

2 **Message in a Bottle – Metabarcoding Enables Biodiversity Comparisons Across**
3 **Ecoregions**

4

5

6 Steinke D^{1,2*}, deWaard SL¹, Sones, JE¹, Ivanova NV^{1,2}, Prosser SWJ¹, Perez K¹,
7 Braukmann TWA¹, Milton M¹, Zakharov EV^{1,2}, deWaard JR^{1,3}, Ratnasingham S^{1,2}
8 Hebert PDN^{1,2}

9

10 Affiliations:

11 ¹Centre for Biodiversity Genomics, University of Guelph, 50 Stone Road East, Guelph,
12 Ontario, N1G 2W1, Canada

13 ²Department of Integrative Biology, University of Guelph, 50 Stone Road East, Guelph,
14 Ontario, N1G 2W1, Canada

15 ³School of Environmental Sciences, University of Guelph, 50 Stone Road East, Guelph,
16 Ontario, N1G 2W1, Canada

17

18

19

20

21

22 *Corresponding author: Dirk Steinke (dsteinke@uoguelph.ca)

23

24

25 **Abstract**

26 Background

27 Traditional biomonitoring approaches have delivered a basic understanding of
28 biodiversity, but they cannot support the large-scale assessments required to manage and
29 protect entire ecosystems. This study employed DNA metabarcoding to assess spatial and
30 temporal variation in species richness and diversity in arthropod communities from 52
31 protected areas spanning three Canadian ecoregions.

32 Results

33 This study revealed the presence of 26,263 arthropod species in the three ecoregions and
34 indicated that at least another 3,000–5,000 await detection. Results further demonstrate
35 that communities are more similar within than between ecoregions, even after controlling
36 for geographical distance. Overall α -diversity declined from east to west, reflecting a
37 gradient in habitat disturbance. Shifts in species composition were high at every site with
38 turnover greater than nestedness, suggesting the presence of many transient species.

39 Conclusions

40 Differences in species composition among their arthropod communities confirm that
41 ecoregions are a useful synoptic for biogeographic patterns and for structuring
42 conservation efforts. The present results also demonstrate that metabarcoding enables
43 large-scale monitoring of shifts in species composition, making it possible to move
44 beyond the biomass measurements that have been the key metric employed in prior
45 efforts to track change in arthropod communities.

46

47

48

49

50

51

52

53

54

55

56 **Background**

57 Terrestrial organisms are exposed to diverse anthropogenic stressors, including
58 climate change, resource extraction, and agriculture. Habitat degradation, pesticide usage,
59 invasive species, and associated shifts in food webs have provoked major reductions in
60 the abundance of terrestrial arthropods [1-4]. These declines have led to calls for more
61 comprehensive biosurveillance to inform environmental management and conservation.
62 Long-term monitoring of species composition is essential to quantify biological change,
63 but efforts employing morphological diagnostics have targeted a small set of indicator
64 species because of the need for taxonomic experts for each group. As a consequence, they
65 cannot support the broad assessments needed to manage and protect ecosystems, let alone
66 forecast human impacts on them by integrating statistical modelling. The latter methods
67 demand comprehensive data on species distributions and abundance [5], information that
68 is currently unavailable because of the prior focus on selected biotic compartments at
69 limited geographic scale.

70

71 Two methodological advances promise to meet the need for comprehensive
72 biodiversity data. Firstly, identification systems based on the analysis of sequence
73 variation in short, standardized gene regions (i.e., DNA barcodes) enable species
74 discrimination [6]. Secondly, high-throughput sequencers (HTS) permit the inexpensive
75 acquisition of millions of DNA barcode records [7]. These advances now enable
76 biodiversity surveys at speeds and scales that were previously inconceivable. In
77 particular, the coupling of HTS with DNA barcoding, known as metabarcoding [8], has a
78 compelling advantage over traditional approaches for tracking shifts in species presence.
79 It can generate georeferenced occurrence data from bulk samples at low cost, and a single
80 instrument can process hundreds of bulk samples each week. Because the sequencing
81 output of HTS is doubling every nine months [9,10], analytical costs are certain to
82 sharply decline, allowing production to soar. This augmented capacity for data generation
83 has already enabled large-scale biotic surveys of aquatic and terrestrial arthropods [11-
84 14], vertebrates [15], pollen [16], diatoms [17], and fungi [18-20].

85

86 Access to large collections of specimens is essential to capitalize on the analytical
87 capacity provided by DNA metabarcoding. Among the many approaches used to sample
88 terrestrial arthropods, Malaise traps [21] have gained wide adoption because they collect
89 large, diverse samples with little effort [22]. Although most-effective for sampling flying
90 insects, they also collect ground-active arthropods. By coupling DNA barcoding with
91 Malaise trapping [23,24], high-resolution monitoring networks for arthropods are within
92 reach, but there are challenges. Data interpretation requires a well-parameterized DNA
93 barcode reference library for the region under investigation, creating the need for a
94 system to aid site selection. Ecoregions represent an obvious candidate [25-28] although
95 their boundaries are rarely sharply defined, and they are based on distributional data for a
96 narrow range of taxa. Despite these limitations, ecoregions have been widely used to
97 guide management decisions and to explore species and community diversity patterns
98 [29,30]. As a result, they are a good candidate to serve as the backbone for a large-scale
99 monitoring network. The most widely adopted schema partitions the world's 14 terrestrial
100 biomes into 846 ecoregions [30].

101

102 This study demonstrates the feasibility of employing metabarcoding for large-scale
103 bio-surveillance by comparing the temporal and spatial patterning of arthropod
104 communities in three of Canada's 47 terrestrial ecoregions: the Eastern Canadian Forest –
105 Boreal Transition (ECF – 75,000 km²), the Eastern Great Lakes Lowland Forests (EGL –
106 63,000 km²), and the Southern Great Lakes Forests (SGL – 22,000 km²) (Figure 1).
107 Forest cover declines from 77.7% in the ECF to 30.1% in the EGL and just 12.1% in the
108 SGL while cropland/pastures cover 78% of the SGL, 57% of the EGL, and 3% of the
109 ECF [31]. The EGL and SGL are the most populated ecoregions in Ontario with
110 developed land (e.g., urban, road networks) encompassing more than 7% of the SGL
111 [31]. As such, these ecoregions provide a good basis for assessing the impacts of varied
112 disturbance regimes on biodiversity.

113

114 **Data Description**

115 Collections were made by deploying a Malaise trap at 52 sites in these three
116 ecoregions and samples were metabarcoded to examine variation in their species

117 richness, community composition, phylogenetic diversity, as well as alpha (α) and beta
118 (β)-diversity. Malaise traps were deployed for 20 weeks at 15 sites in the ECF, 24 sites in
119 the EGL, and 13 sites in the SGL. Catches were harvested at two-week intervals and 410
120 of the resultant 520 samples were designated for metabarcoding (the others were reserved
121 for single specimen barcoding). Analysis began with non-destructive lysis of the
122 specimens in each bi-weekly sample, followed by DNA extraction using a membrane-
123 based protocol [32]. A 463 bp amplicon of cytochrome *c* oxidase I (COI) was then PCR
124 amplified and the amplicon pools from each set of 10 samples were sequenced on an Ion
125 Torrent S5 using a 530 chip. The sequences were subsequently analyzed using the
126 Multiplex Barcode Research And Visualization Environment (mBRAVE – mbrave.net).
127 All raw HTS datasets were deposited in the Sequence Read Archive (SRA –
128 www.ncbi.nlm.nih.gov/sra/) under the BioProject accession number PRJNA629553.

129

130 **Analyses**

131 Sequence analysis of the 410 samples produced 367,823,207 reads across 41 S5 runs
132 (mean reads per run = 8.97 million, see **Table S1**). Two thirds were filtered, leaving
133 126,253,260 reads that could be assigned to a BIN (Barcode Index Number; [33]) on
134 BOLD [34] (**Figure S1**). Nearly all reads (99.3 %) found a BIN match on BOLD, but
135 those that failed were *de novo* clustered using mBRAVE with a 99% similarity threshold.
136 The latter analysis recognized an average of 28 additional OTUs per sample, but >96% of
137 them reflected sequencing/PCR errors (e.g., chimeras, sequences with multiple indels) or
138 NUMTs so they were excluded from further analysis. Consideration of the assigned reads
139 revealed 26,263 BINs among the 52 sites with more than a third (9,301) found at only
140 one site (**Figure 2b**).

141

142 The Chao 1 [35] estimate for the total number of BINs present at the 52 sites was
143 29,640 (**Figure 2a**) while species richness extrapolation based on the lognormal
144 distribution (**Figure 2c**, [36]) suggested the presence of 31,516 BINs. On average, 0.3
145 million sequences were recovered per sample, and they revealed the presence of an
146 average of 2,352 BINs per site (range 996–4,581 BINs, **Table S2**) with bi-weekly
147 samples containing an average of 619 ± 14.3 S.E. BINs (range 60–1666, **Table S3**). Most

148 low BIN counts occurred in spring (May) or fall (September) with diversity peaking in
149 mid-summer (June/July). Taxonomic composition at an ordinal level was similar among
150 samples with over half of the BINs being flies (Diptera), followed by Hymenoptera,
151 Lepidoptera, Hemiptera, and Coleoptera.

152

153 Overlap in BIN composition was higher among parks in an ecoregion than among
154 those in different ecoregions, even after geographical distance was considered (**Figure**
155 **3a**). Sites in the ECF had the highest mean phylogenetic diversity followed by EGL and
156 finally SGL (**Figure 3b**), differences that were significant (KW and Dunn's posthoc $p <$
157 0.003). More BINs were collected in the ECF (14,001) than in the EGL (12,787) or SGL
158 (10,958) (Figure 3c). The Chao 1 estimates for the number of BINs present in each
159 ecoregion were 15,401 for ECF, 14,577 for EGL, and 12,602 for SGL. The three
160 ecoregions shared 4,133 BINs while about a third of those in each region were not
161 collected elsewhere. A two-dimensional NMDS Ordination plot revealed that BIN
162 assemblages for sites in each ecoregion formed cohesive groupings (Figure 3d).
163 PERMANOVA analysis also suggested that community structure varied between
164 ecoregions ($R^2 = 0.141$, $P = 0.0001$) and decreased site elevation ($R^2 = 0.035$,
165 $P = 0.03$).

166

167 Overall, α -diversity was highest in the ECF, intermediate in the EGL, and lowest in
168 SGL (**Figure 4**). The α -diversity patterns for the varied insect orders followed the overall
169 trend, but BIN richness for Collembola showed the opposite trend as it peaked in the
170 SGL, while spider α -diversity was highest in the EGL.

171

172 Levels of turnover (**Figure 5**) were generally high among sites (species replacement
173 by new species not found elsewhere) as well as high nestedness levels (gain and loss of
174 species also found elsewhere). Lower levels of both turnover and nestedness were
175 observed for most taxa at sites in the ECF while the highest values were found in the
176 SGL.

177

178 **Discussion**

179 This study used metabarcoding to examine the species represented in 410 Malaise
180 trap samples derived from 52 protected sites in three juxtaposed Canadian ecoregions.
181 Metabarcoding revealed 26,263 species of arthropods while Chao 1 and Preston
182 lognormal extrapolations indicated that another 3,000–5,000 species await detection. As
183 just 52 sites were surveyed, a more comprehensive sampling program in these ecoregions
184 might reveal as many as 50,000 species of arthropods. Nearly 5-fold variation (996–
185 4,581) in BIN counts were detected among sites; counts showed a similar range for the 30
186 sites where all samples were analyzed (996–4,508) and the 22 where just half were
187 metabarcoded (1,312–4,581). On average, 619 BINs were recovered from each
188 metabarcoded sample, a count that was 52.5% higher than the mean BIN count (406) for
189 samples that were barcoded (Steinke et al. in prep). This difference suggests that more
190 than half the BINs recovered from metabarcoded samples derive from environmental
191 DNA attached to specimens in the sample or from their gut contents.

192

193 The three ecoregions examined in this study collectively span 160,000 km², just
194 1.6% of Canada's land surface, but two (SGL, EGL) are among the most heavily
195 populated areas in the country [31]. The ecoregions showed considerable overlap in
196 species composition; 33.1% of the BINs recorded from three or more sites were shared by
197 the three ecoregions. BIN richness was lowest in the southernmost ecoregion (SGL) and
198 highest in the most northerly (ECF). This difference coincided with a disturbance
199 gradient -- from forested regions with low human density in the ECF (78% forest cover)
200 to disturbed landscapes dominated by farmland/cities in the SGL (12% forest cover). The
201 decline in species richness in response to disturbance is consistent with earlier studies
202 [37-39], even though our collections all derived from protected areas. [40] reported that
203 protected sites contain significantly higher species counts than adjacent disturbed areas,
204 perhaps because communities in protected areas include representatives of original
205 habitats and generalists from adjacent disturbed landscapes [41]. However, protected
206 areas in the SGL were small islands of remnant forest in a landscape dominated by
207 agricultural activity so they were undoubtedly heavily exposed to pesticides with
208 agricultural fields creating dispersal barriers which further reduced diversity.

209

210 Our results indicate that α -diversity for major insect orders of flying insects (Diptera,
211 Hymenoptera, Hemiptera, Lepidoptera) peaked in the least disturbed ecoregion (ECF).
212 By contrast, two groups of arthropods (Araneae, Collembola) lacking flight showed a
213 different trend with their diversity peaking in other ecoregions. This difference might
214 reflect the fact that Malaise traps only sample flightless taxa with resident populations
215 near the trap but capture flying insects from distant habitats. As such, biodiversity
216 patterns for flying insects provide a regional perspective while those for taxa without
217 flight provide a local perspective. If so, the reduction in diversity of Collembola from the
218 most southerly (SGL) to northerly (ECF) ecoregion might reflect the expected latitudinal
219 gradient in biodiversity, undisrupted by disturbance because of the local source of
220 specimens in each sample.

221

222 The present study establishes the feasibility of monitoring temporal changes in
223 species composition of arthropod communities [42,43]. For all three ecoregions, temporal
224 turnover was high, reflecting the seasonal succession of species. β -diversity was lowest
225 for most taxonomic groups at sites in the ECF and highest in the SGL. Species turnover
226 was generally higher than nestedness, suggesting the presence of many transient species
227 [44]. As many species were only collected at one or two sites, many samples likely
228 included transients passively transported by the wind [45].

229

230 Metabarcoding can already provide cost-effective biosurveillance as the present
231 study analyzed about 856,000 specimens and generated 223,860 species occurrence
232 records for \$82,000, an analytical cost of less than \$0.50 per record. By adopting simpler
233 analytical protocols (e.g., destructive processing of samples) with ongoing reductions in
234 sequencing costs [10], costs can be reduced by an order of magnitude, delivering species
235 occurrence records for \$0.04 apiece in the ecoregions targeted in this study. In settings
236 with higher α -diversity, the cost could be halved. Aside from its cost-effectiveness for
237 data acquisition, the digital format of metabarcoding results aids their curation,
238 validation, and preservation. Although current metabarcoding protocols cannot estimate
239 the abundance of each species in a sample, the situation shifts when multiple samples are
240 analyzed as the abundance of a species can then be estimated from its frequency of

241 occurrence in these samples (rare species will be recovered less frequently than abundant
242 taxa).

243

244 As the 846 currently recognized ecoregions [30] were largely delineated based on
245 distributional data for vascular plants and vertebrates, there remains a need to ascertain
246 how well they represent diversity patterns in other taxa. [46] found that arthropods
247 showed weak adherence to ecoregion boundaries and proposed this might reflect
248 dispersal limitations linked to their small body size or to the biased assemblage of
249 arthropod species with data. Our much larger dataset shows evidence of structuring by
250 ecoregion as both phylogenetic diversity and BIN composition were significantly
251 different among ecoregions, even when comparisons extended to widely separated sites.
252 This result suggests that ecoregions do provide a useful structural framework, reinforcing
253 results from earlier studies [47,48]. However, a third of species in this study crossed
254 ecoregion boundaries and more extensive sampling would raise the incidence of shared
255 species. The latter results make it clear that high sampling effort is required to better
256 understand species distributions. In looking to the future, it is apparent that there is an
257 immediate need for a more detailed understanding of the levels of species overlap
258 between adjacent ecoregions. Is, for example, the pattern of high overlap in species
259 composition among neighbouring ecoregions detected in this study a general pattern or
260 are some ecoregion boundaries sharply delineated? Such information is critical in
261 designing an effective global biomonitoring network to inform conservation efforts
262 [49,50].

263

264

265 **Potential Implications**

266 Past monitoring programs have provided limited insights into the shifting
267 distributions and abundances of arthropod species [51]. By coupling the use of an
268 efficient collection method with the capacity of DNA metabarcoding to determine the
269 species composition of bulk samples, this study has shown that compositional shifts in
270 arthropod communities can be tracked [52]. The present results also indicate that the
271 ecoregion concept not only furthers understanding of foundational biogeographic

272 principles and improves their potential application to conservation efforts, but also
273 provides a logical scaffold for large-scale monitoring networks.

274

275

276 **Methods**

277 *Sample collection*

278 An ez-Malaise trap (BioQuip Products) was deployed to collect arthropods at one
279 site in each of 50 provincial parks while two sites were sampled in the final park
280 (Algonquin) because of its large size. Trap catches were harvested every second week
281 from early May through September, producing 10 samples per site for a total of 520
282 samples. These samples were preserved in 95% ethanol and held at -20° C until DNA
283 extraction. Five samples (weeks 1+2, 5+6, 9+10, 13+14, 17+18) from each of 22 sites
284 were employed for single specimen barcoding (Steinke et al., in prep) while the other 410
285 samples were analyzed in this study. A direct count indicated that 230,000 specimens
286 were present in the 21.2% of the samples that were barcoded. On this basis, the remaining
287 samples (78.8%), those examined in this study, included about 856,000 specimens.

288

289 *DNA extraction and PCR*

290 DNA extraction employed a membrane-based protocol [32] modified for bulk
291 samples. Specimens were removed from ethanol by filtration through a sterile
292 Microfunnel 0.45 µm Supor Membrane Filter (Pall Laboratory) using a 6-Funnel
293 Manifold (Pall Laboratory). The wet weight of each sample was then ascertained to allow
294 volume adjustment (**Table S4**) of the lysis buffer [32]. Each sample was then incubated
295 overnight at 56°C while gently mixed on a shaker. Eight 50 µl aliquots (technical
296 replicates) from each of the 410 lysates were then transferred into 3,280 separate wells in
297 96-well microplates and DNA extracts were generated using Acroprep 3.0 µm glass
298 fiber/0.2 µm Bio-Inert membrane plates (Pall Laboratory). Each plate contained 80 lysate
299 samples, 8 technical replicates of a positive control (lysate from a bulk sample whose
300 component specimens were individually Sanger sequenced – public BOLD dataset -
301 dx.doi.org/10.5883/DS-AGAKS) and 8 negative controls. Each lysate was mixed with
302 100 µl of binding mix, transferred to a column plate, and centrifuged at 5000 g for 5 min.

303 DNA was then purified with three washes; the first employed 180 μ l of protein wash
304 buffer centrifuged at 5000 g for 5 min. Each column was then washed twice with 600 μ l
305 of wash buffer centrifuged at 5000 g for 5 min. Columns were transferred to clean tubes
306 and spun dry at 5000 g for 5 min to remove residual buffer before their transfer to clean
307 collection tubes followed by incubation for 30 min at 56°C to dry the membrane. DNA
308 was subsequently eluted by adding 60 μ l of 10 mM Tris-HCl pH 8.0 followed by
309 centrifugation at 5000 g for 5 min.

310

311 PCR reactions employed a standard protocol [53]. Briefly, each reaction included 5%
312 trehalose (Fluka Analytical), 1 \times Platinum Taq reaction buffer (Invitrogen), 2.5 mM
313 MgCl₂ (Invitrogen), 0.1 μ M of each primer (Integrated DNA Technologies), 50 μ M of
314 each dNTP (KAPA Biosystems), 0.3 units of Platinum Taq (Invitrogen), 2 μ l of DNA
315 extract, and Hyclone ultra-pure water (Thermo Scientific) for a final volume of 12.5 μ l.
316 Two-stage PCR was used to generate amplicon libraries for sequencing on an Ion Torrent
317 S5 platform. The first round of PCR used the primer combination AncientLepF3 [54] and
318 LepR1 [55] to amplify a 463 bp fragment of COI. Prior to the second PCR, first round
319 products were diluted 2x with ddH₂O. Fusion primers were then used to attach platform-
320 specific unique molecular identifiers (UMIs) along with the sequencing adaptors required
321 for Ion Torrent S5 libraries. Both rounds of PCR employed the same thermocycling
322 conditions: initial denaturation at 94 °C for 2 min, followed by 20 cycles of denaturation
323 at 94°C for 40 sec, annealing at 51°C for 1 min, and extension at 72 °C for 1 min, with a
324 final extension at 72°C of 5 min.

325

326 *HTS library construction*

327 For each plate, labelled products were pooled prior to sequencing. In total, 41
328 libraries were assembled. Each included eight technical replicates of 10 samples plus
329 eight technical replicates of a negative and a positive control respectively (i.e., 96
330 samples). The ten samples from each of the 30 sites that were only metabarcoded,
331 together with positive and negative controls, were pooled after UMI tagging to create a
332 library that was analyzed on a 530 chip (30 chips in total). Five samples were available
333 from each the other 22 sites (where half the samples were retained for barcoding). The

334 UMI-tagged amplicons from five samples from each of two sites were pooled with
335 positive and negative controls to produce a single library. Amplicon libraries were
336 prepared on an Ion Chef (Thermo Fisher Scientific) following and sequenced on an Ion
337 Torrent S5 platform at the Centre for Biodiversity Genomics following manufacturer's
338 instructions (Thermo Fisher Scientific).

339

340 *Sequence analysis*

341 Reads from the eight replicates for each sample were concatenated using a bash
342 script and uploaded to mBRAVE (<http://mbrave.net/>) for quality filtering and subsequent
343 queries using several reference libraries in an open reference approach. All reads were
344 queried against five system libraries on mBRAVE: bacteria (SYS-CRLBACTERIA),
345 chordates (SYS-CRLCHORDATA), insects (SYS-CRLINSECTA), non-insect
346 arthropods (SYS-CRLNONINSECTARTH), and non-arthropod invertebrates (SYS-
347 CRLNONARTHINVERT). Sequences were only included in this analysis if they
348 possessed a minimum length >350 bp and met the following three quality criteria (Mean
349 QV >20; <25% positions with a QV<20; <5% positions with QV<10). Reads were
350 trimmed 30 bp from their 5' terminus with a set trim length of 450 bp. Reads were
351 matched to the sequences in each reference library with an ID distance threshold of 3%,
352 but were only retained for further analysis when at least three reads matched an OTU in
353 the reference database. All reads failing to match any sequence in the five reference
354 libraries were clustered at an OTU threshold of 1% with a minimum of five reads per
355 cluster. All raw data are available in the NCBI Short Read Archive (PRJNA629553).

356 Using mBRAVE, we generated BIN (and OTU) tables including all library queries
357 for each individual plate/run (10 samples, plus a negative and positive control -
358 [dx.doi.org/10.5883/DS-AGAKS](https://doi.org/10.5883/DS-AGAKS) - for each run). Read counts for any BINs recovered
359 from the negative control on a plate were subtracted from the counts for the same BIN in
360 the 80 non-control wells in the run. When this subtraction reduced the read count for a
361 BIN to zero, its occurrence was removed. This step reduced the effects of rare tag
362 switching on data integrity [56] and removed any background contamination.

363

364 *Ecoregion analysis*

365 To determine the completeness of sampling, we calculated accumulation curves and
366 the Chao-1 estimator for total diversity [35] using the vegan package [57]. For further
367 extrapolation of species richness, we used the lognormal species abundance distribution
368 [36]. The fit of Fisher's Logseries [58] was used to determine relative BIN abundance.
369 Both methods are implemented in vegan (fisherfit, prestonfit) [57]. We calculated
370 Sørensen's similarity coefficient to ascertain if differences in species assemblages were
371 greater between or across ecoregion borders after controlling for distance. Differences in
372 BIN composition among the three ecoregions were examined using non-metric
373 multidimensional scaling (NMDS) with the Bray-Curtis index coefficient as implemented
374 in vegan [57]. The adonis function of the vegan package was used to conduct a
375 Permutational Multivariate Analysis of Variance (PERMANOVA) to partition distance
376 matrices among sources of variation (factors such as latitude, longitude, elevation, and
377 ecoregion).

378 A Maximum likelihood phylogeny was inferred for a BIN sequence alignment using
379 RAxML Black box [59] on XCEDE via the CIPRES portal [60]. The resulting phylogeny
380 comprising 26,263 BIN sequences was used to calculate Faith's phylogenetic distance
381 (PD) [61] using the picante package [62]. Because this measure is influenced by
382 polytomies in a phylogeny [63], only one representative was included per BIN to avoid
383 bias introduced by variation in the number of records for each BIN. A Kruskal-Wallis test
384 followed by a Dunn's posthoc analysis was used to determine if significant PD
385 differences existed between ecoregions.

386 Alpha (α)-diversity was quantified as the number of BINs observed at a site. Beta
387 (β)-diversity was computed as multi-site Sorensen and Simpson indices using the betapart
388 1.3. package [64]. β -diversity calculations between pairs of ecoregions were computed
389 using 12 random sites from the total sites for each ecoregion, and resampled 1000 times.
390 We then decomposed the among-site β -diversity into its turnover (species replacement
391 from site to site) and nestedness (species gain/loss from sites) components. Pairwise BIN
392 diversity among ecoregions was evaluated using the nonparametric multiple comparison
393 function implemented in the R package dunn.test 1.2.4 [65]. dunn.test is equivalent to the
394 Kruskal-Wallis and pair-wise Mann-Whitney post hoc tests with Bonferroni correction.

395 All analyses were performed in R v.3.4.4 [66].

396

397 **Funding**

398 This study was enabled by awards to PDNH from the Ontario Ministry of Research,
399 Innovation and Science, the Canada Foundation for Innovation, and by a grant from the
400 Canada First Research Excellence Fund to the University of Guelph's "Food From
401 Thought" research program.

402 **Author contributions**

403 DS, EVZ, JRDW, PDNH designed the study. DS, JRDW, JES, KP coordinated the study.
404 SLDW, NVI, SWJP, TWAB did the bench work and contributed to analyses. SR and
405 MM oversaw database organisation. DS did the analyses and wrote the manuscript.
406 PDNH, JRDW, EVZ, TWAB revised the manuscript.

407

408 **Acknowledgements**

409 We thank the collections and sequencing staff at the Centre for Biodiversity Genomics
410 for acquiring and processing the specimens analyzed in this study. We are very grateful to
411 Suz Bateson for improving the figures and to staff at the participating Ontario Provincial
412 Parks for facilitating collections.

413

414 **References**

415

- 416 1. Hallmann CA, Sorg M, Jongejans E, Siepel H, Hofland N, Schwan H, Stenmans W,
417 Müller A, Sumser H, Hörren T, Goulson D, de Kroon H. More than 75 percent decline
418 over 27 years in total flying insect biomass in protected areas. *PLoS ONE*. 2017; 12(10):
419 e0185809.
- 420 2. Lister BC, Garcia A. Climate-driven declines in arthropod abundance restructure a
421 rainforest food web. *Proceedings of the National Academy of Sciences of the United*
422 *States of America*. 2018; 115(44): E10397–E10406.
- 423 3. Macgregor CJ, Williams JH, Bell JR, Thomas CD. Moth biomass increases and
424 decreases over 50 years in Britain. *Nature Ecology and Evolution*. 2019; 3: 1645–1649.
- 425 4. Seibold S, Gossner MM, Simons NK, Blüthgen N, Müller J, Ambarli D, Ammer C,
426 Bauhus J, Fischer M, Habel JC, Linsenmair KE, Nauss T, Penone C, Prati D, Schall P,

- 427 Schulze E-D, Vogt J, Wöllauer S, Weisser WW. Arthropod decline in grasslands and
428 forests is associated with drivers at landscape level. *Nature*. 2019; 574: 671–674
- 429 5. Bush A, Sollmann R, Wilting A, Bohmann K, Cole B, Balzter H, Martius C, Zlinszky
430 A, Calvignac-Spencer S, Cobbold CA, Dawson TP, Emerson BC, Ferrirer S, Gilbert
431 MTP, Herold M, Jones L, Leendertz FH, Matthews L, Millington JDA, Olson JR,
432 Ovaskainen O, Raffaelli D, Reeve R, Rödel M-O, Rodgers TW, Snape S, Visseren-
433 Hamakers I, Vogler AP, White PCL, Wooster MJ, Yu DW. Connecting Earth observation
434 to high-throughput biodiversity data. *Nature Ecology & Evolution*. 2017; 1: 0176.
- 435 6. Hebert PDN, Cywinska A, Ball SL, deWaard JR. Biological identifications through
436 DNA barcodes. *Proceedings of the Royal Society B: Biological Science*. 2003; 270: 313–
437 321.
- 438 7. Hebert PDN, Braukmann TWA, Prosser SWJ, Ratnasingham S, deWaard JR, Ivanova
439 NV, Janzen DH, Hallwachs W, Naik S, Sones JE, Zakharov EV. A Sequel to Sanger:
440 amplicon sequencing that scales. *BMC Genomics*. 2018; 19: 219.
- 441 8. Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E. Towards next-
442 generation biodiversity assessment using DNA metabarcoding. *Molecular Ecology*. 2012;
443 21(8): 2045-2050.
- 444 9. O’Driscoll A, Daugelaite J, Sleator RD. ‘Big Data’, Hadoop and cloud computing in
445 genomics. *Journal of Biomedical Informatics*. 2013; 46(5): 774–781.
- 446 10. Lightbody G, Haberland V, Browne F, Taggart L, Zheng H, Parkes E, Blayney JK.
447 Review of applications of high-throughput sequencing in personalized medicine: barriers
448 and facilitators of future progress in research and clinical application. *Briefings in*
449 *Bioinformatics*. 2019; 20(5): 1795–1811.
- 450 11. Ji C, Chng KR, Hui Boey EJ, Ng AHQ, Wilm A, Nagarajan N. INC-Seq: accurate
451 single molecule reads using nanopore sequencing. *Gigascience*. 2016; 5: 34.
- 452 12. Beng KC, Tomlinson KW, Shen XH, Surget-Groba Y, Hughes AC, Corlett RT, Slik
453 JWF. The utility of DNA metabarcoding for studying the response of arthropod diversity
454 and composition to land-use change in the tropics. *Scientific Reports*. 2016; 6: 1–13.
- 455 13. Elbrecht V, Vamos EE, Meissner K, Aroviita J, Leese F. Assessing strengths and
456 weaknesses of DNA metabarcoding-based macroinvertebrate identification for routine
457 stream monitoring. *Methods in Ecology and Evolution*. 2017; 8: 1–21.

- 458 14. D'Souza ML, van der Bank M, Zandisile S, Rattray RD, Stewart R, van Rooyen J,
459 Govender D, Hebert PDN. Biodiversity baselines: tracking insects in Kruger National
460 Park with DNA barcodes. *Biological Conservation*. 2021; 256: 109034.
- 461 15. Sato H, Sogo Y, Doi H, Yamanaka H. Usefulness and limitations of sample pooling
462 for environmental DNA metabarcoding of freshwater fish communities. *Scientific*
463 *Reports*. 2017; 7: 14860.
- 464 16. Bell KL. Applying pollen DNA metabarcoding to the study of plant-pollinator
465 interactions. *Applications in Plant Sciences*. 2017; 5: apps.1600124
- 466 17. Vasselon V, Bouchez A, Rimet F, Jacquet S, Trobajo R, Corniquel M, Tapolczai K,
467 Domaizon I. Avoiding quantification bias in metabarcoding: Application of a cell
468 biovolume correction factor in diatom molecular biomonitoring (A. Mahon, Ed.).
469 *Methods in Ecology and Evolution*. 2018; 9: 1060–1069.
- 470 18. Bellemain E, Davey ML, Kauserud H, Epp LS, Boessenkool S, Coissac E, Geml J,
471 Edwards M, Willerslev E, Gussarova G, Taberlet P, Haile J, Brochmann C. Fungal
472 palaeodiversity revealed using high-throughput metabarcoding of ancient DNA from
473 arctic permafrost. *Environmental Microbiology*. 2012; 15: 1176–1189.
- 474 19. Aas AB, Davey ML, Kauserud H. ITS all right mama: investigating the formation of
475 chimeric sequences in the ITS2 region by DNA metabarcoding analyses of fungal mock
476 communities of different complexities. *Molecular Ecology Resources*. 2017; 17: 730–
477 741.
- 478 20. Tedersoo L, Tooming-Klunderud A, Anslan S. PacBio metabarcoding of Fungi and
479 other eukaryotes: errors, biases, and perspectives. *New Phytologist*. 2018; 217: 1370–
480 1385.
- 481 21. Malaise R. A new insect trap. *Entomologisk Tidskrift*. 1937; 58: 148–160.
- 482 22. Karlsson D, Pape T, Johanson KA, Liljeblad J, Ronquist F. The Swedish Malaise
483 Trap Project, or how many species of Hymenoptera and Diptera are there in Sweden?
484 *Entomologisk Tidskrift*. 2005; 126: 43–53.
- 485 23. deWaard JR, Levesque-Beaudin V, deWaard SL, Ivanova NV, McKeown JTA,
486 Miskie R, Naik S, Perez KHJ, Ratnasingham S, Sobel CN, Sones JE, Steinke C, Telfer
487 AC, Young A, Young MR, Zakharov EV, Hebert PDN. Expedited assessment of

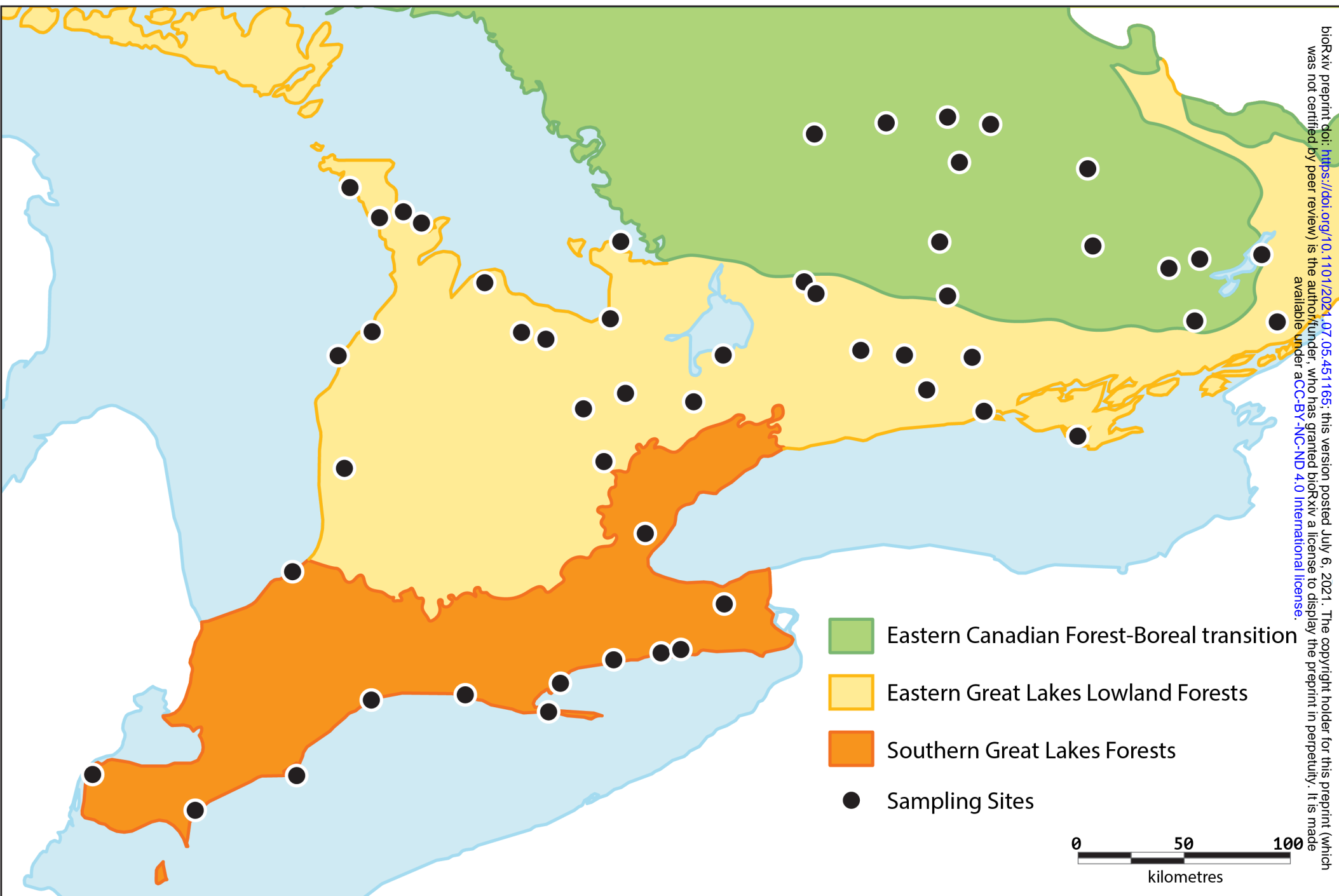
- 488 terrestrial arthropod diversity by coupling Malaise traps with DNA barcoding. *Genome*.
489 2019; 62: 85–95.
- 490 24. Steinke D, Braukmann TWA, Manerus L, Woodhouse A, Elbrecht V. Effects of
491 Malaise trap spacing on species richness and composition of terrestrial arthropod bulk
492 samples. *Metabarcoding and Metagenomics*. 2021; 5: 43–50.
- 493 25. Holdridge LR. Determination of world plant formations from simple climatic data.
494 *Science*. 1947; 105: 367–368.
- 495 26. Whittaker RH. Classification of natural communities. *Botanical Reviews*. 1962; 28:
496 1–239.
- 497 27. Olson DM, Dinerstein E, Wikramanayake ED, Burgess ND, Powell GVN,
498 Underwood EC, D’amico JA, Itoua I, Strand HE, Morrison JC, Loucks CJ, Allnutt TF,
499 Ricketts TH, Kura Y, Lamoreux JF, Wettengel WW, Hedao P, Kassem KR. Terrestrial
500 ecoregions of the world: a new map of life on earth. *Bioscience*. 2001; 51: 933–938.
- 501 28. Bailey RG. *Ecoregions*. Springer, New York; 2014.
- 502 29. Giakoumi S, Sini M, Gerovasileiou V, Mazor T, Beher J, Possingham HP, Abdulla A,
503 Cinar ME, Dendrinou P, Gucu AC, Karamanlidis AA, Rodic P, Panayotidis P, Taskin E,
504 Jaklin A, Voultziadou E, Webster C, Zenetos A, Katsanevakis S. Ecoregion-based
505 conservation planning in the Mediterranean: Dealing with large-scale heterogeneity.
506 *PLoS ONE*. 2013; 8(10): e76449.
- 507 30. Dinerstein E, Olson D, Joshi A, Vynne C, Burgess ND, Wikramanayake E, Hahn N,
508 Palminteri S, Hedao P, Noss R, Hansen M, Locke H, Ellis EC, Jones B, Barber CV,
509 Hayes R, Kormos C, Martin V, Crist E, Sechrest W, Price L, Baillie JEM, Weeden D,
510 Suckling K, Davis C, Sizer N, Moore R, Thau D, Birch T, Potapov P, Turubanova S,
511 Tyukavina A, de Souza N, Pintea L, Brito JC, Llewellyn OA, Miller AG, Patzelt A,
512 Ghazanfar SA, Timberlake J, Klöser H, Shennan-Farpón Y, Kindt R, Barnekow Lillesø
513 J-P, van Breugel P, Graudal L, Vogé M, Al-Shammari KF, Saleem M. An ecoregion-
514 based approach to protecting half the terrestrial realm. *Bioscience*. 2013; 67: 534–545.
- 515 31. Crins WJ, Gray PA, Uhlig PWC, Wester MC. *The Ecosystems of Ontario, Part 1:*
516 *Ecozones and Ecoregions*. Technical Report SIB TER IMA TR-01, Ministry of Natural
517 Resources, Ontario; 2009.

- 518 32. Ivanova NV, deWaard JR, Hebert PDN. An inexpensive, automation-friendly
519 protocol for recovering high-quality DNA. *Molecular Ecology Resources*. 2006; 6: 998–
520 1002.
- 521 33. Ratnasingham S and PDN Hebert. A DNA-based registry for all animal species: The
522 Barcode Index Number (BIN) System. *PLoS ONE*. 2013; 8: e66213.
- 523 34. Ratnasingham S and PDN Hebert. BOLD: The Barcode of Life Data System
524 (www.barcodinglife.org). *Molecular Ecology Notes*. 2007; 7: 355–364.
- 525 35. Magurran AE. *Measuring Biological Diversity*. Wiley-Blackwell, Malden,
526 Massachusetts; 2003.
- 527 36. Preston FW. The canonical distribution of commonness and rarity: Part I. *Ecology*.
528 1962; 43: 185–215.
- 529 37. Luke SH, Fayle TM, Eggleton P, Turner EC, Davies RG. Functional structure of ant
530 and termite assemblages in old growth forest, logged forest and oil palm plantation in
531 Malaysian Borneo. *Biodiversity Conservation*. 2014; 23: 2817–2832.
- 532 38. Newbold T, Hudson LN, Phillips HRP, Hill SLL, Contu S, Lysenko I, Blandon A,
533 Butchart SHM, Booth HL, Day J, De Palma A, Harrison MLK, Kirkpatrick L, Pynegar E,
534 Robinson A, Simpson J, Mace GM, Scharlemann JPW, Purvis A. A global model of the
535 response of tropical and sub-tropical forest biodiversity to anthropogenic pressures.
536 *Proceedings of the Royal Society B*. 2014; 281: 20141435.
- 537 39. Phalan B, Onial M, Balmford A, Green RE. Reconciling food production and
538 biodiversity conservation: Land sharing and land sparing compared. *Science*. 2011; 333:
539 1289–1291.
- 540 40. Gray CL, Hill SLL, Newbold T, Hudson LN, Boerger L, Contu S, Hoskins AJ, Ferrier
541 S, Purvis A, Scharlemann JPW. Local biodiversity is higher inside than outside terrestrial
542 protected areas worldwide. *Nature Communications*. 2016; 7: 12306.
- 543 41. Lingbeek BJ, Higgins CL, Muir JP, Kattes DH, Schwertner TW. Arthropod diversity
544 and assemblage structure response to deforestation and desertification in the Sahel of
545 western Senegal. *Global Ecology and Conservation*. 2017; 11: 165–176.
- 546 42. Tschamtko T, Tylianakis JM, Rand TA, Didham RK, Fahring L, Batary P, Bengtsson
547 J, Clough Y, Crist TO, Dormann CF, Ewers RM, Freund J, Holt RD, Holzschuh A, Klein
548 AM, Kleijn D, Kremen C, Landis DA, Laurance W, Lindenmayer D, Scherber C, Sodhi

- 549 N, Steffan-Dewenter I, Thies C, van der Putten WM, Westphal C. Landscape moderation
550 of biodiversity patterns and processes – eight hypotheses. *Biological Reviews*. 2012; 87:
551 661–685.
- 552 43. Myers JA, Chase JM, Jiminez I, Jorgensen PM, Araujo-Murakami A, Paniagua-
553 Zambrana N, Seidel R. Beta-diversity in temperate and tropical forests reflects dissimilar
554 mechanisms of community assembly. *Ecology Letters*. 2013; 16: 151–157.
- 555 44. Snell Taylor SJ, Evans BS, White EP, Hurlbert AH. The prevalence and impact of
556 transient species in ecological communities. *Ecology*. 2018; 99(8): 1825–1835.
- 557 45. D’Souza ML, Hebert PDN. Stable baselines of temporal turnover underlie beta
558 diversity in tropical arthropod communities. *Molecular Ecology*. 2018; 27: 2447–2460.
- 559 46. Smith JR, Letten AD, Ke P-J, Anderson CB, Hendershot JN, Dhami MK, Dlott GA,
560 Grainger TN, Howard ME, Morrison BML, Routh D, San Juan PA, Mooney HA,
561 Mordecai EA, Crowther TW, Daily GC. A global test of ecoregions. *Nature Ecology &*
562 *Evolution*. 2018; 2: 1889–1896.
- 563 47. Lightfoot DC, Brantely SL, Allen CD. Geographic patterns of ground-dwelling
564 arthropods across an ecological transition in the North American southwest. *Western*
565 *North American Naturalist*. 2008; 68: 83–102.
- 566 48. Gonzales-Reyes AX, Corronca JA, Arroyo NC. Differences in alpha and beta
567 diversities of epideous arthropod assemblages in two ecoregions of northwestern
568 Argentina. *Zoological Studies*. 2012; 51: 1367–1379.
- 569 49. Watson JEM, Venter O. Ecology: a global plan for nature conservation. *Nature*. 2017;
570 550: 48–49.
- 571 50. Wilson EO. *Half-Earth: Our Planet’s Fight for Life*, Liveright, New York; 2017.
- 572 51. Díaz S, Settele J, Brondízio ES, Ngo HT, Guèze M, Agard J, Arneth A, Balvanera P,
573 Brauman KA, Butchart SHM, Chan KMA, Garibaldi LA, Ichii K, Liu J, Subramanian
574 SM, Midgley GF, Miloslavich P, Molnár Z, Obura D, Pfaff A, Polasky S, Purvis A,
575 Razaque J, Reyers B, Chowdhury RR, Shin YJ, Visseren-Hamakers IJ, Willis KJ, Zayas
576 CN (eds.). *Summary for policymakers of the global assessment report on biodiversity and*
577 *ecosystem services of the Intergovernmental Science-Policy Platform on Biodiversity and*
578 *Ecosystem Services*. IPBES secretariat, Bonn, Germany; 2019.

- 579 52. Hobern D. BIOSCAN: DNA barcoding to accelerate taxonomy and biogeography for
580 conservation and sustainability. *Genome*. 2021; 64: 161–164.
- 581 53. Braukmann TWA, Prosser SJR, Ivanova NV, Elbrecht V, Steinke D, Ratnasingham
582 R, deWaard JR, Sones JE, Zakharov EV, Hebert PDN. Metabarcoding a diverse
583 arthropod mock community. *Molecular Ecology Resources*. 2019; 19: 711–727.
- 584 54. Prosser SWJ, deWaard JR, Miller SE, and PDN Hebert. DNA barcodes from century-
585 old type specimens using next-generation sequencing. *Molecular Ecology Resources*.
586 2016; 16: 487–497.
- 587 55. Hebert PDN, Penton EH, Burns JM, Janzen DH, Hallwachs W. Ten species in one:
588 DNA barcoding reveals cryptic species in the neotropical skipper butterfly *Astraptes*
589 *fulgerator*. *Proceedings of the National Academy of Sciences of the United States of*
590 *America*. 2004; 101: 14812–14817.
- 591 56. Elbrecht V, Steinke D. Scaling up DNA metabarcoding for freshwater
592 macrozoobenthos monitoring. *Freshwater Biology*. 2018; 64: 380–387.
- 593 57. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlenn D, Minchin PR,
594 O'Hara RB, Simpson GL, Solymos P, Stevens MHH, Szoecs E, Wagner H. *vegan*:
595 *Community Ecology Package*. R package version 2.5-1. [https://CRAN.R-](https://CRAN.R-project.org/package=vegan)
596 [project.org/package=vegan](https://CRAN.R-project.org/package=vegan); 2018
- 597 58. Fisher RA, Corbet AS, Williams CB. The relation between the number of species and
598 the number of individuals in a random sample of animal population. *Journal of Animal*
599 *Ecology*. 1943; 12: 42–58.
- 600 59. Stamatakis A, Hoover P, Rougemont J. A rapid bootstrap algorithm for the RAxML
601 web servers. *Systematic Biology*. 2008; 57(5): 758–771.
- 602 60. Miller MA, Pfeiffer W, Schwartz T. The CIPRES science gateway. In: *Proceedings*
603 *of the 2011 TeraGrid Conference on Extreme Digital Discovery—TG '11*. New York,
604 USA: ACM Press; 2011.
- 605 61. Faith DP. Conservation evaluation and phylogenetic diversity. *Biological*
606 *Conservation*. 1992; 61: 1–10.
- 607 62. Kembel SW, Cowan PD, Helmus MR, Cornwell WK, Morlon H, Ackerly DD,
608 Blomberg SP, Webb CO. Picante: R tools for integrating phylogenies and ecology.
609 *Bioinformatics*. 2010; 26(11): 1463–1464.

- 610 63. Swenson NG. Phylogenetic resolution and quantifying the phylogenetic diversity and
611 dispersion of communities. *PLoS ONE*. 2009; 4(2): e4390.
- 612 64. Baselga A, Orme CDL. betapart: an R package for the study of beta diversity.
613 *Methods Ecology and Evolution*. 2012; 3: 808–812.
- 614 65. Dinno A. *dunn.test: Dunn's Test of Multiple Comparisons Using Rank Sums*. R
615 *package version 1.3.2*. <http://CRAN.R-project.org/package=dunn.test>; 2016.
- 616 66. R Core Team. *R: A language and environment for statistical computing*. R
617 Foundation for Statistical Computing, Vienna, Austria; 2018. URL [https://www.R-](https://www.R-project.org/)
618 [project.org/](https://www.R-project.org/).



bioRxiv preprint doi: <https://doi.org/10.1101/2021.07.05.451165>; this version posted July 6, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

