corrections are made (Box
229 S6.2). The performance of measures must be evaluated for each evolutionary or ecological
230 question, such as change of molecular diversity with latitude, and estimation of dispersal from
231 genetic markers (where mutual information outperforms other methods [12]). Here we point out
232 some major differences between ($q=1$) versus ($q=0$) or ($q=2$).

233

234 Shannon results generally accord with biological predictions, but in some cases the full q -profile
235 aids detection (or rejection) of a predicted pattern, because the ($q=1$) measure and other measures
236 diverge. Sometimes ($q=1$) has greater sensitivity than do ($q=0$) or ($q=2$) measures (Box S3), as was
237 seen in two studies of alpha (within-population) diversity due to adaptation and drift. In an invasive
238 mosquito *Aedes japonicus japonicus* [80], and in crop carrots affecting nearby wild carrots *Daucus*
239 *carota* [81], it was found that 2H (heterozygosity $q=2$) was less sensitive than 1H (Shannon $q=1$) to
240 the predicted loss of variability resulting from small population size and colonisation. This is
241 presumably because 2H emphasizes common alleles that have relatively low chance of being lost
242 during a bottleneck that was caused by either massive mortality or small numbers of invaders. A
243 **hierarchical** partition of Shannon information (similar to Box 2-IV) in three-species hybrid moss,
244 *Sphagnum x falcatulurm*, revealed differentiation due to rare recent mutations to which ($q=2$)
245 metrics would not have been sensitive [45]. However, ($q=1$) does not always show the greatest
246 sensitivity, in temporal studies of zooplankton, ($q=0$) was more powerful than ($q=1$), which was
247 more powerful than ($q=2$) [28]. Examples in this paragraph demonstrate the importance of
248 analysing the entire q -profile including ($q=1$).

249

250 While recognising the importance of the entire q -profile, Table 1 and Boxes 3 and 4 show that only
251 ($q=1$) measures combine a number of desirable characteristics for tracking evolutionarily important
252 phenomena. Firstly, Shannon measures are sensitive to alleles according to their frequency, yet
253 have minimal sampling problems (Table 1, Boxes 1,3). Secondly, measures of each q -order must be
254 able to distinguish levels of diversity that are individually important for evolution and conservation
255 [19, 20]: within-locality (α), among-locality differentiation (β), and total (γ).

256 Eliminating dependencies between these levels for many common ($q=2$) methods requires more
257 complex equations or algorithms [42], which might complicate predictions. In contrast, the explicit
258 hierarchical nesting of information measures (Figure 6 in [10]), yields independent estimates of
259 within-population (α), among-population (β), and total (γ) levels of diversity (Table 1, Box 4).

260 Disagreements between beta $q=1$ and $q=2$ measures may sometimes derive from these
261 dependencies ([81, 82]. Thirdly, some ($q=1$) measures respond in an intuitively appealing fashion to
262 addition of new alleles, behaving in a predictable way with the number of unshared alleles, and
263 satisfying the principles of strong monotonicity and replication (Table 1, Box 4). ‘Replication’ [15,
264 83], means that a measure increases linearly when equally diverse and completely distinct groups
265 are pooled in equal proportions, as seen in effective numbers (qD , Boxes 1,S1). Information
266 measures also have minimal sampling problems if appropriate estimation methods are used (Box
267 S6.2). In contrast, a drawback of ($q=0$) measures is that although they can be forecast under
268 specified conditions [5], the sampling to test these forecasts can be hampered by their extreme
269 sensitivity to rare alleles (Boxes 1,3,4,S1,S2).

270 [Insert Table 1 and Boxes 3 and 4 about here]

271 Additionally, we have shown above that Shannon measures are catching up in areas where they
272 have lagged behind ($q=0$) and ($q=2$): forecasting methods are now available for ($q=1$) measures of
273 neutral and adaptive variation, and the frequency of use of ($q=1$) methods in molecular ecology is
274 rising to meet their already heavy use in community ecology. Also, we have shown that Shannon
275 ($q=1$) methods outperform others for genetic estimation of dispersal [12], and have great utility in
276 detecting selection and gene expression patterns [75, 76]. Thus we should use information
277 methods as important contributors to the diversity profile.

278

279 **The Near Future: Challenges, Opportunities and Integration**

280 [Box 5 Outstanding questions as sidebar about here]

281 The major components of molecular information theory are now established, but outstanding
282 questions remain (Box 5). In the long term, the most useful measure will undoubtedly be the
283 whole qD profile (0D , 1D , 2D) at α , β , and γ levels. This will maximise understanding of patterns,
284 and allow meta-analysis of the q -profile's performance under a wide variety of conditions.

285

286 To complete an Information-based strategy for forecasting and analysis in molecular ecology and
287 evolution, Shannon's strong performance with dispersal, mutation and drift (above) must be
288 extended to include asymmetric dispersal, unequal or changing population sizes, and variants with
289 mutation not described by models such as infinite (IAM) stepwise (SMM) or biallelic (SNP). This
290 could be through new theory, or *via* Approximate Bayesian Computation using the q -profile [12,
291 27]. It would also be ideal to develop an analogue to AMOVA [44], based on information-theoretic

292 methods. This might build upon a molecular phylogenetic differentiation measure for all q [84],
293 based on the neutral coalescent of evolutionary biology [4].

294

295 As well as neutral predictions, we are only beginning to explore the capacity of information science
296 to integrate all analyses of adaptation and selection into a common scale, possibly exploiting the
297 connection to logit (log-linear) modeling. Moreover, sensitivity of tests for loci under selection is
298 likely to increase with a ($q=1$) approach, because they are more sensitive to novel rare variants that
299 are crucial in adaptation [59], and because many of these tests rely upon partition of among-
300 population (beta) *versus* within-population (alpha) measures of diversity [54, 55, 85], and so should
301 benefit from the complete independence of α and β in the ($q=1$) scale.

302

303 Understanding adaptation requires integrating all aspects of biology from sub-cellular biology to
304 habitat tolerance, and information theory is particularly well suited to this challenge [6, 8, 9], being
305 both explicitly additive and hierarchical, as well as being a general forecasting method [1], already
306 used in community ecology [2]. Community ecologists have borrowed genetic ($q=2$) theory for both
307 neutral [13] and adaptive genes [21], so community ecology might also benefit from deploying the
308 ($q=1$) neutral and selective theory outlined in this article. They will likely be able to make even
309 more forecasts than with ($q=2$), by virtue of the fact that Shannon applies to all types of
310 information [86], including physical habitat, behaviour [87], genetics [12, 35-37], information within
311 each sequence including codon usage biases [88], and various non-genetic inheritance modalities
312 including epigenetics [56, 57]. Moreover, information has been proposed as the driver of any
313 correlations between species and molecular diversity [89]. Incorporating information on genetic
314 and/or functional similarity or difference of alleles (or species) greatly improves the utility of

315 diversity measures [83, 90], and ongoing attempts to incorporate function, without violating
316 fundamental diversity principles, should include Shannon methods [83, 91].

317

318 We have only presented a limited aspect of information theory, but we could paint on a much
319 broader canvas. Other potentially fruitful aspects of information theory include: Kullback-Leibler or
320 Relative entropy used in selection analysis; 'Fisher information' used to compare ($q=2$) measures
321 [92]; fractals [93]; plus compression algorithms proposed for both selection analysis [63] and for
322 joint estimation of sequence alignments and phylogenies [94]. Finally, predictive equations for
323 molecular Shannon information might be used to advance evolutionary computing (machine
324 learning), in which code's performance is optimised by a process similar to biological drift and
325 selection [76, 87, 95, 96]

326

327 **Conclusions**

328 Four common processes - Dispersal, Random drift, Adaptation, and Generation of novelty
329 (speciation, mutation, recombination) – unify community and molecular ecology [13,21] . The
330 unification of these fields will be accelerated by cross-fertilisation between their predictive
331 equations and measures, at levels from bioinformatics to physical landscapes and beyond,
332 especially when using a diversity profile of $q=0,1,2$, and adhering to criteria for robust measures
333 (see "Choosing Measures ..." section). Shannon Information measures ($q=1$) are a vital part of this
334 profile, now have predictive equations available, and excel at providing intuitive results.

Supplementary Material

335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358

Prediction, Sampling, Estimation, and Examples - abbreviated titles

Box S1 Measuring and Predicting Molecular Entropies and Diversities for ($q=1,2,3$), Within-populations (α) and Between-populations (β)

Box S2 Derivation of the q -profile

Box S3 Studies using Molecular Information Measures plus q -profile

Box S4 Selection

Box S5 Selection and Linkage Disequilibrium (LD)

Box S6.1 Molecular Information Analyses – Example Calculations

Box S6.2 Molecular Information Analyses – Sampling Considerations

Box S6.3 Molecular Information Analyses – Estimation Programs

Box S7 References

Acknowledgements

We thank reviewers and colleagues who read drafts and made valuable comments: Bragg, Brandenburger, Byrnes, Carmel, Chabanne, Daly, Jabot, Kidd, Kopps, Kosman, Leinster, Manlik, Maslen, Nichols, O'Reilly, Peakall, Raphael, Reeve, Rollins, Shofner, Stear, Sunnucks, Tanaka. TREE editor Craze also made extremely valuable suggestions. WS, AC and LJ first met in 2012 at a workshop coordinated by Leinster Reeve *et al.*, on “The Mathematics of Biodiversity” at the Centre de Recerca Matemàtica, Universitat Autònoma de Barcelona, funded by that centre, the Spanish government, and the BBSRC (International Workshop Grant BB/J020567/1).

359 **References**

360

- 361 1. Cover, T.M. and Thomas, J.A. (1991) Elements of information theory Wiley.
- 362 2. Elith, J. et al. (2011) A statistical explanation of MaxEnt for ecologists. *Diversity and Distributions* 17, 43–
- 363 57.
- 364 3. Bessa, R.J. et al. (2009) Entropy and correntropy against minimum square error in offline and online three-
- 365 day ahead wind power forecasting. *IEEE Transactions on Power Systems* 24 (4), 1657-1666.
- 366 4. Nielsen, R. and Slatkin, M. (2013) An introduction to population genetics theory and applications, Sinauer.
- 367 5. Halliburton, R. (2004) Introduction to population genetics, Pearson.
- 368 6. Searls, D.B. (2010) The roots of bioinformatics. *PLoS Computational Biology* 6 (6).
- 369 7. Moore, J.H. and Hu, T. (2015) Epistasis analysis using information theory. In *Epistasis: methods and*
- 370 *protocols, methods in molecular biology* (Moore, J.H. and Williams, S.M. eds).
- 371 8. Glazebrook, J.F. and Wallace, R. (2012) ‘The frozen accident’ as an evolutionary adaptation: a rate
- 372 distortion theory perspective on the dynamics and symmetries of genetic coding mechanisms. *Informatica*
- 373 36, 53–73.
- 374 9. Skene, K.R. (2015) Life’s a gas: a thermodynamic theory of biological evolution. *Entropy* 17, 5522-5548.
- 375 10. Shannon, C.E. (1948) A mathematical theory of communication. *The Bell System Technical Journal* 27,
- 376 379–423, 623–656.
- 377 11. Buddle, C.M. et al. (2004) The importance and use of taxon sampling curves for comparative biodiversity
- 378 research with forest arthropod assemblages. *Canadian Entomologist* 137:120-127.
- 379 12. Sherwin, W.B. et al. (2006) Measurement of biological information with applications from genes to
- 380 landscapes. *Molecular Ecology* 15, 2857-2869.
- 381 13. Hubbell, S.P. (2001) *The unified neutral theory of biodiversity and biogeography*, Princeton University
- 382 Press.
- 383 14. Jost, L. (2006) Entropy and diversity. *Oikos* 113, 363–375.
- 384 15. Chao, A. et al. (2014) Unifying species diversity, phylogenetic diversity, functional diversity, and related
- 385 similarity and differentiation measures through Hill numbers. *Annual Review of Ecology, Evolution, and*
- 386 *Systematics* 45, 297-324.
- 387 16. Gregorius, H.R. and Kosman, E. (2016) On the notion of dispersion: from dispersion to diversity. *Methods*
- 388 *in Ecology and Evolution*, DOI: 10.1111/2041-210X.12665.
- 389 17. Pielou, E.C. (1966) The measurement of diversity in different types of biological collections. *Journal of*
- 390 *Theoretical Biology* 13, 131-144.
- 391 18. Hill, M.O. (1973) Diversity and evenness: a unifying notation and its consequences. *Ecology* 54, 427-432.
- 392 19. Anderson, M. et al. (2011) Navigating the multiple meanings of beta diversity: a roadmap for the
- 393 practising ecologist. *Ecology Letters* 14, 19-28.
- 394 20. Socolar, J.B. et al. (2016) How should beta-diversity inform biodiversity conservation? *Trends in Ecology*
- 395 *& Evolution* 31, 67-80.
- 396 21. Vellend, M. (2016) *The Theory of Ecological Communities* (MPB-57), Princeton University Press.
- 397 22. Luikart, G. et al. (1998) Distortion of allele frequency distributions provides a test for recent population
- 398 bottlenecks. *Journal of Heredity* 89 (3), 238-247.
- 399 23. Tajima, F. (1989) Statistical method for testing the neutral mutation hypothesis by DNA polymorphism.
- 400 *Genetics* 123, 585-595.
- 401 24. Fu, Y.-X., Li, W-H. (1993) Statistical tests of neutrality of mutations. *Genetics* 133, 693-709.
- 402 25. Greenbaum, G. et al. (2014) Allelic richness following population founding events – a stochastic modeling
- 403 framework incorporating gene flow and genetic drift. *PLoS ONE* 9 (12), e115203.
- 404 26. Lloyd, M.W. et al. (2013) The power to detect recent fragmentation events using genetic differentiation
- 405 methods. *PLoS ONE* 8 (4), e63981.
- 406 27. Huang, W. et al. (2011) Approximate Bayesian comparative phylogeographic inference from multiple
- 407 taxa and multiple loci with rate heterogeneity. *BMC Bioinformatics* 12, 1.

- 408 28. Iacchei, M. et al. (2017) It's about time: Insights into temporal genetic patterns in oceanic zooplankton
409 from biodiversity indices. *Limnol. Oceanogr.* 00, 00-00.
- 410 29. Jost, L. (2008) G_{st} and its relatives do not measure differentiation. *Mol. Ecol.* 17, 4015-4026.
- 411 30. Crow, J.F. (2001) Shannon's brief foray into genetics. *Genetics* 159, 915-917.
- 412 31. Ewens, W.J. (1972) The sampling theory of selectively neutral alleles. *Theoretical Population Biology* 3,
413 87-112.
- 414 32. Lewontin, R.C. (1972) The apportionment of human diversity. *Evolutionary Biology* 6, 381-398.
- 415 33. Brown, A.H.D. and Weir, B.S. (1983) Measuring genetic variability in plant populations. In *Isozymes in*
416 *Plant Genetics and Breeding* (Tanksley, S.D. and Orton, T.J. eds), pp. 219–239, Elsevier Science Publications.
- 417 34. Meirmans, P.G. and Van Tienderen, P.H. (2004) GENOTYPE and GENODIVE: two programs for the analysis
418 of genetic diversity of asexual organisms. *Molecular Ecology Notes* 4, 792–794.
- 419 35. Sherwin, W.B. (2010) Review: entropy and information approaches to genetic diversity and its
420 expression: genomic geography. *Entropy* 12, 1765-1798.
- 421 36. Dewar, R.C. et al. (2011) Predictions of single-nucleotide polymorphism differentiation between two
422 populations in terms of mutual information. *Molecular Ecology* 20, 3156–3166.
- 423 37. Chao, A. et al. (2015) Expected Shannon entropy and Shannon differentiation between subpopulations
424 for neutral genes under the finite island model. *PLoS ONE* 10 (6), e0125471.
- 425 38. Day, T. (2015) Information entropy as a measure of genetic diversity and evolvability in colonization.
426 *Molecular Ecology* 24, 2073–2083.
- 427 39. Frank, S.A. (2017) Universal expressions of population change by the Price equation: natural selection,
428 information, and maximum entropy production. *Ecology and Evolution* 7 (10), 3381–3396.
- 429 40. Chin, S.-C. et al. (2016) Phased diploid genome assembly with single molecule real-time sequencing
430 bioRxiv, dx.doi.org/10.1101/056887.
- 431 41. Rossetto, M. et al. (2008) Dispersal limitations, rather than bottlenecks or habitat specificity, can restrict
432 the distribution of rare and endemic rainforest trees. *American Journal of Botany* 95, 321–329.
- 433 42. Hedrick, P.W. (2005) A standardized genetic differentiation measure. *Evolution* 59 (8), 1633-1638.
- 434 43. Pritchard, J.K. et al. (2000) Inference of population structure using multilocus genotype data. *Genetics*
435 155, 945–959.
- 436 44. Excoffier, L. et al. (1992) Analysis of molecular variance inferred from metric distances among DNA
437 haplotypes: application to human mitochondrial DNA restriction data. *Genetics* 131 (2), 479-491.
- 438 45. Karlin, E.F. and Smouse, P.E. (2017) Allo-triploid *Sphagnum x falcatulum*: single individuals contain most
439 of the Holantarctic diversity for ancestrally indicative markers. *Ann. Botany In Press*
440 dx.doi.org/doi:10.7282/T3FX7CRH.
- 441 46. Smouse, P.E. et al. (2015) An informational diversity analysis framework, illustrated with sexually
442 deceptive orchids in early stages of speciation. *Molecular Ecology Resources* 15, 1375-1384.
- 443 47. Smouse, P.E. and Ward, R.H. (1978) A comparison of the genetic infrastructure of the Yecua and the
444 Yanomama: a likelihood analysis of genotypic variation among populations. *Genetics* 33, 611-631.
- 445 48. Sokal, R.R. and Rohlf, F.J. (1995) *Biometry: the principles and practice of statistics in biological research.*,
446 Freeman.
- 447 49. Rossetto, M. et al. (2011) The impact of distance and a shifting temperature gradient on genetic
448 connectivity across a heterogeneous landscape. *BMC Evolutionary Biology* 11, 126.
- 449 50. Karlin, E.F. et al. (2011) One haploid parent contributes 100% of the gene pool for a widespread species
450 in northwest North America. *Molecular Ecology* 20, 753-767
- 451 51. Houde, A.L.S. (2016) Restoring species through reintroductions: strategies for source population
452 selection. *Restoration Ecology* 23, 746–753.
- 453 52. Chanda, P. et al. (2009) Information-theoretic gene-gene and gene-environment interaction analysis of
454 quantitative traits. *BMC Genomics* 10, 509.
- 455 53. Chanda, P. et al. (2008) Ambience: A novel approach and efficient algorithm for identifying informative
456 genetic and environmental associations with complex phenotypes. *Genetics* 180, 1191-1210.
- 457 54. Schlötterer, C. et al. (2015) Combining experimental evolution with next-generation sequencing: a
458 powerful tool to study adaptation from standing genetic variation. *Heredity* 114, 431-440.

459 55. Vandepitte, K. et al. (2014) Rapid genetic adaptation precedes the spread of an exotic plant species.
460 Molecular Ecology 23, 2157-2164.

461 56. Danchin, E. et al. (2011) Beyond DNA: integrating inclusive inheritance into an extended theory of
462 evolution. Nature Reviews Genetics 12, 475-486.

463 57. Bonduriansky, R. (2012) Rethinking heredity, again. Trends in Ecology & Evolution 27, 330-336.

464 58. Jensen, J.D. et al. (2005) Distinguishing between selective sweeps and demography using DNA
465 polymorphism data. Genetics 170, 1401-1410.

466 59. Rollins, L.A. et al. (2016) Selection on mitochondrial variants occurs between and within individuals in an
467 expanding invasion. Molecular Biology and Evolution 33, 995-1007.

468 60. Vuong, H.B. et al. (2017) Influences of host community characteristics on *Borrelia burgdorferi* infection
469 prevalence in blacklegged ticks. PLoS ONE 12 (1), e0167810.

470 61. Nourmohammad, A. et al. (2014) Adaptive evolution of molecular Phenotypes. Curr. Opin. Genet. Dev.
471 23.

472 62. Riedel, N. et al. (2015) Multiple-line inference of selection on quantitative traits. Genetics 201, 305-322.

473 63. Kobayashi, T.J. and Sughiyama, Y. (2015) Fluctuation relations of fitness and information in population
474 dynamics. Physical Review Letters 115, 238102.

475 64. de Vladar, H.P. and Barton, N.H. (2011) The contribution of statistical physics to evolutionary biology.
476 Trends in Ecology and Evolution 26, 424-432.

477 65. Saito, N. et al. (2014) Evolution of genetic redundancy: the relevance of complexity in genotype-
478 phenotype mapping. New Journal of Physics 16, 063013

479 66. Nourmohammad, A. et al. (2013) Evolution of molecular phenotypes under stabilizing selection. Journal
480 of Statistical Mechanics: Theory and Experiment stacks.iop.org/JSTAT/2013/P01012 doi:10.1088/1742-
481 5468/2013/01/P01012 arXiv:1301.3981v1 [q-bio.PE].

482 67. Kosman, E. (2014) Measuring diversity: from individuals to populations. Eur J Plant Pathol 138, 467-486.

483 68. Wegmann, D. et al. (2011) Recombination rates in admixed individuals identified by ancestry-based
484 inference. Nature Genetics 43, 847-853.

485 69. Holleley, C.E. et al. (2014) Testing single-sample estimators of effective population size in genetically
486 structured populations. Conservation Genetics 15, 23-35.

487 70. Wu, C. et al. (2012) Genetic association studies: an information content perspective. Current Genomics
488 13, 566-573.

489 71. Yee, J. et al. (2015) Detecting genetic interactions for quantitative traits using spacing entropy measure.
490 BioMed. Research International 2015, ID-523641.

491 72. Zhang, L. et al. (2009) A multilocus linkage disequilibrium measure based on mutual information theory
492 and its applications. Genetica 137, 355-364.

493 73. Isir, A.B. et al. (2016) An information theoretical study of the epistasis between the *CNR1 1359 G/A*
494 polymorphism and the *Taq1A* and *Taq1B DRD2* polymorphisms: assessing the susceptibility to cannabis
495 addiction in a Turkish population. J. Mol. Neurosci. 58, 456-460.

496 74. Smouse, P. (1974) Likelihood analysis of recombinational disequilibrium in multiple locus gametic
497 frequencies. Genetics 76, 557-565.

498 75. Yoder, J.B. et al. (2014) Genomic signature of adaptation to climate in *Medicago truncatula*. Genetics
499 196, 1263-1275.

500 76. Moore, J.H. et al. (2017) Grid-based stochastic search for hierarchical gene-gene interactions in
501 population-based genetic studies of common human diseases. BioData Mining 10, 19.

502 77. Wang, S. et al. (2017) Integrative information theoretic network analysis for genome-wide association
503 study of aspirin exacerbated respiratory disease in Korean population. BMC Medical Genomics 10(Suppl 1),
504 31.

505 78. von Kodolitsch, Y. et al. (2006) Predicting severity of haemophilia A and B splicing mutations by
506 information analysis. Haemophilia 12, 258-262.

507 79. Wang, K. et al. (2014) Differential Shannon entropy and differential coefficient of variation: alternatives
508 and augmentations to differential expression in the search for disease-related genes. Int. J. Computational
509 Biology and Drug Design 7, 183.

- 510 80. Egizi, A. and Fonseca, D.M. (2015) Ecological limits can obscure expansion history: patterns of genetic
511 diversity in a temperate mosquito in Hawaii. *Biological Invasions* 17, 123-132.
- 512 81. Mandel, J.R. et al. (2016) Patterns of gene flow between crop and wild carrot, *Daucus carota* (Apiaceae)
513 in the United States. *PLoS ONE* 11 (9), e0161971.
- 514 82. Cooke, G.M. et al. (2016) Understanding the spatial scale of genetic connectivity at sea: unique insights
515 from a land fish and a meta-analysis. *PLoS ONE* 11 (5), e0150991.
- 516 83. Leinster, T. and Cobbold, C. (2012) Measuring diversity: the importance of species similarity. *Ecology* 93,
517 477-489.
- 518 84. Chiu, C.-H. et al. (2014) Phylogenetic beta diversity, similarity, and differentiation measures based on Hill
519 numbers. *Ecological Monographs* 84, 21–44.
- 520 85. Beaumont, M. and Nichols, R.A. (1996) Evaluating loci for use in the genetic analysis of population
521 structure. *Proceedings of the Royal Society of London* 263, 1619-1626.
- 522 86. Bergstrom, C.T. and L., L. (2005) The fitness value of information. *arXiv q-bio/0510007v1 [q-bio.PE]*.
- 523 87. Dodig-Crnkovic, G. (2017) Nature as a network of morphological infocomputational processes for
524 cognitive agents. *Eur. Phys. J. Special Topics* 226, 181.
- 525 88. Zeeberg, B. (2002) Shannon information theoretic computation of synonymous codon usage biases in
526 coding regions of human and mouse genomes. *Genome Research* 12, 944-955.
- 527 89. Mashayekhi, M. et al. (2014) A machine learning approach to investigate the reasons behind species
528 extinction. *Ecological Informatics* 20, 58–66.
- 529 90. Chao, A. et al. (2010) Phylogenetic diversity measures based on Hill numbers. *Philosophical Transactions*
530 *of the Royal Society B* 365, 3599-3609.
- 531 91. Scheiner, S.M. et al. (2016) Decomposing functional diversity. *Methods in Ecology and Evolution* DOI:
532 10.1111/2041-210X.12696.
- 533 92. Verity, R. and Nichols, R.A. (2014) What is genetic differentiation, and how should we measure it—GST,
534 D, neither or both? *Molecular Ecology* 23, 4216–4225.
- 535 93. Zhou, Q. and Yu, Y.-M. (2014) Information dimension analysis of bacterial essential and nonessential
536 genes based on chaos game representation. *Journal of Physics D: Applied Physics* 47 (46), 465401.
- 537 94. Ané, C. and Sanderson, M.J. (2005) Missing the forest for the trees: phylogenetic compression and its
538 implications for inferring complex evolutionary histories. *Syst. Biol.* 54, 146-157.
- 539 95. Hadka, D. and Reed, P. (2013) Borg: an auto-adaptive many-objective evolutionary computing
540 framework. *Evolutionary Computation* 21, 231-259.
- 541 96. Sole, R. (2016) Synthetic transitions: towards a new synthesis. *Phil. Trans. Roy. Soc. B.* 371, 2015.0438.
- 542 97. Mather, A.E. et al. (2012) An ecological approach to assessing the epidemiology of antimicrobial
543 resistance in animal and human populations. *Proc. R. Soc. B.* 279, 1630-1639.
- 544 98. Jost, L. et al. (2010) Partitioning diversity for conservation analyses. *Diversity and Distributions* 16, 65–76.
- 545 99. Chiu, C.-H. and Chao, A. (2014) Distance-based functional diversity measures and their decomposition: A
546 framework based on Hill numbers. *PlosONE* 9 (7), e100014.
- 547 100. Chao, A. and Chiu, C.-H. (2016) Bridging the variance and diversity decomposition approaches to beta
548 diversity via similarity and differentiation measures. *Methods in Ecology and Evolution* 7, 919-928.

549

550

551

552

553

554

****** FIGURE 1 CAPTIONS (NB FIGURES SUPPLIED AS .pptx)******

555

Figure 1. Relationship of Mutual Information I to Population Size and Dispersal, from Simulations

556

and Living Populations.

557

(A) Observed Mutual information I per locus from simulated microsatellite data, used to estimate

558

dispersal ($N_e m$) via F_{st} or via I . Root mean square error (RMSE) of $N_e m$ is plotted against

559

dispersal rate, for several different effective population sizes (N_e). RMSE is similar to variance,

560

except it assesses departure of (simulated) observations from the equation's prediction, rather

561

than from the mean of the observations. In every case, $N_e m$ calculated from I had the lower RMSE

562

(from [12]).

563

(B) Comparison of analytical predictions of mutual information I with observed SNP mutual

564

information in *Drosophila* fly dispersal experiments with known $N_e m$. For the predictions, a

565

mutation rate of $\mu=10^{-6}$ was assumed, but using $\mu=10^{-3}$ to 10^{-9} made little difference to the

566

predicted values of I (modified from [36]).

567

568

569

570

*****BEGIN TABLE 1*****

571 **Table 1. Relative Sensitivities, Strengths, and Weaknesses of Each Element of the Diversity Profile**
 572 **From ($q=0$ to 2) (Boxes 1, S1); citations are in text.**

573

Value of q	Interpretation of alpha (between-population) and beta (within-population) measures	Advantages	Disadvantages	Possible fixes for disadvantages
$q=0$	Within-population measure ${}^0H_\alpha$: Number of different types of alleles (or species) (Box 1).	More sensitive to rare alleles than $q=1$ or $q=2$.	Serious sampling problems, because very sensitive to rare alleles.	Sampling problems somewhat alleviated by method in Box S6.2
	Between-population (β): Several measures (Box S1), all related to the number of allelic types that are NOT shared between the localities.	More sensitive to rare alleles than $q=1$ or $q=2$.	Serious sampling problems, because very sensitive to rare alleles.	Sampling problems somewhat alleviated by method in Box S6.2
$q=1$	Within-population ${}^1H_\alpha$: Number of ways the array of different alleles could be set out, given the relative fractions of different alleles, p_i available (Box 1). More formally, higher 1H means that there is reduced certainty about what type to expect when a single allele is randomly sampled.	The most commonly used abundance-sensitive measure in community diversity. Sensitivity to each allele according to its frequency, <i>ie.</i> : each copy of each allele is treated equally.	Some sampling problems	Sampling problems fixed by method in Box S6.2
	Between-population (β): Several measures (BoxS1), all related to whether knowing the allelic type of ONE single individual will help identify the location from which it was sampled (eg, if there is no differentiation, then the allelic information is completely uninformative, but if there is complete	Shannon differentiation satisfies 'monotonicity' (Some other transformations for ($q=1$) do not). Shannon differentiation (Box 2) satisfies 'true dissimilarity', which means that the differentiation	Some sampling problems	Sampling problems fixed by method in Box S6.2

	differentiation, then knowing the allelic type would give certain identification of the location of origin)	measure should represent the actual proportion of non-overlapping alleles, when populations are equally diverse and all alleles have the same frequencies. Some other transformations for ($q=1$) do not satisfy this. Shannon information provides a natural measure of evolvability.		
$q=2$	Within-population $^2H_\alpha$: Chance of choosing TWO different types, given the relative fractions of different types, p_i available (Box 1).	Very sensitive to frequent types. Relatively few sampling problems	Insensitive to rare types that may be important in conservation and long-term evolution	
$q=2$	Between-population (β): Many measures (Box S1), all related to whether TWO individuals sampled from different localities are likely to have different allelic types (eg if there is no differentiation, they MUST be the same allelic type, but if there is complete differentiation, they MUST be different types).	Very sensitive to frequent types. Relatively few sampling problems <i>Jost-D</i> (Box S1) satisfies 'true dissimilarity', which means that the differentiation measure should represent the actual proportion of non-overlapping alleles, when populations are equally diverse and all alleles have the same frequencies. Some other transformations for ($q=2$) do not satisfy this.	Insensitive to rare types that may be important in conservation and long-term evolution Most ($q=2$) measures confound among-population (beta) diversity with one or both of within-population (alpha) diversity and total (gamma) diversity. These are cases of the 'dependence-on-alpha problem' and the 'dependence-on-gamma problem'. Most ($q=2$) measures lack monotonicity, ie apparent difference between two	<i>Jost-D</i> (Box S1) fixes the dependence-on-alpha problem. A measure to fix the dependence-on-gamma problem is cited in the main text.

			<p>populations can decrease when a new unshared allele appears in a population, wrongly implying that this reduces the pace of speciation. The $q = 2$ differentiation measure <i>Jost-D</i> (Box S1), does not possess the expected monotonicity. F_{ST} (Box S1), which is also sometimes used as a beta differentiation measure, does not satisfy monotonicity.</p>	
--	--	--	--	--

574

575

576

577

*****START BOX 1*****

578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599

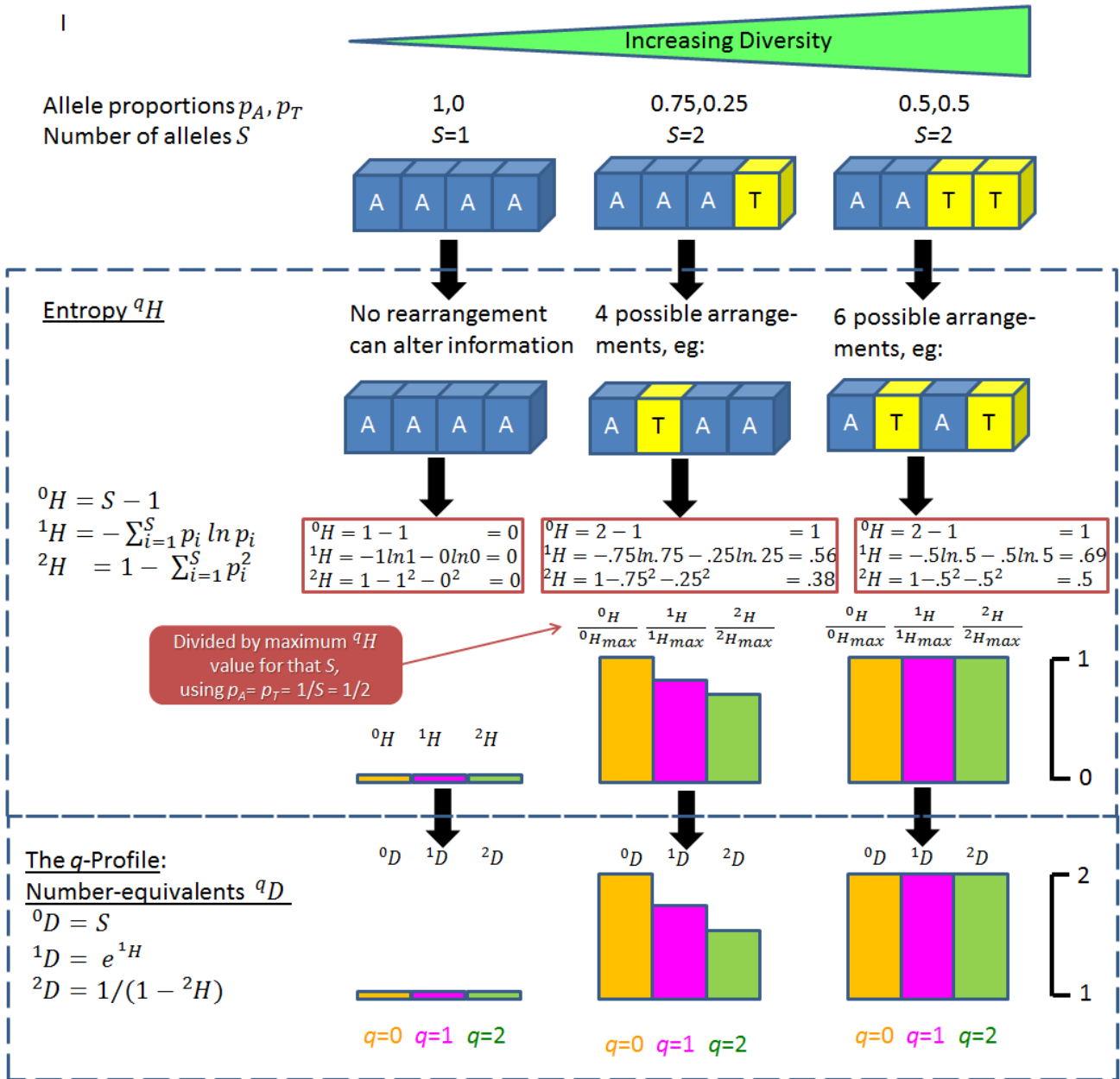
Box 1 Tutorial on the q-Profile: ‘Effective-Number’ Diversities 0D , 1D , 2D and their derivation from the corresponding Entropies 0H , 1H , 2H .

Figure 1 shows three samples of four haploid individuals, each genotyped to identify SNP allele A or T. The first row of histogram bars shows entropy values qH , including heterozygosity 2H , and 1H which can be derived from the number of possible novel arrangements of the alleles carried by the four sampled individuals, as if one was trying to spell out words with the alleles. In the left sample, which has no diversity, all qH measures are zero (NB conventionally, $0 \ln 0 = 0$). In the right sample, where alleles are equally-frequent, each measure is at its maximum possible value with two alleles. The second row of histogram bars shows the q-profile of ‘effective-number’ diversities qD derived from the qH entropies. Note that in moving from the left sample to the middle sample, we are adding a rare allele, a single copy of a new allele T. In this case, 0H and 0D show the greatest response, while heterozygosity 2H and its transform 2D show the smallest response, because in Equation 2 when a rare allele’s proportion is squared, its effect becomes much smaller or even negligible. In contrast, when both alleles are already present, 0H and 0D show no response to changing the numbers of copies of each allele, while 2H and 2D show the greatest response. In both cases 1H and 1D show an intermediate response. All formulas are in Box S1.

Box 1 continues on next page...

600 Box 1 continued.

601 Figure I. Calculating qH and qD



602

603

604

605 Box 1 continues on next page...

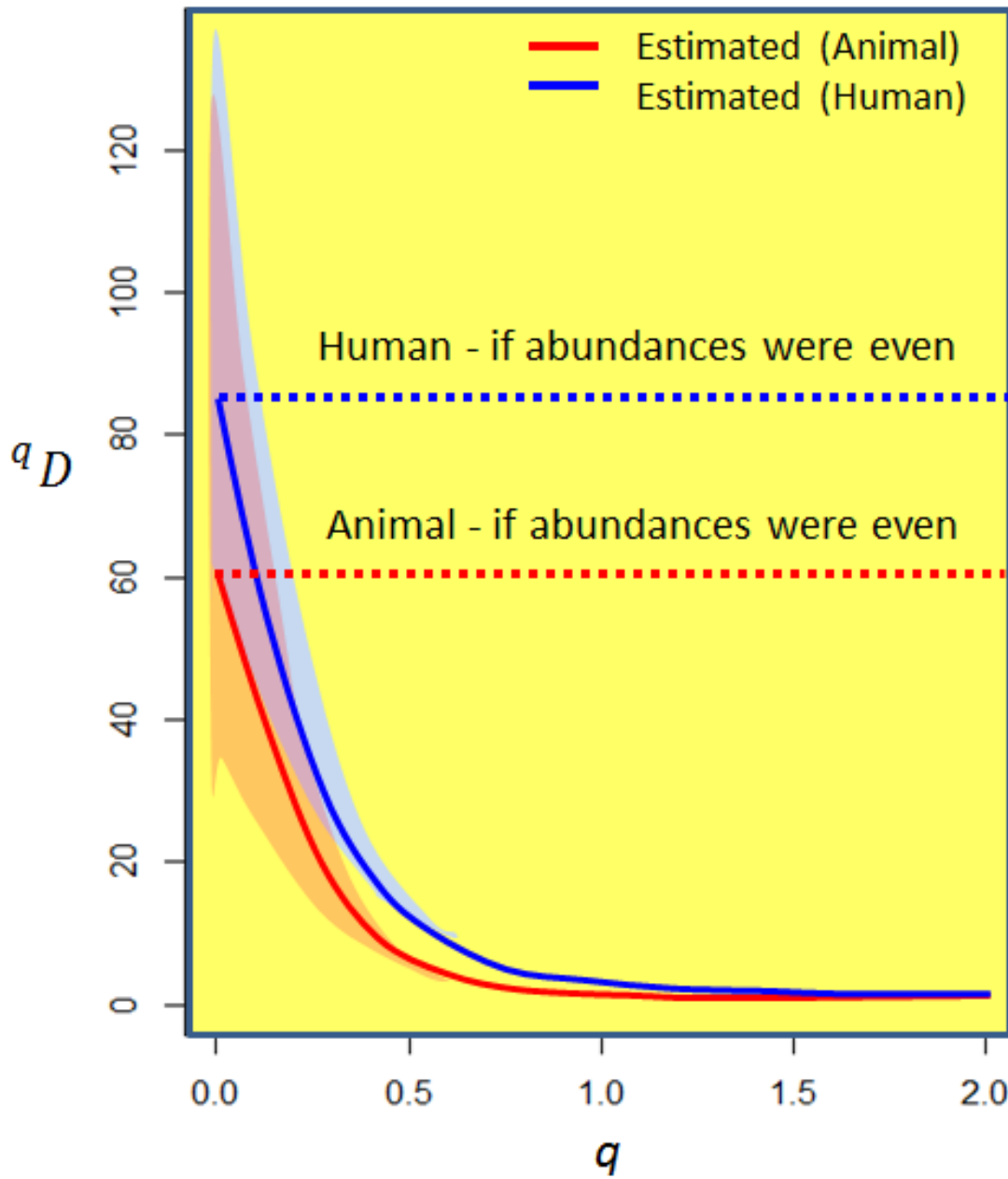
606 Box 1 continued.

607 Figure II shows diversity profiles of antimicrobial resistance variants within *Salmonella typhimurium*
608 in humans (blue) and domestic animals (red) [97], with correction for sampling bias (Box S6.2). The
609 horizontal axis is the value of q , and the vertical axis is qD for each value of q . The shaded areas
610 are 95% confidence limits. Note that these limits overlap, except in the region ($0.3 < q < 1.2$) where
611 the plots can be discriminated. This highlights the difficulty of comparing profiles if only ($q=0$) and
612 ($q=2$) are used. Also note the effect of evenness of observed allelic distributions upon diversity
613 profile. If the resistance types had been equally frequent within each host ($p_1 = p_2 =$
614 $p_3 = \dots = p_S$) then the qD profiles would have been flat (dotted lines). However, in each group
615 (Animal, Human), there is one highly abundant type and a many types only recorded once.
616 Therefore, the profiles drop very sharply and become flat for order ($q=2$) (eg. Heterozygosity or
617 Gini-Simpson), because for this q -value, qD is mainly determined by the abundant type(s), being
618 insensitive to rare types.

619 Box 1 continues on next page...

620 Box 1 continued.

621 Figure II q -Profiles



622

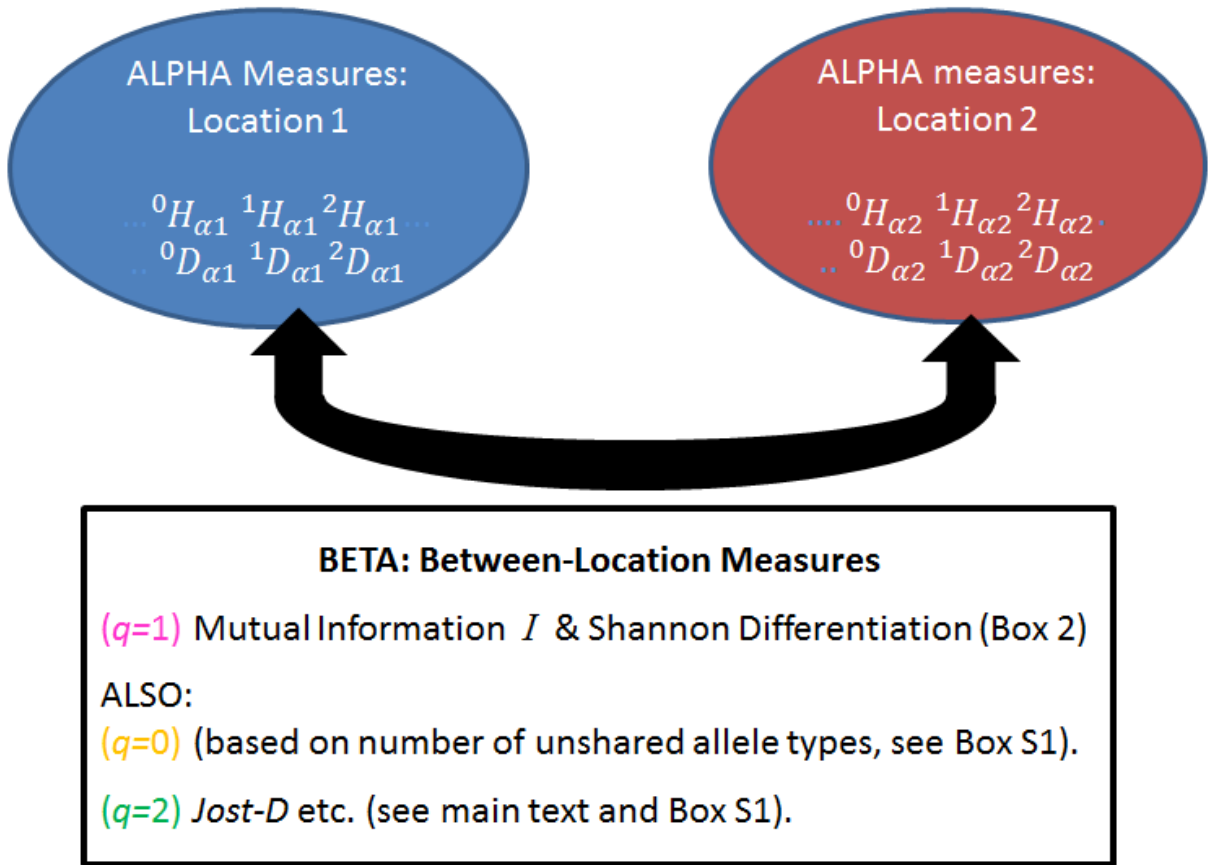
623

624 Box 1 continues on next page...

625 Box 1 continued.

626 **Figure III Alpha and Beta Measures of Entropy and Diversity**

627



628

629

630

631

632

*****END BOX 1*****

*****START BOX 2*****

Box 2 Information Measures of Geographic Differentiation: Mutual

Information I.

‘Mutual information’ I is a ($q=1$) measure that has a similar role to correlation in the ($q=2$) scale; it expresses association between two variables, such as allelic differentiation and population membership ([12], where I is called $^S H_{UA}$). Therefore differentiation among populations can be measured as mutual information I between geographic location and allelic differentiation. For example, we ask: “Does knowing the alleles in a sample, give any information about which location was sampled?” This is not true in Figure I, where the allele proportions p_C and p_T are the same in the two locations, so there is zero mutual information between the variables ‘location’ and ‘allele type’ – in other words allelic data provide no information about population of origin. However in Figure II, mutual information is maximal: knowing the alleles in a sample gives perfect information about which location was sampled, because there are no alleles that are shared between the two locations.

Figure I. Zero Mutual Information I .

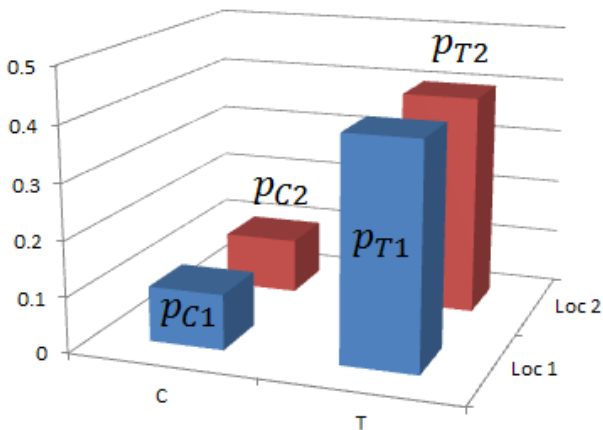
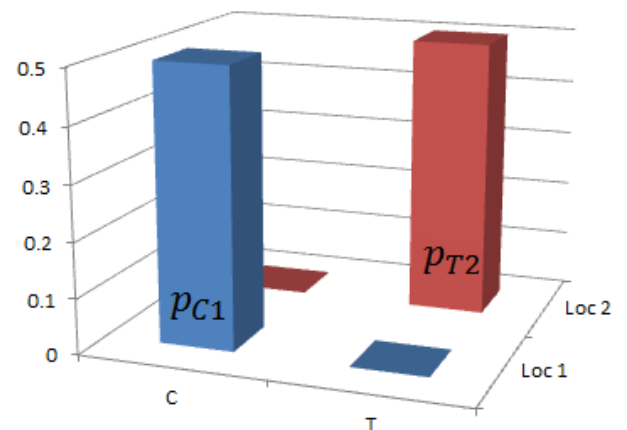


Figure II. Maximum Mutual Information I .



648

649 Box 1 continues on next page...

650 Box 1 continued.

651 Figure III outlines the calculation of mutual information I from observed data for two locations (eg,
 652 alpha-level variation within estuary 1, and within estuary 2) as well as total for both locations
 653 (gamma-level proportions, marginals p_C and p_T). Mutual information can be derived very easily
 654 from the chisquare for allelic differentiation between populations, using $I = \chi^2/2n$, (where n is
 655 the total sample size, I is the mutual information, and χ^2 is log-linear-chisquare, with expectations
 656 for each cell calculated as shown for one example in Figure III [48]).

657 Alternatively, I is calculated as the part of the total information ${}^1H_\gamma$ that is not due to variability
 658 within single locations ${}^1H_\alpha$; using terms from Figure III:

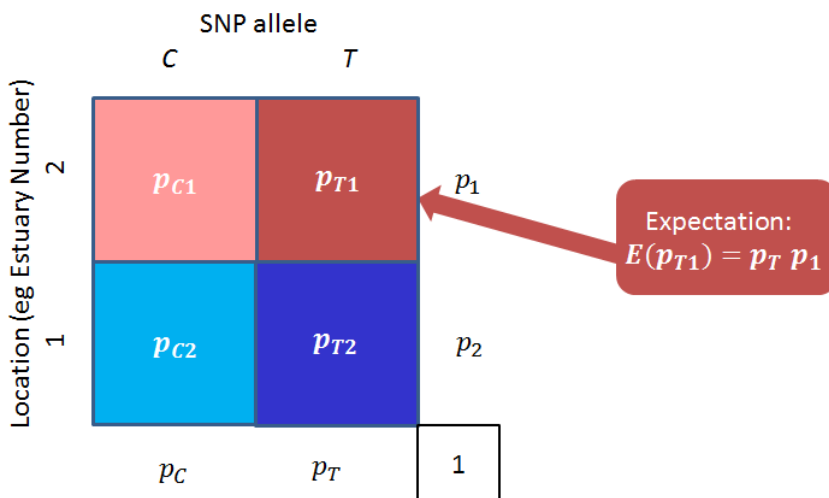
659
$$I = ({}^1H_\gamma - \overline{{}^1H_\alpha})$$

660
$$= (-p_C \ln p_C - p_T \ln p_T) - 0.5 \left(-\frac{p_{C1}}{p_1} \ln \frac{p_{C1}}{p_1} - \frac{p_{T1}}{p_1} \ln \frac{p_{T1}}{p_1} \right) - 0.5 \left(-\frac{p_{C2}}{p_2} \ln \frac{p_{C2}}{p_2} - \frac{p_{T2}}{p_2} \ln \frac{p_{T2}}{p_2} \right)$$
 (Box S1).

661 Also, mutual information adjusted to range from zero to unity is:

662 Shannon Differentiation = $I/\ln K$ (where K is the number of equal-sized populations; Box S1 shows
 663 the formula for other cases). Shannon differentiation has useful properties, discussed in Box 4.

664 Figure III. Data for Calculating Mutual Information between Allele-type and Location.



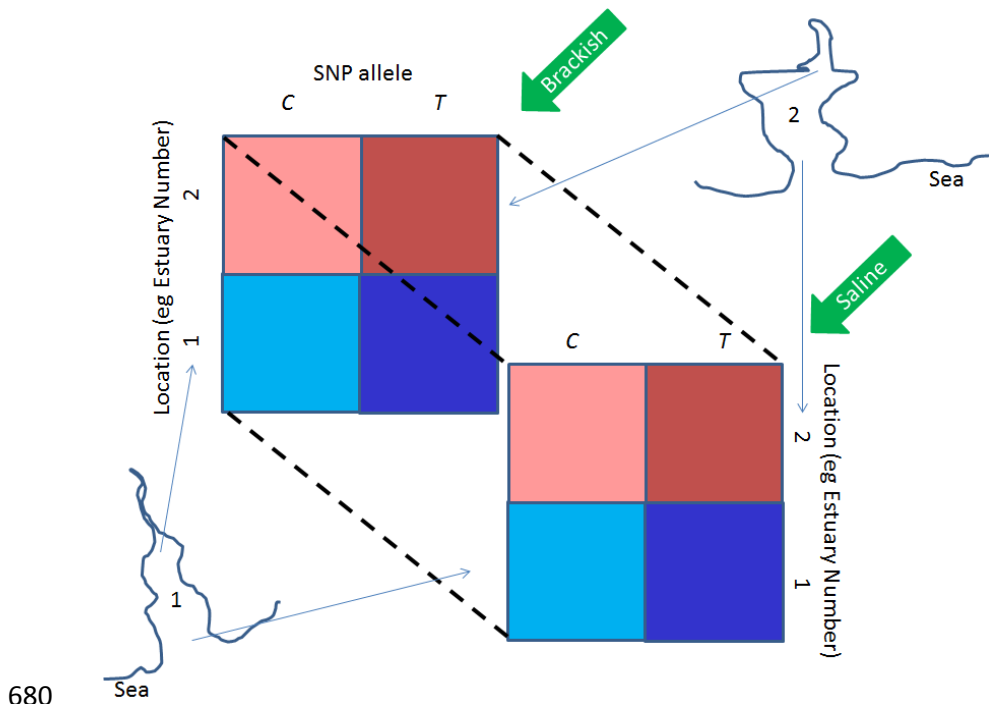
665

666 Box 2 continues on next page....

667 Box 2 continued....

668 Figure IV shows expansion to include two adjacent habitats (brackish and saline) sampled within each
669 estuary. Each estuary has a saline habitat near the mouth, whose data from Figure III are shown in the
670 foreground of the cube, and a similar set of data for a brackish habitat further inland are in the
671 background of the cube [35]. This expansion for multiple variables is standard for contingency-
672 analysis [48], so the equivalent of mutual information can easily be calculated for this situation [35].
673 In fact, indefinite expansion is possible, to accommodate multiple alleles per locus, plus dimensions
674 added to incorporate diversity within and among different habitats, landscapes, etc. This is possible
675 because log-linear χ^2 (and therefore I) are completely additive [35, 48]. The partition strategy will
676 depend upon the hypothesis being tested [45, 46, 48]. Programs, sampling bias corrections, and
677 example calculations are in Box S6.

678 **Figure IV. Partitioning Molecular Information with Two Variables (Location and**
679 **Habitat).**



680

681 Box 2 continues on next page....

682 Box 2 contined....

683 With equal population sizes, I is the ecologists' Horn measure [37]. Mutual information is closely
684 related to the relative entropy (Kullback-Leibler) which is used for tests for Hardy-Weinberg
685 equilibrium (Box S4) and for comparing allele frequencies before vs. after selection, as well as the
686 rate of change of Fisher information [92].

687 *******END BOX 2*******

688 *****BEGIN BOX 3 *****

689 **Box 3 Choosing a Measure for Within-population Diversity (alpha):**

690 **the Relative Merits of Information-based Measures ($q=1$) and Other Measures.**

691 Elements of the diversity profile from ($q=0$) to ($q=2$), have various strengths and weaknesses, each
692 being sensitive to particular aspects of diversity relevant to different questions (Boxes 1,S1).

693 Measures based on counts of different types of alleles (or species, $q=0$) are very sensitive to rare

694 alleles (Box 1), which is obviously important if there is a focus on rarity, either for conservation

695 reasons, or because novel adaptive mutants are initially rare [4, 5, 59]. On the other hand, the

696 ($q=2$) measures such as heterozygosity (${}^2H_\alpha$ Equation 2 in main text, also Box 1) give very little

697 weight to rare alleles, because they are based on the chance of choosing two different types, given

698 the proportions p_i available in the population. Estimating heterozygosity involves squaring

699 p_i values, so that values close to zero (*ie.* rare alleles) become vanishingly small, relative to more

700 common alleles. The intermediate value of q ($q=1$, ${}^1H_\alpha$ Shannon information, Equation 1 in main

701 text, also Box 1) assesses diversity as the number of ways the array of different alleles could be set

702 out, given the p_i values. Shannon weights each allele by its frequency, so its response to addition of

703 a single novel allele is intermediate between those of ($q=2$) and ($q=0$) measures (Box 1-I).

704

705 These sensitivities to rare alleles result in different sampling properties. Counts of different types

706 ($q=0$) are changed considerably by either adding a single individual of a different type, or else failing

707 to sample this single individual. As a result, even after sampling corrections, these measures often

708 cannot distinguish diversity levels of assemblages that can be discriminated by either ($q=1$) or ($q=2$)

709 scales (Box 1-II in main text). In contrast, ($q=2$) methods such as heterozygosity/Simpson's have

710 relatively few sampling problems. Sampling for Shannon information ($q=1$) can be well addressed by
711 modern corrections (Box S6.2), so that it can distinguish between alternative assemblages (Box1-II).

712

713 Frequency of use is important for comparability between studies. In community diversity, Shannon
714 is the most commonly used abundance-sensitive measure ($q = 1 \ ^1H_\alpha$ [11]), whereas heterozygosity
715 ($q=2$) is most commonly used in molecular ecology [4, 5]. As the two fields gradually unify [13, 21],
716 it will become important to use measures that are common to both fields, and we suggest that a
717 profile of $q = 0, 1, 2$ will achieve this best.

718

719 Finally, there is extensive literature on neutral and adaptive processes in molecular ecology for both
720 ($q=0$) and ($q=2$) alpha-measures [4, 5]. More recently, it has been proposed that the ($q=1$) measure
721 Shannon information $^1H_\alpha$ provides a natural measure of evolvability [38, 39], and Shannon
722 predictive methods are being developed for both neutral [12, 36, 37] and adaptive variants [35, 61-
723 64, 66].

724 *******END BOX 3*******

725

726

727

*******BEGIN BOX 4*******

728

729

Box 4 Choosing an Among-Population Differentiation Measure (Diversity).

730

Among-population measures ($q=0,1,2$) inherit virtues and shortfalls of corresponding alpha

731

measures (Box 3), plus properties specific to the beta level. For three reasons, the ($q=1$) measures

732

might be the best tools to track and understand evolutionary processes of divergence.

733

734

First, many measures used as beta-differentiation measures actually confound this with within-

735

population (alpha) or total (gamma) diversity, *e.g.* F_{ST} and G_{ST} ($q=2$) tend towards zero as diversity

736

increases [98, 99] (Figure I). In contrast, Shannon differentiation ($q = 1$, Box 2) and *Jost-D* ($q=2$, Box

737

S1) are zero only when allele frequencies are identical across localities, and unity only when there

738

are no shared alleles. Chao and Chiu [100] proposed a measure to avoid dependence-on-gamma.

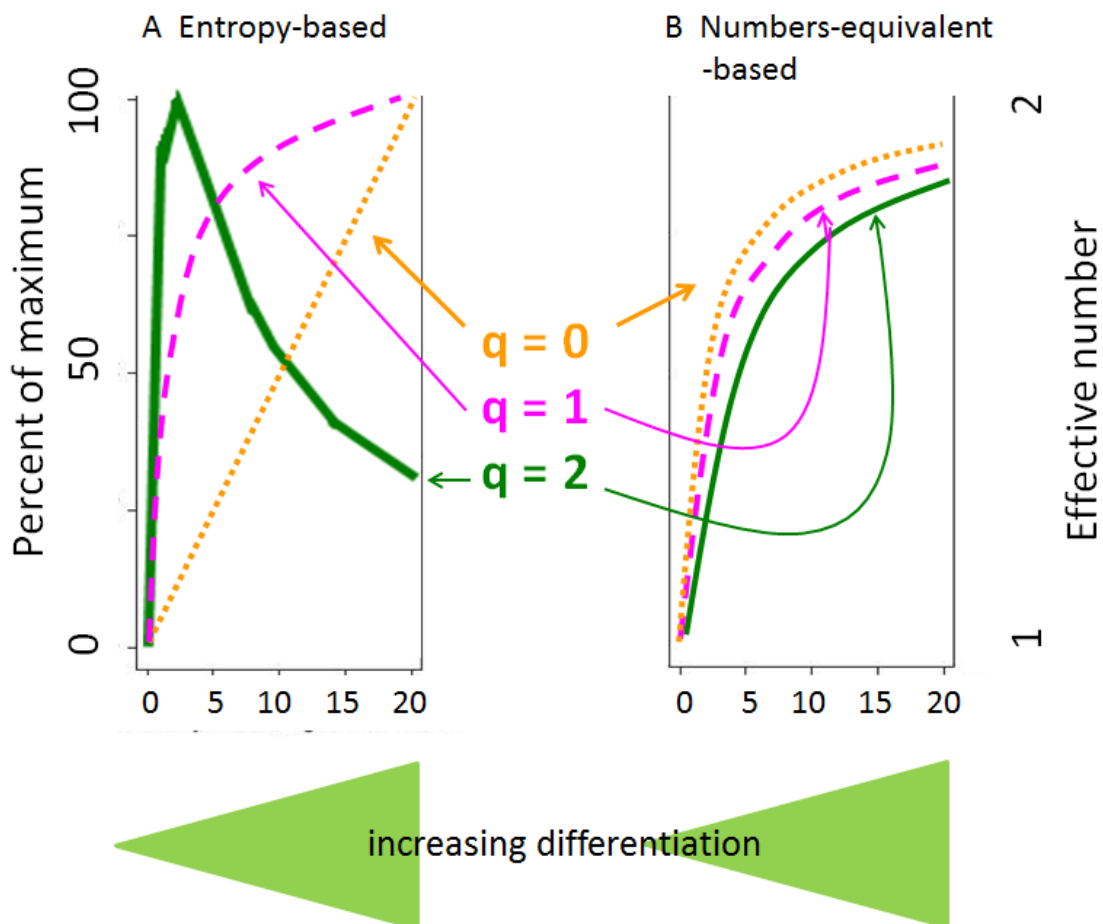
739

Box 4 continues on next page....

740 Box 4 continued.

741 **Figure I Effect of increasing allelic differentiation upon differentiation measures**
742 **for ($q=0,1,2$).**

743 Initially two populations shared the same two alleles, and the horizontal axis shows the addition of
744 extra unshared alleles to each locality. Equations and symbols are as shown in Box S1, except the q
745 $=0$ entropy-based measure is $S_Y - (S_{\alpha_1} + S_{\alpha_2})/2$ and the ($q=2$) entropy-based measure is
746 $({}^2H_Y - \overline{{}^2H_\alpha}) = G_{ST} * {}^2H_Y$. In panel A, the vertical axis for each plot is adjusted to percent of
747 the maximum observed differentiation value for that order of q (modified from [98]).



748

749 Box 4 continues on next page....

750 Box 4 continued.

751 Second, to track evolving divergence among localities, a differentiation measure should never
752 decrease when shared alleles are replaced or augmented by new unshared alleles, called
753 'monotonicity'. Surprisingly, most differentiation measures derived from diversity partitioning do
754 not meet this basic criterion, wrongly implying that new unshared alleles can reduce the pace of
755 differentiation. For example, the additive ($q=2$) differentiation (Figure 1, green solid line, panel A)
756 increases then decreases as differentiation increases, as do most ($q=2$) measures (e.g., F_{ST} and
757 *Jost-D* Box S1), and even some ($q=1$) measures. The only differentiation measures with strong
758 monotonicity are ($q=1$) Shannon differentiation, and those based on ($q=0$), though the latter are
759 less desirable because of their sampling problems.

760

761 A third important property of differentiation measures is 'true dissimilarity', which means that the
762 differentiation measure should represent the actual proportion of non-overlapping alleles, when
763 populations are equally diverse and all alleles have the same frequencies. This is untrue for many
764 ($q=1$ or 2) measures, but is true for the ($q=1$) measure Shannon differentiation (Box 2) and for ($q=2$)
765 differentiation *Jost-D* [29].

766 *******END BOX 4*******

767

768

769

770

771 *****BEGIN BOX 5 OUTSTANDING QUESTIONS *****

772 **Box 5 Outstanding Questions**

773 Regular use of a formal qD profile of at least ($q = 0, 1, 2$) measures at each level ($\alpha, \beta,$ and γ) will
774 achieve maximum understanding of patterns, and will later allow meta-analysis of the performance
775 of the molecular q -profile under a wide variety of conditions.

776

777 Neutral forecasts need to be extended to complex scenarios, such as asymmetric dispersal, unequal
778 population sizes, and bottlenecks including possible separate effects of actual and effective
779 population size.

780

781 We need an analogue to AMOVA based on information-theoretic methods.

782

783 There is a need for new predictive theory for variants with mutation not described by models such
784 as infinite (IAM) stepwise (SMM) and biallelic SNPs.

785

786 Further analyses of adaptation and selection with ($q=1$) approaches will profit from three attributes
787 of Shannon: similarity to logit (log-linear) modeling, sensitivity to rare novel variants that are crucial
788 in adaptation, and independence of α and β .

789

790 Analyses of linkage disequilibrium and expression, which already use measures related to mutual
791 information, might benefit from using its transform to Shannon differentiation.

792

793 Box 5 Outstanding Questions continues on next page....

794

795 Box 5 Outstanding Questions continued.

796 Integrating all biological levels from community ecology and evolution, through to sub-cellular

797 biology, can capitalise on existing protocols for using information and entropy methods as general

798 forecasting methods. This will include incorporating information on genetic and/or functional

799 similarity or difference of alleles (or species).

800

801 Information and entropy theory is very broad, with an abundance of other possible connections

802 within and outside biology.

803 *******END BOX 5 OUTSTANDING QUESTIONS*******

*******BEGIN GLOSSARY BOX*******

804

805 **Glossary**

806 **Adaptation:** This refers to the evolutionary process due to natural selection for organisms that are
807 better at surviving and/or reproducing in a particular environment.

808

809 **Additivity:** This refers to a multidimensional table (e.g., dimension 1 is allele type, dimension 2 is
810 location, dimension 3 is an environmental variable, etc). In this type of table, measures such as log-
811 linear contingency chisquare (i.e., mutual information) can be partitioned into completely additive
812 sub-investigations, unlike Pearson's chisquare.

813

814 **Allele proportions and frequencies:** For conformity to conventions in information and entropy
815 theory, statistics, and all other science (except population genetics!), we refer to 'allele-
816 frequencies' when dealing with counts ranging from zero to infinity, and to 'allele-proportions',
817 when these frequencies have been converted to p_i ranging from zero to unity.

818

819 **Alpha, beta, and gamma diversity measures (α , β , γ):** These indicate within-locality (alpha)
820 diversity, among-locality differentiation (beta) diversity, and total (gamma) diversity. In other
821 publications, alpha values are often assumed to be averaged over many locations. Where an
822 average is made, we will indicate this by an overbar, and describe any unequal weighting.

823

824 **Bottleneck:** A period of reduced population size, which usually will alter entropy and diversity
825 levels away from the previous expectations.

826 Glossary Box continues on next page....

827 Glossary Box continued

828 **Drift:** Random genetic drift is caused by the chance nature of transmission of alleles within each
829 family. In a finite population, this chance in transmission results in fluctuations of allele proportions
830 in the entire population. Drift erodes genetic variation summarised by any measure ($q=1,2,3$).

831

832 **Effective numbers (or 'true diversities') qD :** These are conversions of entropies (qH) into qD the
833 number of equally-frequent alleles (or species) that would give the actual observed qH , derived
834 from a possibly uneven array of alleles in the observed sample.

835

836 **Entropy (qH):** Entropies include Heterozygosity or Gini-Simpson (${}^2H = H_e$), Shannon (1H), and the
837 number of allelic types minus one ($S-1$, or 0H). The 'H' symbol is used for entropy throughout
838 science.

839

840 **Epigenetics:** This is one type of non-genetic inheritance, due to chemical modification of DNA
841 sequence, e.g., methylation, which is sometimes transmissible through at least one generation, and
842 may have phenotypic effects. We do not use the term 'epigenetics' to cover all types of non-
843 genetic inheritance. Frequently confused with epistasis.

844

845 **Epistasis:** Functional interactions between multiple loci, at any level from transcription onwards, to
846 affect a measured phenotype. Frequently confused with epigenetics.

847

848 **Haplotype:** A collective genetic combination located on a single molecule of DNA, typically
849 including two or more SNPs showing variation of bases between alternative haplotypes.

850 Glossary Box continues on next page....

851 Glossary Box continued

852 **Hardy-Weinberg-Equilibrium (HWE):** For a single locus, HWE is the condition where the
853 combinations of alleles in diploid genotypes are as expected from random combination of the
854 population's pool of alleles. HWE can be disrupted by mutation, selection, dispersal, non-random
855 mating (including inbreeding), or random genetic drift.

856

857 **IAM Infinite allele model of mutation:** This is suitable for long genomic regions with rare base
858 substitutions, so that each mutation creates an allele that has never existed before. This is
859 approximately suitable for protein-coding regions, and much other non-repetitive DNA.

860

861 **Information theory:** This has multiple strands, but we focus on Shannon's Information Index 1H ,
862 which is a measure of the number of different ways that a group of objects (e.g., individuals of
863 different species, or DNA molecules with different sequences) can be rearranged. The information
864 index is also a measure of 'surprise' or uncertainty, increasing in groups where we are less certain
865 of what we would sample at random.

866

867 **Linkage disequilibrium (LD):** This is when the combinations of alleles at two or more diploid loci
868 depart from those expected by random sampling from the population. The loci might be carried on
869 the same molecule of DNA (true LD), or on different molecules (sometimes called genotypic
870 disequilibrium or GD). LD results from random genetic drift, mutation, selection, dispersal, non-
871 random mating (including inbreeding), clonal or asexual reproduction.

872

873 Glossary Box continues on next page....

874 Glossary Box continued

875 **Monotonicity:** This means that increases of variable 'x' are EITHER always associated with an
876 increase of 'y', OR always associated with a decrease of 'y'. 'Weak monotonicity' allows plateaus of
877 'y'.

878

879 **Neutral:** This refers to genetic variants that do not affect fitness (survival and/or reproduction), so
880 are not involved in natural selection and adaptation.

881

882 **Non-genetic inheritance:** Any type of inheritance that does not rely upon variation of the DNA
883 sequence of A,T,C,G. Examples include microbiome inheritance, epigenetic inheritance, niche
884 inheritance, etc. All of these can have profound fitness effects.

885

886 **Replication:** This means that a diversity measure increases linearly when equally diverse and
887 completely distinct groups are pooled in equal proportions.

888

889 **Selection:** This is based on individuals that possess heritable characteristics (eg alleles) that give
890 them higher relative survival and/or reproduction in a particular environment. As a result, there
891 will be increased representation of those heritable characteristics in the next generation. Such
892 individuals are said to have higher fitness in that environment.

893

894 **SMM Stepwise mutation model:** An approximation of microsatellite evolution.

895

896 **SNP Single Nucleotide Polymorphism:** A single-basepair locus that varies in the population.

897 *******END GLOSSARY BOX*******