

GigaScience

Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale

--Manuscript Draft--

Manuscript Number:	GIGA-D-18-00245R2	
Full Title:	Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale	
Article Type:	Research	
Funding Information:	National Science Foundation (DEB 1457845.)	Not applicable
Abstract:	<p>Background</p> <p>In light of the current biodiversity crisis, DNA barcoding is developing into an essential tool to quantify state shifts in global ecosystems. Current barcoding protocols often rely on short amplicon sequences, which yield accurate identification of biological entities in a community, but provide limited phylogenetic resolution across broad taxonomic scales. However, the phylogenetic structure of communities is an essential component of biodiversity. Consequently, a barcoding approach is required that unites robust taxonomic assignment power and high phylogenetic utility. A possible solution is offered by sequencing long ribosomal DNA (rDNA) amplicons on the MinION platform (Oxford Nanopore Technologies).</p> <p>Findings</p> <p>Using a dataset of various animal and plant species, with a focus on arthropods, we assemble a pipeline for long rDNA barcode analysis and introduce a new software (MiniBar) to demultiplex dual indexed nanopore reads. We find excellent phylogenetic and taxonomic resolution offered by long rDNA sequences across broad taxonomic scales. We highlight the simplicity of our approach by field barcoding with a miniaturized, mobile laboratory in a remote rainforest. We also test the utility of long rDNA amplicons for analysis of community diversity through metabarcoding and find that they recover highly skewed diversity estimates.</p> <p>Conclusions</p> <p>Sequencing dual indexed, long rDNA amplicons on the MinION platform is a straightforward, cost effective, portable and universal approach for eukaryote DNA barcoding. Although bulk community analyses using long-amplicon approaches may introduce biases, the long rDNA amplicons approach signifies a powerful tool for enabling the accurate recovery of taxonomic and phylogenetic diversity across biological communities.</p>	
Corresponding Author:	Henrik Krehenwinkel UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Henrik Krehenwinkel	
First Author Secondary Information:		
Order of Authors:	Henrik Krehenwinkel	

	Aaron Pomerantz
	James B. Henderson
	Susan R. Kennedy
	Jun Ying Lim
	Varun Swamy
	Juan Diego Shoobridge
	Natalie Graham
	Nipam H. Patel
	Rosemary G. Gillespie
	Stefan Prost
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Dear Dr. Edmunds, thank you and the two reviewers for your effort in improving our manuscript. We have addressed the few additional changes raised by Reviewer 1 and hope the manuscript is now deemed acceptable for publication. Sincerely, Henrik Krehenwinkel</p> <p>Reviewer reports: Reviewer #1: Most of the issues have been well addressed except for the ones regarding to mutation saturations of COI gene and the biased description of ITS regions.</p> <p>The authors stated that "The phylogenetic resolution offered by short barcodes is very limited, as they contain only a restricted number of informative sites. This problem is exacerbated by the fast evolutionary rate of mitochondrial DNA, which leads to a quick saturation with mutations, increasing the probability of homoplasy." For DNA sequences, homoplasy can hardly be avoided due to its four-state nature. However, as I mentioned before, COI gene has > 600 sites, it is going to be extremely rare or, I would say, impossible for the entire gene getting saturated. The authors may want to provide citations here to illustrate how mutation saturation affect phylogenetic resolution?</p> <ul style="list-style-type: none"> •We do not question the utility of COI as a barcode marker, for which it is very well suited. Our main intention was to highlight the utility of backing up mitochondrial information with nuclear data. The sentence in question is really not essential to communicate that point. We thus deleted the according sentence. <p>For the ITS, as another reviewer also mentioned, ITS2 region is widely utilized to serve as barcode sequences for fungi and plants and less variable in length than other ITS regions. The authors may want to add an unbiased description of ITS regions in their main text.</p> <ul style="list-style-type: none"> •We have now added some additional detail on the differences between ITS1 and ITS2. <p>Reviewer #2: The authors of the manuscript entitled 'Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale' have addressed all my concerns and suggestions and I am happy to recommend this manuscript for publication.</p> <p>I would recommend to move Figure 3 to supplementary material, as this is merely a confirmation that the sequences do not produce false phylogenetic signal.</p>

	•We feel this figure is important, as it nicely highlight the broad taxonomic utility of our method. We would thus prefer to leave it in the main text.
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	Yes
Resources A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible. Have you included the information requested as detailed in our Minimum Standards Reporting Checklist ?	Yes
Availability of data and materials All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials”	Yes

section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

[Click here to view linked References](#)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple**
2 **biodiversity assessments with high phylogenetic resolution across broad taxonomic**
3 **scale**

4
5 Henrik Krehenwinkel^{1,4}, Aaron Pomerantz², James B. Henderson^{3,4}, Susan R. Kennedy¹, Jun
6 Ying Lim^{1,2}, Varun Swamy⁵, Juan Diego Shoobridge⁶, Natalie Graham¹, Nipam H. Patel^{2,7},
7 Rosemary G. Gillespie¹, Stefan Prost^{2,8}

8
9 ¹ Department of Environmental Science, Policy and Management, University of California,
10 Berkeley, USA

11 ² Department of Integrative Biology, University of California, Berkeley, USA

12 ³ Institute for Biodiversity Science and Sustainability, California Academy of Sciences, San
13 Francisco, USA

14 ⁴ Center for Comparative Genomics, California Academy of Sciences, San Francisco, USA

15 ⁵ San Diego Zoo Institute for Conservation Research, Escondido, USA

16 ⁶ Applied Botany Laboratory, Research and development Laboratories, Cayetano Heredia
17 University, Lima, Perú

18 ⁷ Department of Molecular and Cell Biology, University of California, Berkeley, USA

19 ⁸ Research Institute of Wildlife Ecology, Department of Integrative Biology and Evolution,
20 University of Veterinary Medicine, Vienna, Austria

21
22 Corresponding authors: Henrik Krehenwinkel (krehenwinkel@berkeley.edu, ORCID: 0000-0001-
23 5069-8601) and Stefan Prost (stefan.prost@berkeley.edu, ORCID: 0000-0002-6229-3596).

24
25
26 **Abstract**

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

27

Background: In light of the current biodiversity crisis, DNA barcoding is developing into an essential tool to quantify state shifts in global ecosystems. Current barcoding protocols often rely on short amplicon sequences, which yield accurate identification of biological entities in a community, but provide limited phylogenetic resolution across broad taxonomic scales. However, the phylogenetic structure of communities is an essential component of biodiversity. Consequently, a barcoding approach is required that unites robust taxonomic assignment power and high phylogenetic utility. A possible solution is offered by sequencing long ribosomal DNA (rDNA) amplicons on the MinION platform (Oxford Nanopore Technologies).

36

Findings: Using a dataset of various animal and plant species, with a focus on arthropods, we assemble a pipeline for long rDNA barcode analysis and introduce a new software (MiniBar) to demultiplex dual indexed nanopore reads. We find excellent phylogenetic and taxonomic resolution offered by long rDNA sequences across broad taxonomic scales. We highlight the simplicity of our approach by field barcoding with a miniaturized, mobile laboratory in a remote rainforest. We also test the utility of long rDNA amplicons for analysis of community diversity through metabarcoding and find that they recover highly skewed diversity estimates.

44

Conclusions: Sequencing dual indexed, long rDNA amplicons on the MinION platform is a straightforward, cost effective, portable and universal approach for eukaryote DNA barcoding. Although bulk community analyses using long-amplicon approaches may introduce biases, the long rDNA amplicons approach signifies a powerful tool for enabling the accurate recovery of taxonomic and phylogenetic diversity across biological communities.

50

Keywords

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
61
62
63
64
65

52 Biodiversity, ribosomal, eukaryotes, long DNA barcodes, Oxford Nanopore Technologies,
53 MinION, metabarcoding

54

55 **Background**

56

57 The world is changing at an unprecedented rate, threatening the integrity of biological
58 communities [1, 2]. To understand the impacts of change, whether a system is close to a regime
59 shift, and how to mitigate the impacts of a given environmental stressor, it is important to
60 consider the biological community as a whole. In recognition of this need, there has been a shift
61 in emphasis from studies that focus on single indicator taxa, to comparative studies across
62 multiple taxa and metrics that consider the properties of entire communities [3]. Such efforts
63 require accurate information on the identity of the different biological entities within a
64 community, as well as the phylogenetic diversity that they represent.

65

66 Comparative ecological studies across multiple taxa have been greatly simplified by molecular
67 barcoding [4], where species identifications are based on short PCR amplicon “barcode”
68 sequences. Different barcode marker genes have been established across the tree of life [5, 6],
69 with mitochondrial cytochrome oxidase subunit I (COI) commonly used for animal barcoding [4].

70 The availability of large sequence reference databases and universal primers, together with its
71 uniparental inheritance and fast evolutionary rate, make COI a useful marker to distinguish even
72 recently diverged taxa. In recent years, DNA barcoding has greatly profited from the emergence
73 of next generation sequencing (NGS) technology. Current NGS platforms enable the parallel
74 generation of barcodes for hundreds of specimens at a fraction of the cost of Sanger
75 sequencing [7]. Furthermore, NGS technology has enabled metabarcoding, the sequencing of
76 bulk community samples, which allows scoring the diversity of entire ecosystems [8].

77

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

78 However, despite their undeniable advantages, barcoding approaches using short,
79 mitochondrial markers have several drawbacks. The phylogenetic resolution offered by short
80 barcodes is very limited, as they contain only a restricted number of informative sites. While this
81 does not affect the taxonomic utility of COI, it causes problems in phylogenetic analyses of
82 divergent lineages. The accurate estimation of phylogenetic diversity across wide taxonomic
83 scales, however, is an important component of biodiversity research [9]. Moreover,
84 mitochondrial DNA is not always the best marker to reflect species differentiation, as different
85 factors are known to inflate mitochondrial differentiation in relation to the nuclear genomic
86 background. For example, male biased gene flow [10] or infections with reproductive parasites
87 [11] (e.g. *Wolbachia*) can lead to highly divergent mitochondrial lineages in the absence of
88 nuclear differentiation. In contrast, introgressive hybridization can cause the complete
89 replacement of mitochondrial genomes (see e.g. [12, 13]), resulting in shared mitochondrial
90 variation between species.

91
92 Considering this background, it would be desirable to complement mitochondrial DNA based
93 barcoding with additional information from the nuclear genome. An ideal nuclear barcoding
94 marker should possess sufficient variation to distinguish young species pairs, but also provide
95 support for phylogenetic hypotheses between divergent lineages. Moreover, the marker should
96 be present across a wide range of taxa and amplification should be possible using universal
97 primers. A marker that fulfils all the above requirements is the nuclear ribosomal DNA (rDNA).
98 As an essential component of the ribosomal machinery, rDNA is a common feature across the
99 tree of life from microbes to higher eukaryotes [14]. All eukaryotes share homologous
100 transcription units of the 18S, 5.8S and 28S-rDNA genes, which include two internal transcribed
101 spacers (ITS1 and ITS2) [15]. Due to varying evolutionary constraints acting on different parts of
102 the rDNA, it consists of regions of extreme sequence conservation, which are interrupted by
103 highly variable stretches [16]. While some rDNA gene regions are entirely conserved across all

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

104 eukaryotes, the two ITS sequences are distinguished by such rapid evolutionary change that
105 they separate even lineages within species [5, 17]. rDNA markers thus offer taxonomic and
106 phylogenetic resolution at a very broad taxonomic scale. As an essential component of the
107 translation machinery, nuclear rDNA is required in large quantities in each cell. It is thus present
108 in multiple copies across the genome [15] and is readily accessible for PCR amplification. Due
109 to the above advantages, rDNA already is a popular and widely used marker for molecular
110 taxonomy and phylogenetics in many groups of organisms [5, 6, 15, 17, 18]. However, its
111 presence in multiple copies across the genome may also make rDNA susceptible to the
112 emergence of paralogs and pseudogenization, which could affect taxonomic and phylogenetic
113 utility.

114 Spanning about 8 kb, the ribosomal cluster is fairly large, and current barcoding protocols, e.g.
115 using Sanger sequencing or Illumina amplicon sequencing, can only target short sequence
116 stretches of 150 – 1,000 bp. Such short stretches of 28S and 18S are often too conserved to
117 identify young species pairs [19]. The ITS regions, on the other hand, are so variable that they
118 cannot be properly aligned across deeply divergent lineages. Moreover, ITS sequences can
119 show considerable length variation between taxa. This holds particularly true for the ITS1 region,
120 whose length can vary between few 100 and more than 1000 bp [20]. Considerable, but less
121 pronounced length variation can also be observed in ITS2. Short amplicon-based sequencing
122 approaches are limited to a maximum fragment length of about 500 bp. As ITS priming sites
123 have to rest in the conserved flanking rDNA gene sequences, the resulting amplicon often
124 exceeds this length and thus can not be used for short amplicon-based barcoding in some taxa.
125 Consequently, it would be ideal to amplify and sequence a large part of the ribosomal cluster in
126 one fragment. A solution to sequence the resulting long amplicons is offered by recent
127 developments in third generation sequencing platforms, which now enable researchers to
128 generate ultra-long reads, of up to 800 kb [21]. Recently, amplicons of several kilobases of the
129 rDNA cluster were sequenced using Pacific Bioscience (PacBio) technology, to explore fungal

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

130 community composition [22, 23]. With its circular consensus sequencing technology, PacBio
131 allows the generation of very accurate consensus reads. But while PacBio sequencing is well
132 suited for long amplicon sequencing, it is currently not readily available to every laboratory due
133 to the high cost and limited distribution of sequencing machines. PacBio sequencers are also
134 bulky and cannot be used outside of conventional laboratory settings.

135 A cost-efficient and readily available alternative is provided by nanopore sequencing
136 technology. The MinION sequencer (Oxford Nanopore Technologies) is small in size,
137 lightweight, allows for sequencing of several Gb's of DNA with average read lengths over 10 kb
138 on a single flow cell [24] and is available starting at \$1,000. Despite a raw read error rate of
139 about 12-22 % [21], highly accurate consensus sequences can be called from nanopore data
140 [25, 26], by assembling multiple sequences for individual specimens. The MinION is well suited
141 for amplicon sequencing, and a simple dual indexing strategy can be used to demultiplex
142 amplicon samples [27]. This technology offers tremendous potential for long-amplicon barcoding
143 applications, as recently shown in an analysis in fungi [26]. Oxford Nanopore Technologies'
144 MinION is a portable sequencer, and Nanopore based DNA barcoding can be applied with
145 mobile laboratories in remote sites outside of conventional labs (see e.g. [25, 28, 29]). However,
146 current analyses are still exploratory or limited in taxonomic focus, and streamlined analysis
147 pipelines to establish the method across the eukaryote tree of life are still missing.

148
149 Considering this background, we explore the feasibility of nanopore sequencing of long rDNA
150 amplicons as a simple, cost efficient DNA barcoding approach for animals and other eukaryote
151 taxa. We compile a workflow from PCR amplification, to library preparation, to demultiplexing
152 and consensus calling (see Fig. 1 for an overview). We explore the error profile of nanopore
153 consensus sequences and introduce MiniBar, a new software to demultiplex dual indexed
154 nanopore amplicon sequences. We test the utility of the ribosomal cluster for molecular
155 taxonomy and phylogenetics across divergent plant and animal taxa. A particular focus of our

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

156 analysis are arthropods, the most diverse group in the animal kingdom [30], which are highly
157 threatened by current mass extinctions [31]. Using a dataset of spiders, we compare the
158 taxonomic resolution of the ribosomal cluster with that offered by molecular barcoding using
159 mitochondrial COI, the currently preferred barcode marker for arthropods. Oxford Nanopore
160 Technologies' MinION is a portable sequencer, and Nanopore based DNA barcoding has been
161 applied in remote sites outside of conventional labs (see eg. [25, 30, 31]). As mentioned above,
162 the MinION is portable and can be used for DNA barcoding in field settings. Such field-based
163 applications confront researchers with additional complexities and challenges. To highlight the
164 simplicity of our approach, we tested it under field conditions and generated long rDNA barcode
165 sequences using a miniaturized mobile laboratory in a Peruvian rainforest.

166

167 We also tested the efficacy of long-amplicon rDNA sequencing for metabarcoding of bulk
168 community samples. A study of bacterial communities [32] suggests Nanopore long-amplicon
169 sequencing as a powerful tool for community characterization, but also found pronounced
170 biases in the recovered taxon abundance. Currently, little is known about the utility of long-
171 amplicon sequencing for animal community analysis. Metabarcoding protocols for community
172 samples need to be carefully optimized, as they can suffer from pronounced taxonomic biases,
173 e.g. due to primer binding or polymerase efficiency [33]. Well established Illumina based short
174 amplicon metabarcoding protocols can account for these biases and allow for a relatively good
175 qualitative and even quantitative recovery of taxa in communities [34]. However, additional, yet
176 unexplored, biases may affect long-amplicon metabarcoding. We thus also test the utility of
177 long-amplicon rDNA barcoding to recover taxonomic diversity from arthropod mock
178 communities. We compare the qualitative (species richness) and quantitative (species
179 abundance) recovery of taxa in simple mock communities by long-amplicon sequencing with
180 that based on short read Illumina amplicon sequencing of the 18SrDNA.

181

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

182 Overall, we demonstrate that long rDNA amplification and sequencing on the MinION platform is
183 a straightforward, cost effective, and universal approach for eukaryote DNA barcoding. It
184 combines robust taxonomic assignment power with high phylogenetic resolution and will enable
185 future analyses of taxonomic and phylogenetic diversity across wide taxonomic scales.

186
187

188 **Data Description and Analyses**

190 ***DNA extraction, PCR and library preparation***

191 We analyzed 114 specimens of eukaryotes including 17 insect and 42 spider species, two
192 annelid and nine plant species (Supplementary Table 1). Some feeder insects and the annelids
193 were purchased at a pet store. The remaining specimens were collected in oak forest on the
194 University of California Berkeley's campus or in native rainforests of the Hawaiian Archipelago
195 (under the Hawaii DLNR permit: FHM14-349). We particularly focused our arthropod sampling
196 on spiders, which are ubiquitous and essential predators in all terrestrial ecosystems. Recent
197 phylogenomic work [35] provided us with a solid baseline to test the efficiency of rDNA
198 amplicons for phylogenetic and taxonomic purposes. We included a taxonomically diverse
199 collection of 16 spider families from the Araneoidea, the RTA clade and a haplogyne outgroup
200 species. Within spiders, we additionally focused on the genus *Tetragnatha*, which has
201 undergone a striking adaptive radiation on Hawaii.

202
203 DNA was extracted from each sample using the Qiagen Archivepure kit (Qiagen, Valencia, CA,
204 USA) according to the manufacturer's protocols. The DNA integrity was checked on an agarose
205 gel. Only samples with high DNA integrity were used for the following PCRs. All DNA extracts
206 were quantified using a Qubit fluorometer using the high sensitivity dsDNA assay (Thermo
207 Fisher, Waltham, MA, USA) and diluted to concentrations of 20 ng/μl. We designed a primer

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

208 pair of each 27 bases to amplify a ~4,000 bp fragment of the ribosomal DNA, including partial
209 18S and 28S as well as full ITS1, 5.8S and ITS2 sequences (18S_F4
210 GGCTACCACATCYAARGAAGGCAGCAG and 28S_R8
211 TCGGCAGGTGAGTYGTRCACAYTCCT). The primers were designed using alignments of
212 partial 18S and 28S sequences of ~1,000 species of eukaryotes, with a focus on animals
213 (Supplementary Fig. 1). The primers targeted highly conserved regions across all analyzed
214 taxa. Degenerate sites were incorporated to account for variation. We aimed for high annealing
215 temperatures (65-70°C) to impose stringent amplification. These were calculated using the NEB
216 Tm Calculator (<https://tmcaculator.neb.com/#!/main>).

217
218 To index every PCR amplicon separately, we used a dual indexing strategy with each primer
219 carrying a unique 15 bp index sequence at its 5'-tail. Index sequences were designed using
220 Barcode Generator (http://comailab.genomecenter.ucdavis.edu/index.php/Barcode_generator)
221 with a minimum distance of 10 bases between each index. A total of 15 forward and 16 reverse
222 indexes were designed. Every sample was amplified separately using the Q5 Hot Start High-
223 Fidelity 2X Master Mix (NEB, Ipswich, MA, USA) in 15 µl reactions, at 68°C annealing
224 temperature, with 35 PCR cycles and using 50 ng of template DNA per PCR. All PCR products
225 were quantified on an agarose gel, based on band intensity on the gel, using the Gel Doc XR
226 System with the Quantity One software (Bio-Rad, CA, USA) and then pooled.

227
228 100 µl of the final pool were cleaned from residual primers by 0.75 X AMPure Beads XP
229 (Beckman Coulter, Brea, CA, USA). DNA library preparation was carried out according to the
230 1D PCR barcoding amplicons SQK- LSK108 protocol (Oxford Nanopore Technologies, Oxford,
231 UK). Barcoded DNA products were pooled with 5 µl of DNA CS (a positive control provided by
232 ONT) and an end-repair was performed (NEB-Next Ultra II End-prep reaction buffer and
233 enzyme mix), then purified using AMPure XP beads. Adapter ligation and tethering was carried

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

234 out with 20 µl Adapter Mix and 50 µl of NEB Blunt/TA ligation Master Mix. The adapter-ligated
235 DNA library was then purified with AMPure beads XP, followed by the addition of Adapter Bead
236 binding buffer, and finally eluted in 15 µl of Elution Buffer. Each R9 flow cell was primed with
237 1000 µl of a mixture of Fuel Mix and nuclease-free water. Twelve µl of the amplicon library were
238 diluted in 75 µL of running buffer with 35 µL RBF, 25.5 uL LLB, and 2.5 µL nuclease-free water
239 and then added to the flow cell via the SpotON sample port. The “NC_48Hr_sequencing_FLO-
240 MIN107_SQK- LSK108_plus_Basecaller.py” protocol was initiated using the MinION control
241 software, MinKNOW.

242

243 ***Field trial in the Amazon rainforest***

244 A field trial using the protocol described above was conducted in Tambopata, Peru, at the
245 Refugio Amazonas lodge (-12.874797, -69.409669) using two butterflies, a grasshopper, one
246 mosquito, unidentified insect eggs and two plant specimens. Collection permits in Peru were
247 issued by the Servicio Nacional Forestal y de Fauna Silvestre, 403-2016-SERFOR-
248 DGGSPFFS, 019-2017-SERFOR-DGGSPFFS. DNA extractions, PCR and library preparation
249 were performed in the field using a highly miniaturized laboratory consisting of portable
250 equipment. Equipment used for sequencing under remote tropical conditions is described in
251 further detail in Pomerantz, et al. [25]. DNA extractions were carried out with the Quick-DNA
252 Miniprep Plus Kit (Zymo Research, Irvine, CA, USA) according to manufacturer's protocol.
253 PCRs were performed using the Q5 Hot Start High-Fidelity 2X Master Mix and the same primers
254 as described above. A battery operated portable miniPCR device (Ampliyus, Cambridge, MA,
255 USA) was used to run PCRs. The sequencing on the MinION was carried out as described
256 above.

257

258 **Bioinformatics**

259

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

260 **Raw data processing and consensus calling**

261 The fastq files generated by the ONT software MinKNOW were de-multiplexed using MiniBar
262 (see description below), with index edit distances of 2, 3, and 4 and a primer edit distance of 11.
263 Next, the reads were filtered for quality (>13) and size (>3kb) using Nanofilt [36](
264 <https://github.com/wdecoster/nanofilt>). Individual consensus sequences were created using
265 Allele Wrangler (<https://github.com/transplantation-immunology/allele-wrangler/>) for
266 demultiplexed fastq files with a minimum coverage of 30. Error correction was performed using
267 RACON [37] (<https://github.com/isovic/racon>). To do so, we first mapped all the reads back to
268 the consensus using minimap (<https://github.com/lh3/minimap2>). We performed two cycles of
269 running minimap and RACON. Final consensus sequences were compared against the NCBI
270 database using BLASTn to check if the taxonomic assignment was correct.

271
272 We performed multiple tests to validate and optimize the consensus accuracy of long-amplicon
273 barcode sequences. To comparatively assess the accuracy, we used consensus sequences of
274 short 18S and 28SrDNA amplicons, which were previously generated using Illumina amplicon
275 sequencing for the 47 analyzed Hawaiian *Tetragnatha* specimens (Kennedy unpublished data).
276 These sequences were aligned with the respective stretches of our nanopore consensus
277 sequences using ClustalW in MEGA (MEGA Software , RRID:SCR_000667)[38]. All alignments
278 were then visually inspected and edited manually, where necessary. Pairwise distances
279 between Illumina and nanopore consensus were calculated in MEGA.

280
281 To measure consensus accuracy over the whole ribosomal amplicon, we utilized genome
282 skimming data [39] for six Hawaiian *Peperomia* plant species (Lim et al unpublished data). 150
283 bp paired-end TruSeq gDNA shotgun libraries for the six *Peperomia* samples were sequenced
284 on a single HiSeq v4000 lane (Illumina, San Diego, CA, USA). The resulting paired-end reads
285 were trimmed and filtered using Trimmomatic v0.36 (Trimmomatic , RRID:SCR_011848)[40]

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

286 and mapped to their respective nanopore consensus sequences using bowtie2 (Bowtie ,
287 RRID:SCR_005476)[41] under default parameter values and allowing for minimum and
288 maximum fragment size of 200 and 700 bases respectively. Mapping coverage of Illumina reads
289 to nanopore consensus sequences ranged between 150 - 600 X with a mean of ~ 300 X across
290 all six samples. We called Illumina read based consensus sequences for each *Peperomia*
291 species using bcftools [42], and aligned them with the previously generated nanopore
292 consensus sequences. Pairwise genetic distances were then calculated in MEGA as described
293 above. We performed two independent distance calculations: 1) excluding indels, i.e. only using
294 nucleotide substitutions to estimate genetic distances, and 2) including indels as additional
295 characters.

296
297 Our demultiplexing software allows flexible edit distances to identify forward and reverse
298 indexes from Nanopore reads. Due to the high raw read error rate, too large edit distances
299 could lead to crossover between samples during demultiplexing. This crossover could possibly
300 affect the accuracy of the called consensus sequence. On the other hand, too stringent edit
301 distances may result in very large read dropout. Assuming an average error rate of 12-22 %, 3
302 bp of our 15 bp indexes should maximize sequence recovery. We thus tested index edit
303 distances of 2, 3, and 4 bp in MiniBar for the six *Peperomia* specimens for which we had
304 generated Illumina based consensus sequences. We counted the number of recovered reads
305 and estimated the accuracy of the resulting consensus sequence based on the relevant edit
306 distances as described above.

307
308 A recent study [25] showed that accurate consensus sequences from nanopore data can be
309 generated using only 30x coverage. We tested 18 different assembly coverages from 10 to 800
310 sequences for a *Peperomia* species, to explore optimal assembly coverage. We randomly
311 subsampled the quality filtered and demultiplexed fastq file for the relevant specimen 10 times

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

312 for each tested assembly coverage. Consensus sequences were then assembled and genetic
313 distances to the Illumina consensus calculated as described above.

314
315 ***Phylogenetic and taxonomic analysis***

316 We carried out phylogenetic analyses on two hierarchical levels. First, we built a phylogeny for
317 all higher eukaryote taxa in our dataset, which included plants, animals and fungi. Second, we
318 took a closer look into the phylogeny of spiders. The resulting quality checked consensus
319 sequences of all taxa were aligned using ClustalW in MEGA. The alignments were visually
320 inspected and manually edited. The exact position of gene sequences was identified by
321 downloading full length 18S, 5.8S and 28S sequences from GenBank and then aligning them
322 against the amplicons. Due to the deep divergence in the eukaryote data set, the highly variable
323 ITS sequences could not be aligned and were excluded. For the analyses of spiders, we
324 retained both ITS sequences and aligned the whole rDNA amplicon. Appropriate models of
325 sequence evolution for each gene fragment of the rDNA cluster were identified using
326 PartitionFinder [43]. Phylogenies were built using MrBayes [44], with 4 heated chains, a chain
327 length of 1,100,000, subsampling every 200 generations and a burnin length of 100,000.

328
329 Focusing on the endemic Hawaiian *Tetragnatha* species, we also tested the utility of the
330 ribosomal cluster for taxonomic identification, as we also had COI barcodes available for these
331 species. Our dataset contained ribosomal DNA sequences for 47 specimens in 16 species,
332 which had been identified morphologically before barcoding. We calculated pairwise genetic
333 distances between and within all species for the whole ribosomal cluster and for each separate
334 gene region of the rDNA cluster using MEGA. As the 18S and 5.8S did not yield any species
335 level resolution within Hawaiian *Tetragnatha*, they were not analyzed separately. To compare
336 the taxonomic resolution of the ribosomal cluster with that of the commonly used mitochondrial
337 COI, we calculated inter- and intraspecific distances for an alignment of 418 bp of the COI

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

338 barcode region for the same spider specimens (Kennedy et al. unpublished data). We
339 performed a Mantel test using the R package ade4 [45] to test for a significant correlation
340 between COI and ribosomal DNA based distances. A comparison of intraspecific and
341 interspecific distances for mitochondrial COI and ribosomal DNA also allowed us to test for the
342 presence of a barcode gap.

343

344 ***Nanopore based arthropod metabarcoding***

345 To test for the possibility of estimating arthropod community composition from Nanopore
346 sequencing, we prepared four mock communities of different amounts of DNA extracts from 9
347 species of arthropods from different orders (see Supplementary Table 2). It should be noted that
348 with representatives of nine different orders, these community samples were highly simplified
349 and are not necessarily representative of a natural arthropod community. Due to the high error rate of
350 individual reads, we did not know if, and how, the MinION's high error rate would affect taxonomic
351 assignment, hence we decided to limit our current analysis to these simplified communities.

352 The samples were amplified using the Q5 High Fidelity Mastermix as described above at 68 °C
353 annealing temperature and 35 PCR cycles. We additionally tested two variations of PCR
354 conditions: 1) we either reduced the annealing temperature to 63 °C or, 2) reduced the PCR
355 cycle number to 25.

356 In order to compare our results with those from an optimized Illumina short read protocol, we
357 amplified all samples for a ~300 bp fragment of the 18S rDNA using the primer pair
358 18S2F/18S4R [46]. Amplification and library preparation were performed as described in [47]
359 using Qiagen Multiplex PCR kits. The 18S amplicon pools were sequenced on an Illumina
360 MiSeq using V3 chemistry and 2 x 300 bp reads. Sequence quality filtering, read merging and
361 primer trimming were performed as described in [34].

362

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

363 A library of 18S sequences for all included arthropod species (from [34]) was used as a
364 reference database to identify the recovered sequences using BLASTn [48], with a minimum e-
365 value of 10^{-4} and a minimum overlap of 95 %. Despite the high raw error rate of nanopore reads,
366 taxonomic status of sequences could be assigned using BLAST, as our pools contained
367 members of highly divergent orders. We compared the qualitative (number of species) and
368 quantitative (abundance of species) recovery of taxa from the communities by nanopore long-
369 amplicon and Illumina short read data. To estimate the recovery of taxon abundances, we
370 calculated a fold change between input DNA amount and recovered reads for each taxon and
371 mock community. A fold change of zero corresponded to a 1:1 association of taxon abundance
372 and read count, while positive or negative values indicated higher or lower read counts than the
373 taxon's actual abundance.

374

375 ***MiniBar***

376 We created a de-multiplexing software, called MiniBar. It allows customization of search
377 parameters to account for the high read error rates and has built-in awareness of the dual
378 barcode and primer pairs flanking the sequences. MiniBar takes as input a tab-delimited
379 barcode file and a sequence file in either fasta or fastq format. The barcode file contains, at a
380 minimum, sample name, forward barcode, forward primer, reverse barcode, and reverse primer
381 for each of the samples potentially in the sequence file. The software searches for barcodes and
382 for a primer, each permitting a user defined number of errors, an error being a mismatch or
383 indel. Error count to determine a match can either be a percentage of each of their lengths or
384 can be separately specified for barcode and primer as a maximum edit distance [49]. Output
385 options permit saving each sample in its own file or all samples in a single file, with the sample
386 names in the fasta or fastq headers. The found barcode primer pairs can be trimmed from the
387 sequence or can remain in the sequence distinguished by case or color. MiniBar, written in
388 Python 2.7, can also run in Python 3 and has the single dependency of the Edlib library module

for edit distance measured approximate search [50]. MiniBar can be found at <https://github.com/calacademy-research/minibar> along with test data.

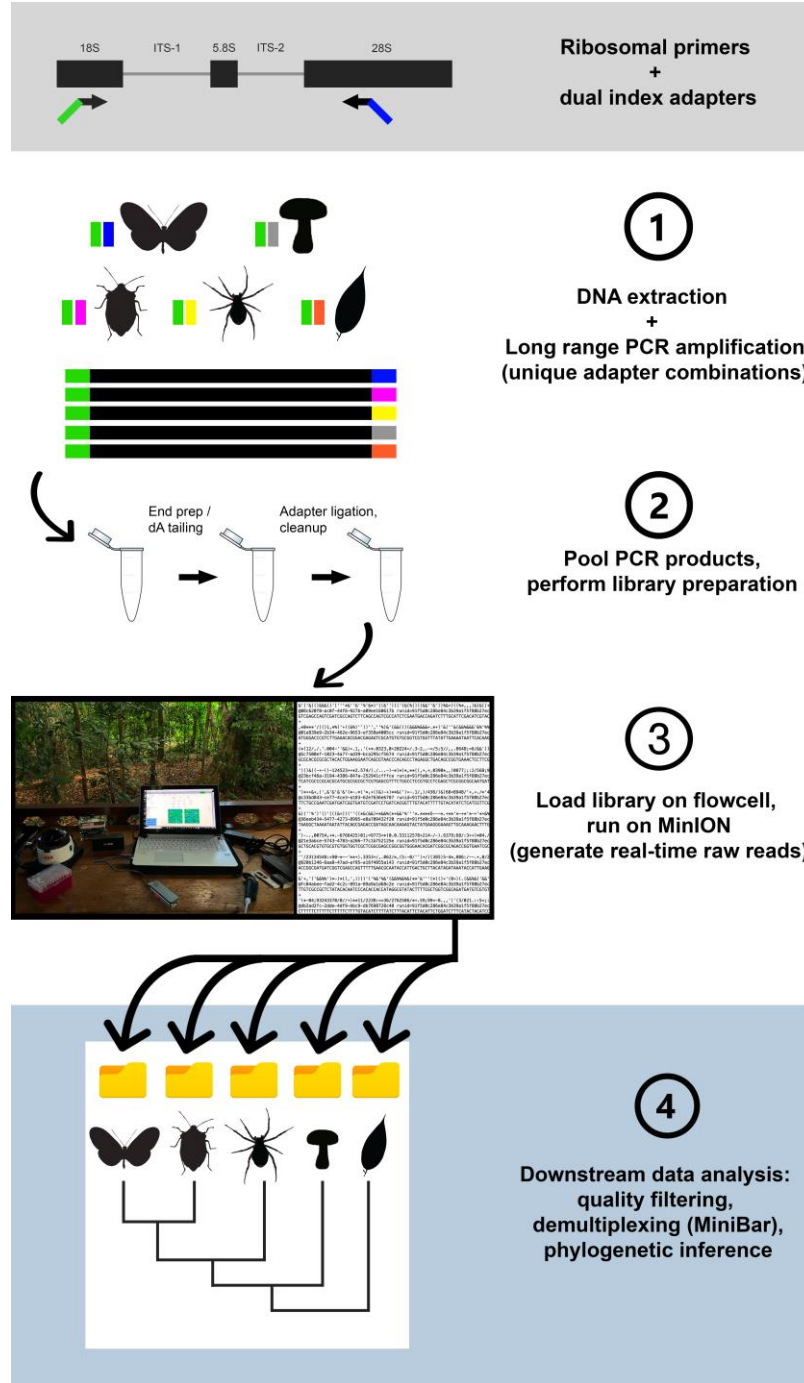


Figure 1. Workflow for the design, amplification, and sequencing of the ribosomal DNA cluster.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

395

396

397 **Results**

398

399 ***Sequencing, specimen recovery and consensus quality***

400 After quality filtering and trimming, our nanopore run yielded 245,433 reads. We tested edit
401 distances of two, three and four bases in MiniBar to demultiplex samples. Increasing edit
402 distances led to a significant increase in read numbers assigned to index combinations
403 (Pairwise Wilcoxon Test, FDR-corrected P -value < 0.05). On average, we found 355 reads per
404 specimen for an edit distance of two, 647 for a distance of three and 1,051 for a distance of four.
405 However, at an edit distance of four, we found a considerable increase of wrongly assigned
406 samples. A relatively high number of index combinations were incorrectly assigned at the
407 highest edit distance. Demultiplexed samples were then mixtures of different taxa, which
408 probably affected consensus accuracy. Using Illumina shotgun sequencing-derived consensus
409 sequences of rDNA from six *Peperomia* plants, we tested the accuracy of the nanopore
410 consensus assemblies based on the three edit distances (Fig. 2). While a distance of four
411 yielded the highest number of assigned reads (1,785 on average), it also led to slightly more
412 inaccurate consensus assemblies, with an average distance of 2.072 % to Illumina based
413 consensus sequences. We found a significant increase of consensus accuracy (Pairwise
414 Wilcoxon Test, FDR corrected $P < 0.05$) for edit distances of two (0.165 % average distance)
415 and three (0.187 % average distance). Despite significant differences in assigned reads (1,091
416 vs. 637 reads on average), there was not a significant difference in consensus accuracy of edit
417 distances of two versus three bases (Pairwise Wilcoxon Test, FDR corrected $P > 0.05$).

418

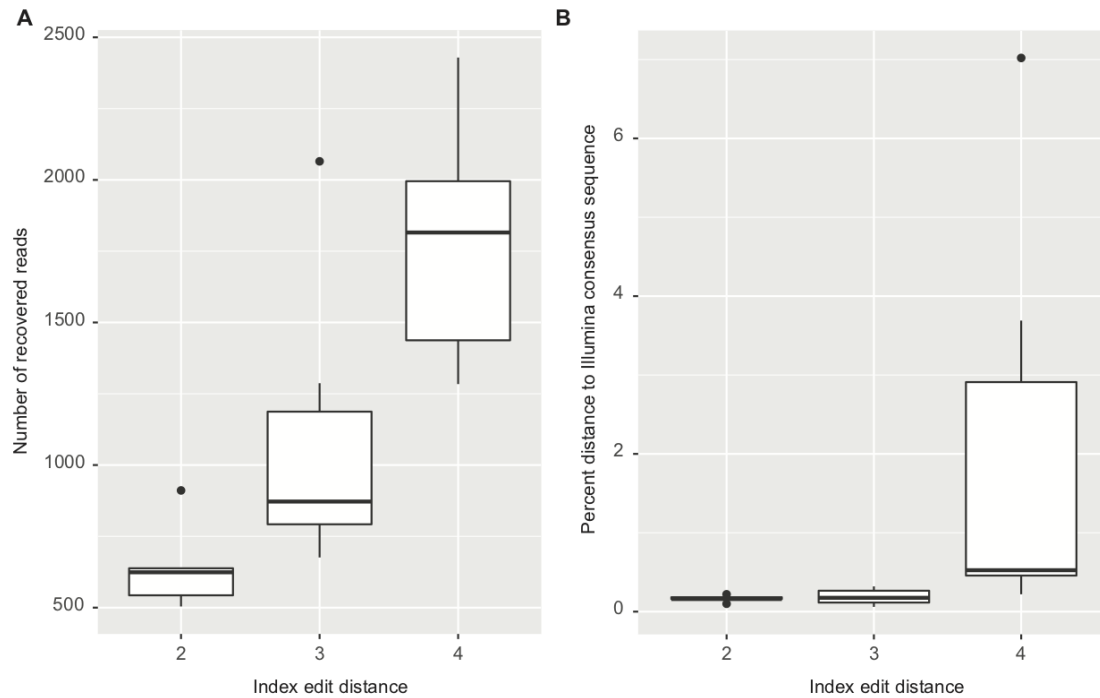


Figure 2: Comparison of recovered sequences and consensus accuracy for different index edit distances in Minibar. A) Number of recovered reads for six *Peperomia* species at index edit distances of two, three and four. B) Pairwise sequence divergence between Illumina and Nanopore based consensus sequences of the same six *Peperomia* specimens at the same index edit distances.

We chose a minimum coverage of 30 (see below) and an edit distance of two (which showed the smallest final consensus error rate) for all subsequent analyses. BLAST analyses suggested a correct taxonomic assignment for the majority of these consensus sequences. However, we found some notable exceptions. For two insect specimens, we amplified mite rDNA sequences. One of these specimens was *Drosophila hydei*, with the mite taxon being a well known phoretic associated with arthropods. A different mite taxon was assembled from an unidentified termite species. A species of isopod and a neuropteran yielded fungal sequences after assembly. The larva of a butterfly and a feeder mealworm (*Zophobas morio*) generated consensus sequences

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

434 for plants. In most of these samples, the targeted arthropod species was either extremely
435 underrepresented among the read populations or completely absent.

436 A comparison of our consensus sequences for 47 Hawaiian specimens of the spider genus
437 *Tetragnatha* with short Illumina amplicon sequencing-derived 18S and 28S rDNA sequences
438 suggests a very high consensus accuracy. Except for a single specimen, with a single
439 substitution error, all nanopore based consensus sequences were completely identical to the
440 Illumina based consensus. However, the corresponding 18S and 28S fragments did not contain
441 long stretches of homopolymer sequences, where nanopore raw read errors are known to
442 accumulate [51]. Despite containing several homopolymers, the nanopore derived *Peperomia*
443 consensus sequences were highly accurate (Supplementary Fig. 2). Including gaps in the
444 alignment, an average distance of 0.165 % to Illumina based consensus sequences was found.
445 Errors were clustered in indel regions (Supplementary Fig. 3). After excluding gaps, the average
446 distance dropped to 0.102 %.

447
448 We found only a small effect of sequence coverage on consensus assembly accuracy
449 (Supplementary Fig. 4). Even at 10-fold coverage, a low average distance of 0.257% to Illumina
450 consensus sequences was observed. However, at 20-fold coverage, the average distance
451 significantly decreased to 0.128 % (Pairwise Wilcoxon Test, FDR corrected $P < 0.05$). A slight,
452 but not significant, decrease of distance was observed with increasing coverage, with optimal
453 consensus accuracy at 300-fold coverage (0.031 % distance). At coverages larger than 300, the
454 consensus accuracy slightly decreased (average distance of 0.103 % at 800 X coverage), but
455 this change was not significant.

456
457 The length of the rDNA amplicon was quite variable between taxa. Arachnids, hexapods and
458 magnoliopsid plant specimens all showed a significantly different amplicon lengths (Pairwise
459 Wilcoxon Test, FDR corrected $P < 0.05$). The length difference was found for the actual gene

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

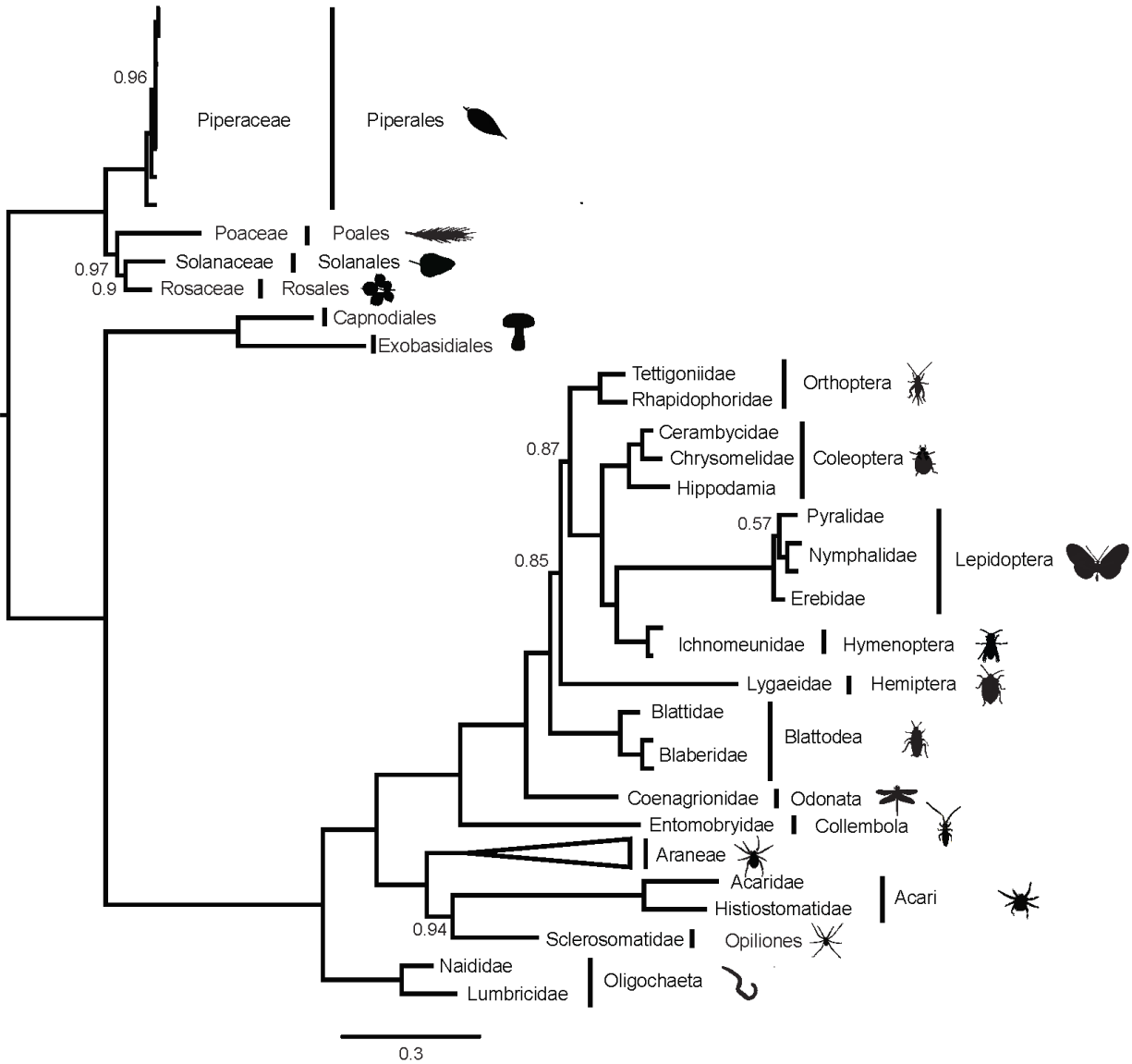
460 sequences (18S, 5.8S, 28S: plants: 2781 ± 4.96 ; hexapods: 3154 ± 50.35 ; arachnids: $3047 \pm$
461 10.77 %; Supplementary Fig. 5A) as well as including the ITS sequences (plants: 3243 ± 11.78 ;
462 hexapods: 4192 ± 498.05 ; arachnids: 3644 ± 129.07 , Supplementary Fig. 5B). While most
463 spiders showed very stable length distributions for the rDNA amplicon length (average length \pm
464 standard deviation across all Araneae: $3,629 \text{ bp} \pm 81$), several hexapod orders had rDNA
465 sequences of more variable length (Coleoptera: $4,488 \text{ bp} \pm 352$; Lepidoptera: $4363 \text{ bp} \pm 603$).

466
467 In contrast to the variable length of the rDNA cluster, we found a very stable GC content across
468 the whole taxonomic spectrum (46.75 ± 2.67 % across all taxa). GC content of magnoliopsid
469 plants, hexapods and arachnids was highly similar (plants: 46.01 ± 1.66 %; hexapods: $46.67 \pm$
470 3.73 %; arachnids: 46.93 ± 2.47 %) (Supplementary Fig 5c).

471
472

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

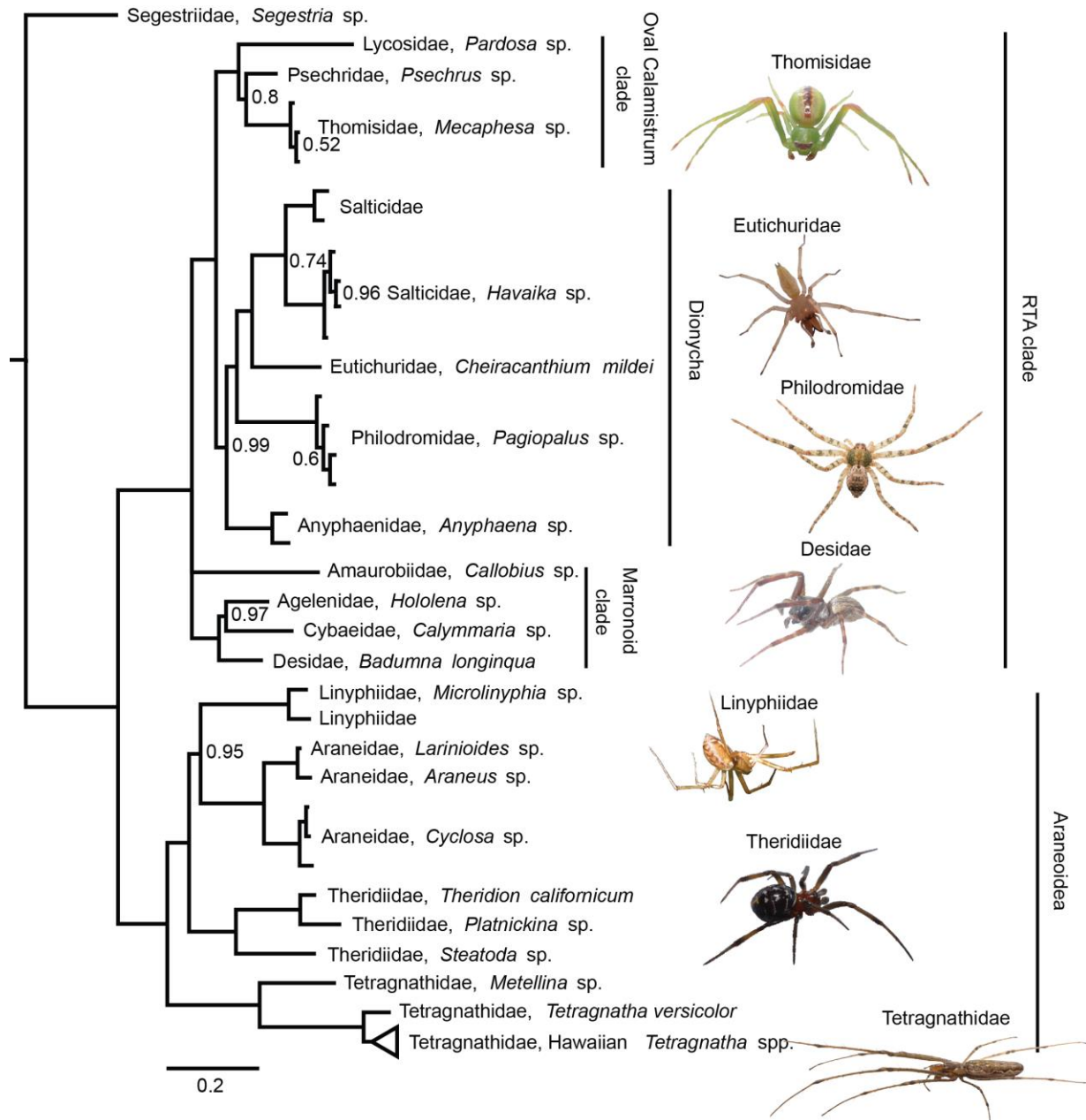
473 **Phylogenetic reconstruction**



475
476 **Figure 3 Bayesian consensus phylogeny based on a 3,656 bp alignment of 18S, 5.8S and**
477 **28S sequences of 117 animal, fungal and plant taxa.** The phylogeny is rooted using plants
478 as outgroup. Branches are annotated with family and order level taxonomy. The Araneae clade
479 of 83 specimens is collapsed. Only posterior probability values below 1 are displayed.

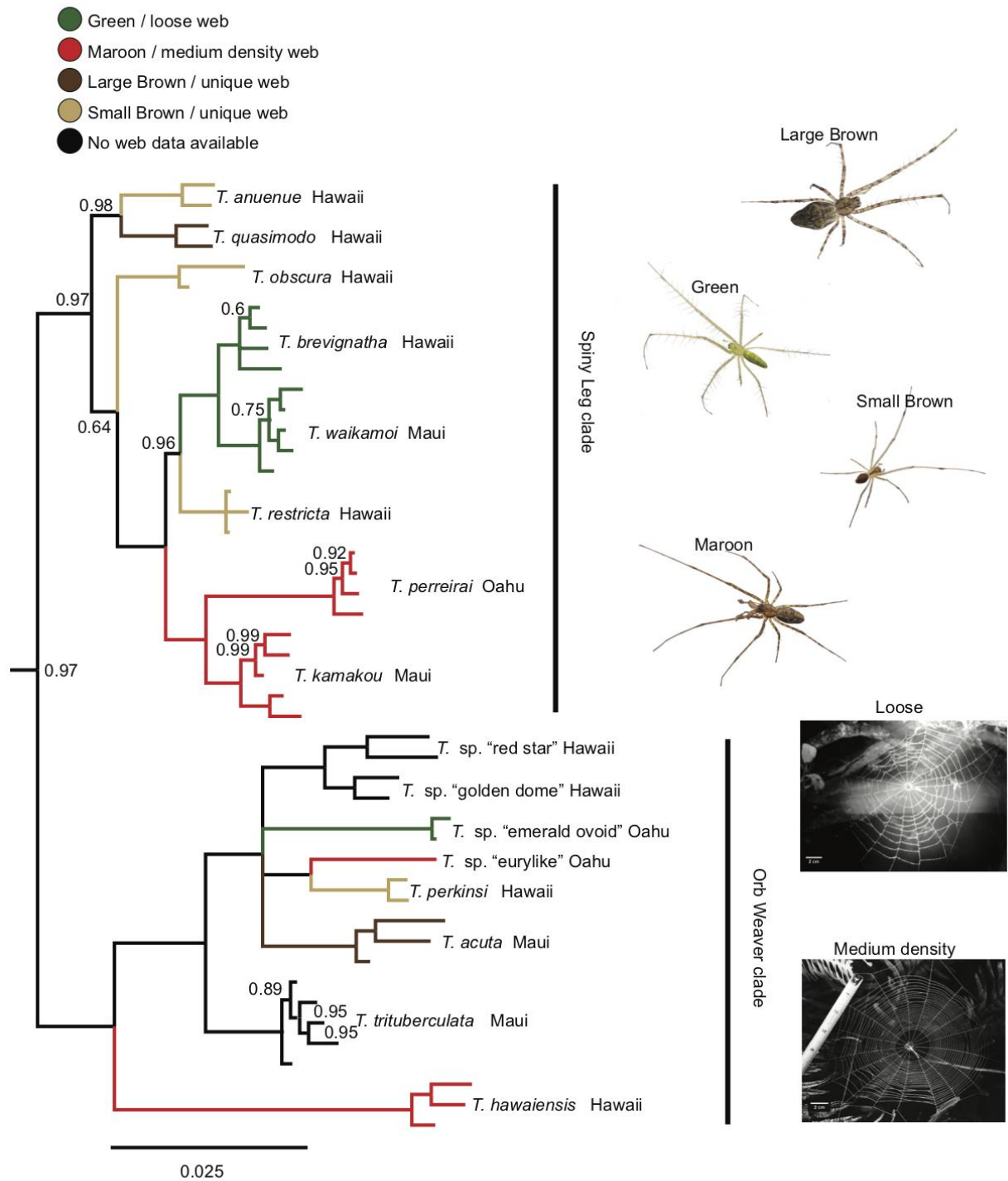
480

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



481
482 **Figure 4. Bayesian consensus phylogeny of 83 spiders in 16 families, based on a 4,214**
483 **bp alignment of 18S, ITS1, 5.8S, ITS2 and 28S.** The phylogeny is rooted using the basal
484 haplogyne *Segestria* sp. The clade containing Hawaiian members of the genus *Tetragnatha* is
485 collapsed (the uncollapsed clade is shown in Fig. 5). Only posterior probability values below 1
486 are displayed.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



487
488 **Figure 5. Section of the same phylogeny as Fig. 4, with expansion of the clade of 16**
489 **Hawaiian *Tetragnatha* species.** Different “Spiny Leg” ecomorphs and web architectures are
490 indicated by branch coloration. Only posterior probability values below 1 are displayed.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

491

492

493

494

495

496

497

498

499

500

501

502

503

504

We generated an alignment of 3,656 bp for 117 concatenated 18S, 5.8S and 28S sequences of plants, fungi, annelids and arthropods. Our phylogeny was well supported (most posterior support values equal one; Fig. 3). A basal split separated plants from fungi and animals. Within plants, the genus *Peperomia* was recovered as monophyletic. Fungi formed the sister group of animals. Within animals, annelids formed a separate clade from arthropods. Arthropods separated into arachnids and hexapods. Each separate arthropod order formed well supported groups. The hexapod phylogeny generally resembled that found in latest phylogenomic work [52]. The Collembola species *Salina* sp. formed the base to the insect tree, followed by the odonate *Argia* sp. A higher branch led to Blattodea, Hemiptera and Orthoptera. However, the support values for the relationships between these three orders were comparatively low (~0.85). Finally, holometabolan insects (Hymenoptera, Coleoptera and Lepidoptera) were recovered as monophyletic. The two Acari species, together with Opiliones, formed the sister clade to the monophyletic Araneae clade.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

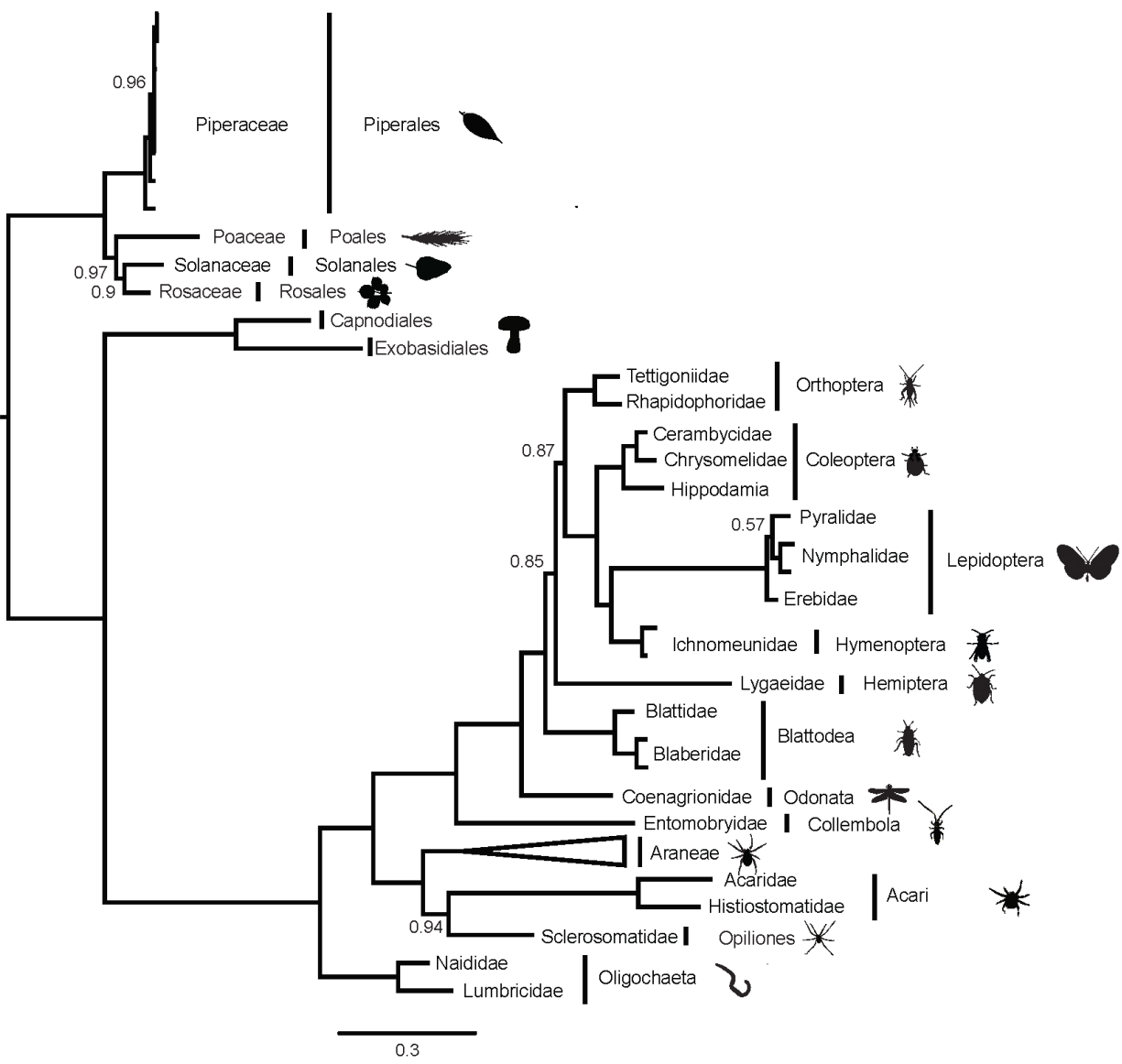


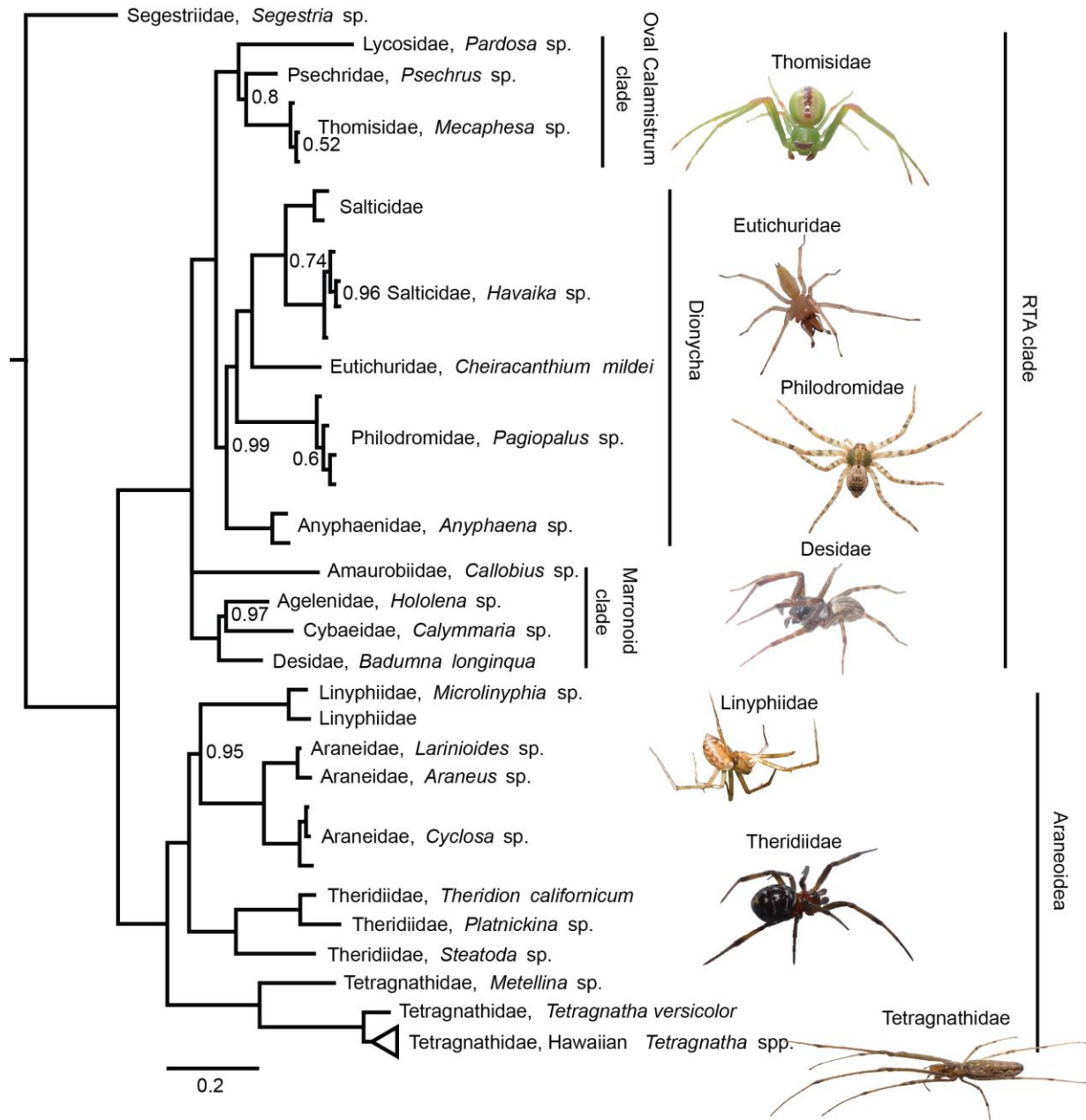
Figure 3 Bayesian consensus phylogeny based on a 3,656 bp alignment of 18S, 5.8S and 28S sequences of 117 animal, fungal and plant taxa. The phylogeny is rooted using plants as outgroup. Branches are annotated with family and order level taxonomy. The Araneae clade of 83 specimens is collapsed. Only posterior probability values below 1 are displayed.

Next, we generated a separate alignment of rDNA sequences for 83 spiders, including both ITS regions (totaling 4,214 bp). The spider phylogeny was also strongly supported (Fig. 4). Overall,

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

our phylogenetic tree topology agreed with the most recent phylogenetic work of [53] and [35].
With the haplogyne *Segestria* sp. (family Segestriidae) forming the root, we recovered the so-called RTA clade (represented in our dataset by families Agelenidae, Amaurobiidae, Anyphaenidae, Cybaeidae, Desidae, Eutichuridae, Lycosidae, Philodromidae, Psechridae, Salticidae and Thomisidae) and the Araneoidea (Araneidae, Linyphiidae, Tetragnathidae, Theridiidae) as two well supported monophyla. Within these clades, all families and genera formed well supported monophyletic groups. Similar to recent studies, we found the Marronoid clade as basal to the rest of the RTA clade; more derived clades were the Oval Calamistrum and the Dionycha clade. Inter-family relationships also closely matched those found in recent work: Lycosidae was basal to the clade formed by Psechridae and Thomisidae; Salticidae was closest to Eutichuridae and Philodromidae, with Anyphaenidae falling basal within Dionycha. Within Araneoidea, our results differed slightly from recent studies in that we recovered Tetragnathidae, rather than Theridiidae, as basal.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



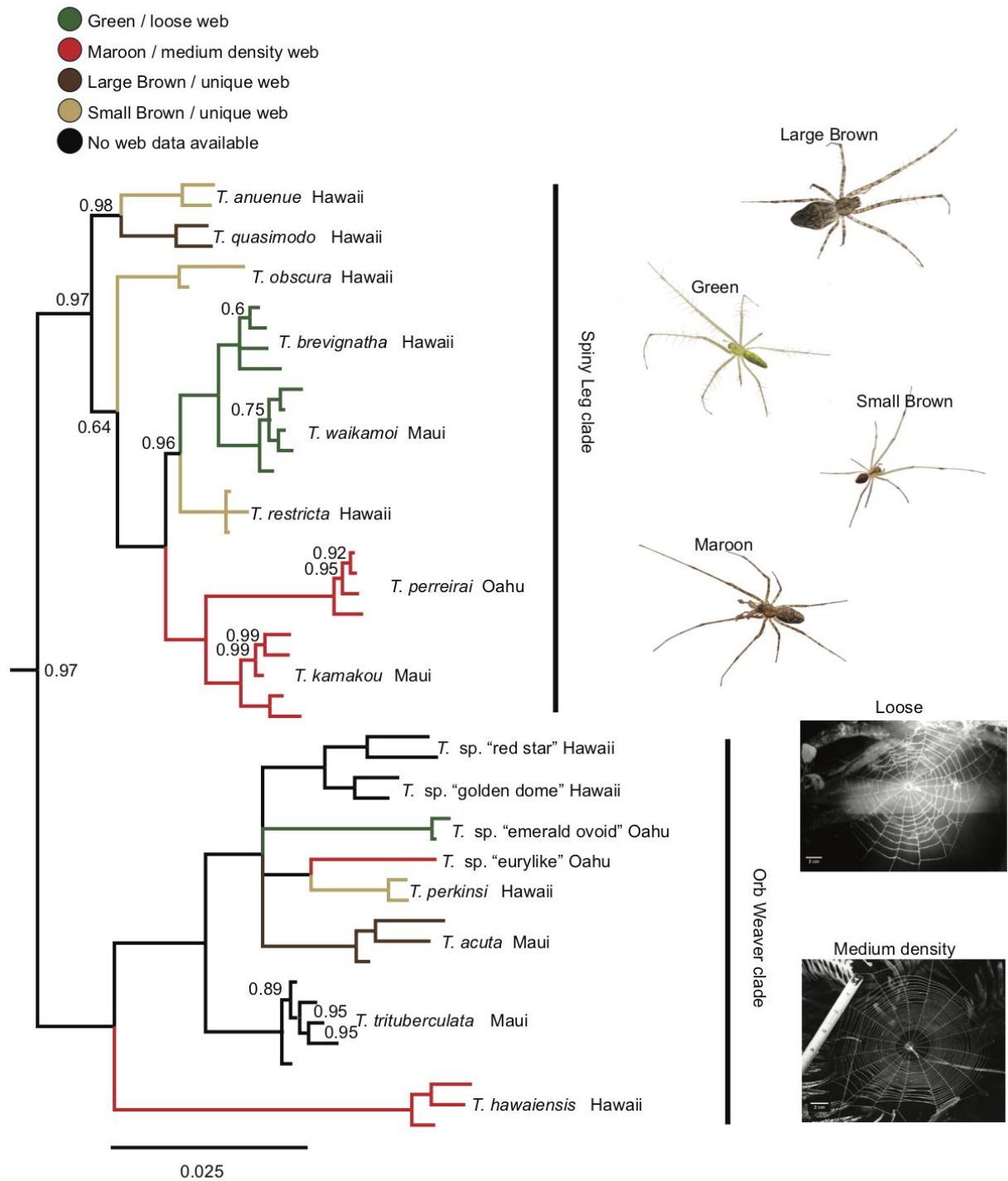
526
527 **Figure 4. Bayesian consensus phylogeny of 83 spiders in 16 families, based on a 4,214**
528 **bp alignment of 18S, ITS1, 5.8S, ITS2 and 28S.** The phylogeny is rooted using the basal
529 haplogyne *Segestria* sp. The clade containing Hawaiian members of the genus *Tetragnatha* is
530 collapsed (the uncollapsed clade is shown in Fig. 5). Only posterior probability values below 1
531 are displayed.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552

We recovered Hawaiian *Tetragnatha* as a well supported monophyletic clade within the Tetragnathidae. We found two main clades of Hawaiian *Tetragnatha* (Fig. 5), both of which have been supported by earlier work [54-57]: the orb weaving clade and the “Spiny Leg clade” of actively hunting species. All *Tetragnatha* species formed monophyletic groups, and the relationships among different species were mostly well supported. Within the Spiny Leg clade, species fell into one of four ecomorphs, each of which is associated with a particular substrate type [58]: “large brown” (*T. quasimodo*) with tree bark, “small brown” (*T. anuenue*, *T. obscura* and *T. restricta*) with twigs, “green” (*T. brevignatha* and *T. waikamoʻi*) with green leaves, and “maroon” (*T. perreirai* and *T. kamakou*) with lichen. While green and maroon ecomorphs clustered phylogenetically, small brown species appeared in three separate clades on the tree. Within the orb weaving clade, *T. hawaiiensis*, a generalist species which occurs on all of the Hawaiian Islands, fell basal. The characteristic web structures of some of these species have been documented [59, 60]. We found a pattern of apparent convergence in web structure for some species. *T. sp.* “emerald ovoid” spins a loose web with widely spaced rows of capture silk. *T. hawaiiensis* and *T. sp.* “eurylike,” which are distant relatives within the Hawaiian *Tetragnatha* clade, both spin webs of medium silk density, i.e. with more rows of capture silk per unit area than *T. sp.* “emerald ovoid.” *T. perkinsi* and *T. acuta* each spin a web structure that is not comparable in its silk density or size to any other known *Tetragnatha* species in this group [60], and are thus classified as “unique”.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



553
554 **Figure 5. Section of the same phylogeny as Fig. 4, with expansion of the clade of 16**
555 **Hawaiian *Tetragnatha* species.** Different “Spiny Leg” ecomorphs and web architectures are
556 indicated by branch coloration. Only posterior probability values below 1 are displayed.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

557
558
559 ***Taxonomic resolution***
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582

Our inferred genetic distances for rDNA sequences within and between Hawaiian *Tetragnatha* species were significantly correlated to those found for COI sequences of the same taxa ($R^2 = 0.70$, $P < 0.001$) (Fig. 6a). A Mantel test also suggested highly significant correlation of mitochondrial COI and nuclear rDNA-based distances (Mantel test, 9999 replicates, $P < 0.001$). Hence, the rDNA cluster supported a very similar pattern of genetic differentiation to COI. However, the faster evolutionary rate of COI was reflected in lower distances for the whole rDNA than for COI. Interspecific distances were significantly higher than intraspecific ones for COI and rDNA (Fig 6b,c). No overlap of intra- and interspecific distances was evident for COI, suggesting the presence of a barcode gap. A small overlap of intra- and interspecific distances was evident for the rDNA (Supplementary Table 3). However, this overlap was caused only by a single undescribed species (*T. sp.* “golden dome”) with unclear status, which showed a high intraspecific divergence in rDNA. Further morphological analyses will be necessary to rule out that the included samples do not actually comprise two species. At the same time, the interspecific rDNA distance of the relevant species was higher than its intraspecific distance. The lowest interspecific distance was found for a complex of closely related species from Maui. Like the combined rDNA cluster, genetic distances for different parts of the rDNA cluster all showed significant correlation with COI based distances, when analyzed separately (R^2 28S = 0.57, R^2 ITS1 = 0.68, R^2 ITS2 = 0.56, $P < 0.001$) (Supplementary Fig. 6). While the 28SrDNA showed considerably lower distances than COI, those for ITS1 and ITS2 were more comparable to COI (Supplementary Fig. 6b-d). Yet, interspecific and intraspecific distances for COI were significantly different from those for any part of the rDNA cluster (Pairwise Wilcoxon Test, FDR corrected $P < 0.05$).

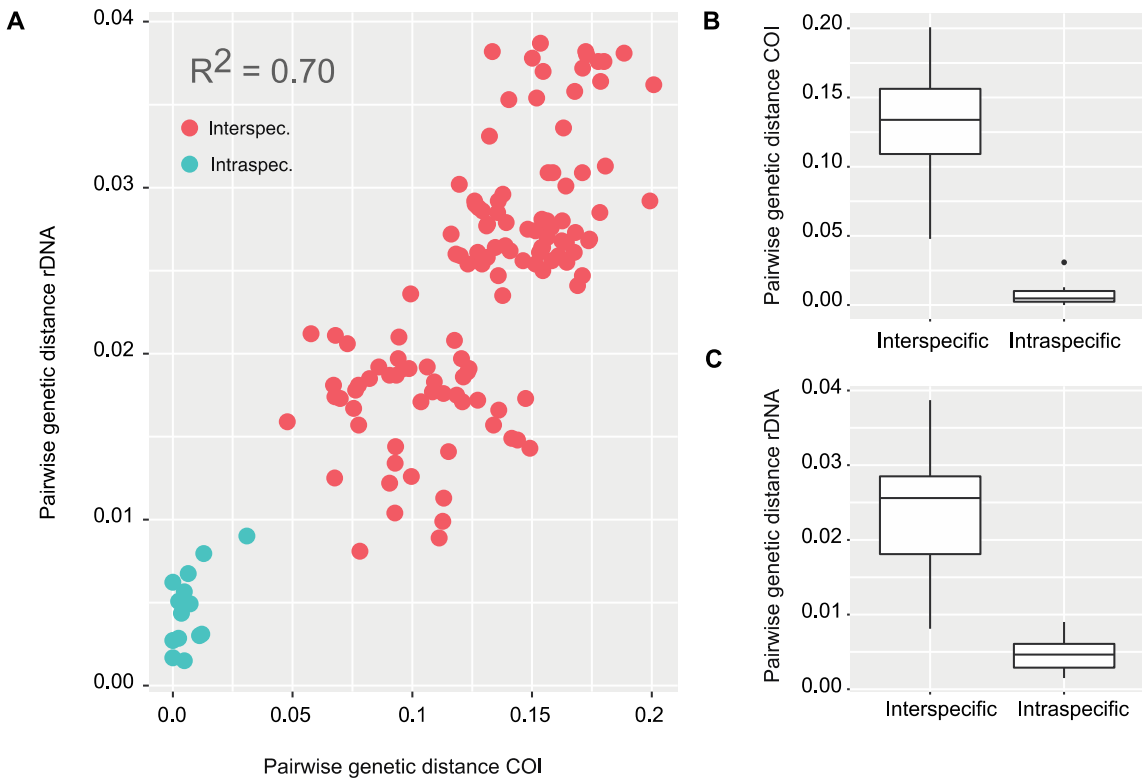


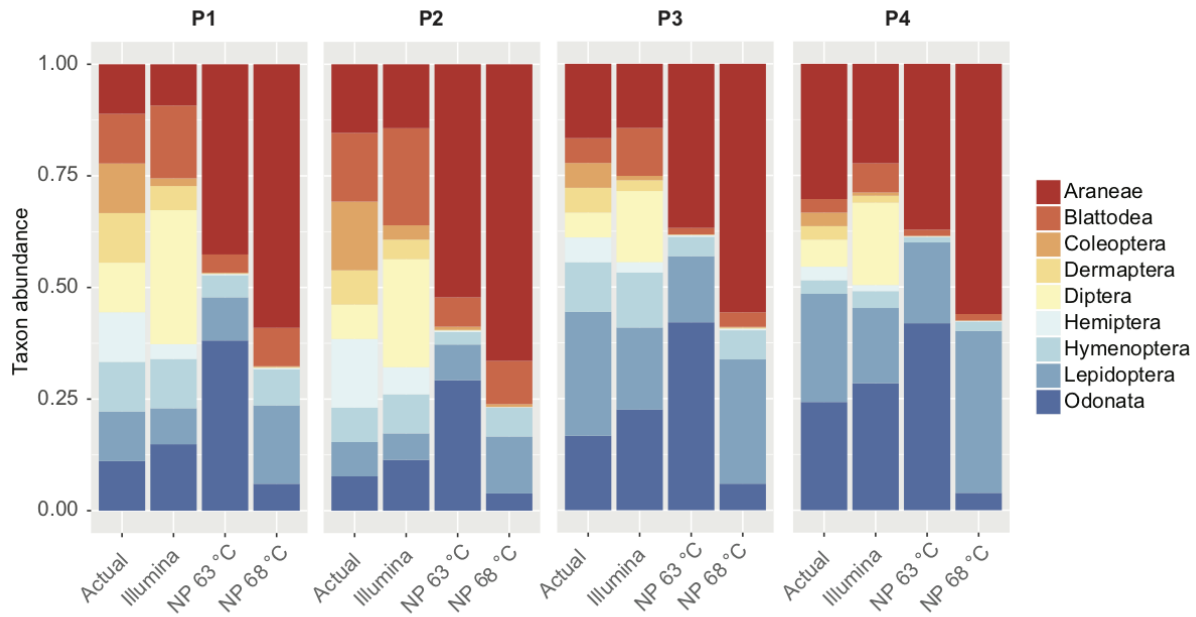
Figure 6 Inter and intraspecific genetic distances for the nuclear rDNA and mitochondrial COI for Hawaiian *Tetragnatha* spiders. A) Correlation of pairwise genetic distance between (red) and within (green) 16 Hawaiian *Tetragnatha* species based on COI and the full rDNA amplicon. B) Interspecific and intraspecific genetic distances for the same spider species based on mitochondrial COI and C) the whole rDNA amplicon.

Field trial in the Amazon rainforest

On March 26, 2018, we set out to test this method and a portable laboratory (as described in Pomerantz, et al. [25]) during an expedition to the Peruvian Amazon at the Refugio Amazonas Lodge (Supplementary Fig. 7). This field site is a “Terra firme” forest in the sector of “Condenado”, approximately two and a half hours by boat up river from the native community of

1
 2
 3
 4 597 Infierno on the buffer zone of the Tambopata National Reserve. We collected plant and insect
 5
 6 598 material, extracted DNA, amplified the rDNA cluster, and sequenced material on the MinION
 7
 8
 9 599 platform using the MinKNOW offline software (provided by ONT). The first run generated 17,149
 10
 11 600 reads and the second one 20,167 reads. We generated consensus sequences for five out of the
 12
 13 601 seven analyzed specimens. One plant sample and the grasshopper could not be assembled
 14
 15 602 due to too low read coverage. Moreover, BLAST analysis of the reads assigned to the
 16
 17 603 grasshopper suggested that we had sequenced a mite, instead of the grasshopper DNA. The
 18
 19
 20 604 unidentified insect eggs resulted in a butterfly consensus sequence, possibly a pierid species.
 21
 22 605
 23
 24 606

24 606 ***Nanopore based arthropod metabarcoding***



49 607
 50
 51 608 **Figure 7: Relative abundances for nine arthropod species in our four mock communities**
 52
 53 609 **(actual), compared to an Illumina amplicon sequencing protocol, and nanopore protocols**
 54
 55 610 **at 63 °C and 68 °C annealing temperature**
 56
 57
 58 611

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

612 On average, we recovered 2,645 reads for each Illumina sequenced mock community and
613 1,149 for each nanopore mock community. The optimized Illumina amplicon sequencing based
614 18SrDNA protocol resulted in a very good taxon recovery. All nine taxa were recovered from all
615 four mock communities (Fig. 7). Moreover, the Illumina based protocol allowed relatively
616 accurate predictions of taxon abundances. Even though no taxon's actual abundance was
617 predicted by Illumina amplicon data, the average fold change between input DNA and recovered
618 read count was closely distributed around zero (Supplementary Fig 8). In contrast, the long-
619 amplicon nanopore protocol showed very biased qualitative and quantitative taxon recovery
620 (Fig. 7). On average, only 83.33 % of taxa were recovered per nanopore sequenced mock
621 community. Moreover, the fold change of input DNA and recovered read count were highly
622 biased between taxa. Some taxa were considerably over or underrepresented among the read
623 population. This led to a significantly higher variation of fold change between input DNA and
624 read count compared to the Illumina amplicon-based protocol (Levene's test $P < 0.05$;
625 Supplementary Fig. 8). A reduction of PCR annealing temperature did result in a considerable
626 increase of Odonata sequences, but overall did not have a strong effect on qualitative (77.78 %
627 of taxa recovered) or quantitative taxon recovery (Fig. 7). The variation of fold change between
628 different PCR annealing temperatures was not significantly different (Levene's test, $P > 0.05$). A
629 reduction of PCR cycle number by 10 also did not yield any significant effect on qualitative
630 (88.89 % of taxa recovered) or quantitative taxon recovery (Supplementary Fig. 9).

631

632 **Discussion and Potential implications**

633

634 ***Phylogenetic and taxonomic utility of long rDNA amplicons***

635 Developments in long-amplicon sequencing hold great promise for molecular taxonomy and
636 phylogenetics across very broad taxonomic scales. We recovered phylogenetic relationships
637 across the eukaryote tree of life, which were mostly consistent with the current state of research

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

638 (e.g. [52]). Separate orders of arthropods all formed well supported monophyletic groups. Our
639 spider phylogeny was highly congruent with recent work based on whole transcriptomes [35]
640 and multi-amplicon data [53]. Moreover, using the rDNA cluster allowed us to resolve young
641 phylogenetic divergences: the relationships within the recent adaptive radiation of the genus
642 *Tetragnatha* in Hawaii confirmed previous research [58, 60].

643
644 Besides their high phylogenetic utility, long rDNA amplicons showed excellent support for
645 taxonomic hypotheses. All morphologically identified species of Hawaiian *Tetragnatha* were
646 recovered as monophyletic groups. The divergence patterns and taxonomic classifications of
647 spiders based on rDNA were strongly correlated to those based on mitochondrial COI, the most
648 commonly used animal barcode marker [4]. rDNA may thus be ideal to complement
649 mitochondrial barcoding. A universal and variable nuclear marker as a supplement to COI
650 barcoding will be particularly useful in cases of mito-nuclear discordance due to male biased
651 gene flow [10, 61], hybridization [12] or infections with reproductive parasites [11].

652
653 Their high phylogenetic utility across very broad taxonomic categories also provides long rDNA
654 amplicons with a distinct advantage over short read barcoding protocols, which are not well
655 suited to support broad scale phylogenetic hypotheses [62]. The inclusion of long amplicons
656 would make it possible to scale up barcoding from simple taxon assignment to community wide
657 phylogenetic inferences [9]. It should be noted that the nuclear rDNA cluster is a single locus
658 and its divergence pattern does not necessarily reflect species divergence. Also, the multiple
659 genomic rDNA copies do not necessarily all evolve in concert. rDNA genes may even be prone
660 to pseudogenization.

661 Taxonomic and phylogenetic analyses based on rDNA may thus be affected by paralogues, and
662 additional information from unlinked genomic regions would therefore be highly desirable to
663 support taxonomic and phylogenetic hypotheses. The mitochondrial genome may be an ideal

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

664 target for this purpose. Recently, the amplification of whole mitochondrial genomes was
665 suggested for animal barcoding [63]. This would increase taxonomic and phylogenetic
666 resolution and alleviate some disadvantages of short COI amplicons. However, it is challenging
667 to develop truly universal primers to target mitochondrial genomes across a very wide range of
668 taxonomic groups [64]. A straightforward way to achieve highly resolved phylogenies may be
669 the combination of long rDNA amplicon sequencing with multiplex PCRs of short mitochondrial
670 amplicons, to amplify multiple mitochondrial DNA fragments [65]. Conserved stretches in
671 mitochondrial rDNA may also allow the design of order- or even phylum-specific primers for long
672 range amplification [65]. A combination of long mitochondrial and nuclear rDNA amplicons,
673 possibly in a multiplex PCR, would be a desirable development for future DNA barcoding. With
674 whole genome sequences of different taxa rapidly accumulating, it may also be possible to
675 identify additional unlinked DNA barcoding markers.

676
677 *Simple, accurate, universal and cost efficient long-amplicon DNA barcoding*

678 Despite the high raw read error of nanopore data, consensus sequences were highly accurate,
679 and library preparation and sequencing for our protocol are simple and cost efficient. Using a
680 single pair of universal primers, long rDNA amplicons can potentially be amplified across
681 diverse eukaryote taxa, here largely demonstrated in arthropods, and in small scale in fungi and
682 plants. A simple dual indexing approach during PCR allows large numbers of samples to be
683 pooled before library preparation [27]. Only a single PCR is required per specimen, while
684 subsequent cleanup and library preparation can be performed on pooled samples. The
685 simplicity of our approach is additionally highlighted by its effectiveness even under field
686 conditions in a remote rainforest site. Nanopore sequencing technology is affordable and
687 universally available to any laboratory. Our ONT MinION generated about 250,000 reads per
688 run. Aiming for about 1,000 reads per amplified specimen, 250 long rDNA barcodes could be
689 generated in single MinION run. Input DNA amounts for different specimens will have to be

1
2
3
4 690 carefully balanced to maximize the recovery. The total reagent costs, including PCR, library
5
6 691 preparation and sequencing, then amount to less than \$4 for each long barcode sequence
7
8 692 generated.

9
10 693

11 694 ***Pitfalls of nanopore based long-amplicon barcoding***

12
13 695 While our protocol was generally straightforward and reliable, we found several drawbacks,
14
15 696 which require further considerations and optimization. First, it needs to be noted that long rDNA
16
17 697 amplification will not be possible with highly degraded DNA molecules, e.g. from historical
18
19 698 specimens [66]. Moreover, amplification success of long range PCRs proved less consistent
20
21 699 than that for amplification of short amplicons. We observed a complete failure of some PCRs
22
23 700 when too high template DNA concentrations were loaded. The long range polymerase may be
24
25 701 more sensitive to PCR inhibitors present in some arthropod DNA extractions [67]. PCR
26
27 702 conditions will have to be carefully optimized for reliable and consistent amplification. We also
28
29 703 found that highly universal eukaryote primers may result in undesired amplification, for example
30
31 704 plants from beetle and butterfly larval guts, phoretic mites, or fungal sequences. However, as
32
33 705 long as the DNA of the target taxon is still dominating the resulting amplicon mixture, this
34
35 706 undesired amplification will not affect consensus calling. It may be advisable to check the
36
37 707 taxonomic composition of amplicon samples before assembly, e.g. by blasting against a
38
39 708 reference library. To reduce non-target amplification, PCR primers could also be redesigned to
40
41 709 exclude certain lineages from amplification.

42
43 710 It should also be noted that our approach results in only a single consensus sequence
44
45 711 for each processed specimen. As a diploid marker, the rDNA cluster can contain heterozygous
46
47 712 positions in some specimens, in particular within the ITS regions. This information is currently
48
49 713 lost, and a different assembly approach may be necessary to recover heterozygosity as well.
50
51 714 Furthermore, index length and edit distance are also important considerations. We used indexes
52
53 715 of 15 bp and with a minimum distance of 10 bp to index both sides of our amplicons. Index edit
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

716 distance of only 4 bp between samples already led to considerable cross-specimen index
717 bleeding. It may thus be better to increase the length and edit distances of indexes. For
718 example, indexes of 20 or 30 bp could be easily attached to the 5'-tails of PCR primers. without
719 strongly affecting PCR efficiency. We have used a relatively crude gel-based approach for
720 pooling amplicon samples. This could have contributed to biased read abundance between
721 some samples. Instead of gel electrophoresis, it may be advisable to use a more precise
722 spectrophotometric quantification.

723

724 ***Nanopore based arthropod metabarcoding***

725 It is well known that Illumina amplicon sequencing of short 18SrDNA fragments can yield
726 accurate taxon recovery in metabarcoding experiments [34], a finding that is confirmed by our
727 results. Except for some outliers (e.g. *Diptera* were overrepresented), even the approximate
728 relative abundance of all taxa was recovered. In contrast, little is known on the performance of
729 long-amplicon nanopore sequencing for community diversity assessments [32]. Our long
730 barcode-based approach resulted in the dropout of several taxa and highly skewed relative
731 taxon abundances. Skewed abundances were already found in microbial community analysis
732 using nanopore [32]. In the simplest case, primer mismatches may be responsible for biased
733 amplification [32, 68]. However, the targeted priming sites in our study were extremely
734 conserved. Also, a change of PCR cycle number and annealing temperature did not have a
735 strong effect on taxon abundances, as would be expected in the case of PCR priming bias [69].
736 Another possibility is the preferential amplification of template molecules with a certain GC
737 content by the DNA polymerase [33]. However, we found the GC content of the rDNA cluster to
738 be very stable across taxa. Yet another potential explanation for the differential recovery of taxa
739 in community samples is taxonomic bias in DNA degradation [70], but we do not expect DNA
740 degradation to have played a role in our experiment because we used only high quality DNA
741 extractions (verified by gel electrophoresis) from fresh specimens. The most plausible

1
2
3
4 742 explanation appears to be that variable rDNA lengths are found between different taxa. It is well
5
6 743 known that shorter sequences are amplified preferentially in a PCR, especially after it reaches
7
8 744 the plateau stage [71]. Such dominance of shorter amplicons could explain the observed biases
9
10 745 very well. In fact, the most abundant taxon in our pools was a spider, which also had the
11
12 746 shortest amplicon length. The dominant amplification of shorter sequences may also explain the
13
14 747 amplification of plant DNA from a butterfly and a flour beetle larva, as plants showed
15
16 748 considerably shorter rDNA amplicons than insects. We found a very high variation of rDNA
17
18 749 amplicon length within many taxonomic groups, which could be a considerable problem for long
19
20 750 read metabarcoding applications. This suggests that it may be worthwhile to focus on narrower
21
22 751 taxonomic groups for long amplicon metabarcoding. For example, all spiders in our study share
23
24 752 rDNA amplicons of very similar size and would probably be less affected by amplification bias.
25
26 753 However, with more closely related taxa in a community, the high error rate of raw reads may
27
28 754 cause problems during read clustering and taxon assignments. It should also be noted that we
29
30 755 used highly simplified mock community samples, not reflecting actual community composition in
31
32 756 nature. Even with those simplified communities, we encountered considerable problems in
33
34 757 taxon recovery. Metabarcoding with MinION sequencing may thus be much less trivial than
35
36 758 single specimen sequencing. More research into the causes and possible mitigation of these
37
38 759 biases will be required before long-amplicon sequencing can be routinely utilized for
39
40 760 metabarcoding applications.
41
42
43
44
45
46
47
48

49 761
49 762 **Conclusion**

50
51 763 Sequencing long dual indexed rDNA amplicons on Oxford Nanopore Technologies' MinION is a
52
53 764 simple, cost effective, accurate and universal approach for eukaryote DNA barcoding. Long
54
55 765 rDNA amplicons offer high phylogenetic and taxonomic resolution across broad taxonomic
56
57 766 scales from kingdom down to species. They also prove to be an excellent complement to
58
59 767 mitochondrial COI based barcoding in arthropods. However, despite the long-amplicon
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

768 advantages in the analysis of separate specimens, we found considerable biases associated
769 with sequencing bulk community samples. The observed taxonomic bias is possibly a result of
770 taxon-specific length variation of the rDNA cluster and preferential amplification of species with
771 shorter rDNA. Further research into the sources of the observed bias is required before long
772 rDNA amplicon sequencing can be utilized as a reliable resource for the analysis of bulk
773 samples.

774

775 **Availability of source code and requirements**

776 1. The program Minibar can be found at <https://github.com/calacademy-research/minibar>

777 Programming language: Python 2.7 (but can be run in Python 3)

778 Operating systems: MacOS, Linux and Windows

779

780 Other requirements: Edlib library module (<https://github.com/Martinsos/edlib>)

781

782

783 **Availbity of supporting data**

784 The following data supporting the results of this article are available in the *GigaScience*
785 repository[72].

786

787 1. Raw fastq read files from Nanopore sequencing runs and Illumina sequencing of arthropod
788 mock communities for short 18S amplicons

789 2. Fasta sequences of rDNA amplicon for all taxa, mitochondrial COI for Hawaiian *Tetragnatha*
790 spp., as well as Illumina derived consensus sequences for Hawaiian *Peperomia* spp.

791 3. Newick tree files

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

792 4. Analysis tables for the mock community sequencing experiment, the comparison of genetic
793 distances within and between Hawaiian *Tetragnatha* species for COI and rDNA and the
794 distance between Nanopore based and Illumina based consensus sequences

795

796 **Author contributions**

797 HK and SP designed the study. HK, AP, SRK, JYL, NG, VS and JDS collected the specimens.
798 Laboratory work was carried out by HK, AP and SRK and the data were subsequently analyzed
799 by HK, AP, JBH, SRK and SP. The paper was written by HK, AP, JBH, SRK, JYL, VS, JDS, NG,
800 NHP, RGG, SP.

801

802

803 **Abbreviations**

804 COI: Cytochrome c oxidase subunit I; ONT: Oxford Nanopore Technologies; PCR: polymerase
805 chain reaction; rDNA: ribosomal DNA; RTA: retrolateral tibial apophysis

806

807

808 **Competing interests**

809 The authors declare that they have no competing interests.

810

811

812 **Acknowledgements**

813 We thank Taylor Liu for help during laboratory work, and Tara Gallant for help during specimen
814 collection. Hitomi Asahara graciously provided access to a laboratory facility and the necessary
815 software for our MinION sequencing run. We thank the State of Hawaii Department of Land and
816 Natural Resources and the Servicio Nacional Forestal y de Fauna Silvestre, who provided
817 collection permits, rainforest Expeditions and Gabriela Orihuela for providing assistance and

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

818 support with fieldwork in Peru. We thank Anna Holmquist for providing the *Psechrus* sp.
819 specimen. The specimen was collected with permits from the Indonesian Ministry of Research
820 and Technology (KEMENRISTEK) and in collaboration with Pungki Lupiyaningdyah and Anang
821 Achmadi from the Museum Zoologicum Bogoriense and funded by an NSF grant DEB
822 1457845.

823
824

References

- 826 1. Sala, O.E., Chapin, F.S., Armesto, J.J., Berlow, E., Bloomfield, J., Dirzo, R., Huber-
827 Sanwald, E., Huenneke, L.F., Jackson, R.B., and Kinzig, A. (2000). Global biodiversity
828 scenarios for the year 2100. *Science* 287, 1770-1774.
- 829 2. Pimm, S.L., Jenkins, C.N., Abell, R., Brooks, T.M., Gittleman, J.L., Joppa, L.N., Raven,
830 P.H., Roberts, C.M., and Sexton, J.O. (2014). The biodiversity of species and their rates
831 of extinction, distribution, and protection. *Science* 344, 1246752.
- 832 3. Rominger, A., Goodman, K., Lim, J., Armstrong, E., Becking, L., Bennett, G., Brewer, M.,
833 Cotoras, D., Ewing, C., and Harte, J. (2016). Community assembly on isolated islands:
834 macroecology meets evolution. *Global ecology and biogeography* 25, 769-780.
- 835 4. Hebert, P.D., Ratnasingham, S., and de Waard, J.R. (2003). Barcoding animal life:
836 cytochrome c oxidase subunit 1 divergences among closely related species.
837 *Proceedings of the Royal Society of London B: Biological Sciences* 270, S96-S99.
- 838 5. Schoch, C.L., Seifert, K.A., Huhndorf, S., Robert, V., Spouge, J.L., Levesque, C.A.,
839 Chen, W., Bolchacova, E., Voigt, K., and Crous, P.W. (2012). Nuclear ribosomal internal
840 transcribed spacer (ITS) region as a universal DNA barcode marker for Fungi.
841 *Proceedings of the National Academy of Sciences* 109, 6241-6246.
- 842 6. China Plant BOL Group, Li, D.-Z., Gao, L.-M., Li, H.-T., Wang, H., Ge, X.-J., Liu, J.-Q.,
843 Chen, Z.-D., Zhou, S.-L., and Chen, S.-L. (2011). Comparative analysis of a large

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

dataset indicates that internal transcribed spacer (ITS) should be incorporated into the core barcode for seed plants. *Proceedings of the National Academy of Sciences* 108, 19641-19646.

7. Shokralla, S., Porter, T.M., Gibson, J.F., Dobosz, R., Janzen, D.H., Hallwachs, W., Golding, G.B., and Hajibabaei, M. (2015). Massively parallel multiplex DNA sequencing for specimen identification using an Illumina MiSeq platform. *Scientific reports* 5, 9687.

8. Yu, D.W., Ji, Y., Emerson, B.C., Wang, X., Ye, C., Yang, C., and Ding, Z. (2012). Biodiversity soup: metabarcoding of arthropods for rapid biodiversity assessment and biomonitoring. *Methods in Ecology and Evolution* 3, 613-623.

9. Graham, C.H., and Fine, P.V. (2008). Phylogenetic beta diversity: linking ecological and evolutionary processes across space in time. *Ecology Letters* 11, 1265-1277.

10. Krehenwinkel, H., Graze, M., Rödder, D., Tanaka, K., Baba, Y.G., Muster, C., and Uhl, G. (2016). A phylogeographical survey of a highly dispersive spider reveals eastern Asia as a major glacial refugium for Palaearctic fauna. *Journal of Biogeography* 43, 1583-1594.

11. Hurst, G.D., and Jiggins, F.M. (2005). Problems with mitochondrial DNA as a marker in population, phylogeographic and phylogenetic studies: the effects of inherited symbionts. *Proceedings of the Royal Society of London B: Biological Sciences* 272, 1525-1534.

12. Bernatchez, L., Glémet, H., Wilson, C.C., and Danzmann, R.G. (1995). Introgression and fixation of Arctic char (*Salvelinus alpinus*) mitochondrial genome in an allopatric population of brook trout (*Salvelinus fontinalis*). *Canadian Journal of Fisheries and Aquatic Sciences* 52, 179-185.

13. Melo- Ferreira, J., Boursot, P., Suchentrunk, F., Ferrand, N., and Alves, P. (2005). Invasion from the cold past: extensive introgression of mountain hare (*Lepus timidus*) mitochondrial DNA into three other hare species in northern Iberia. *Molecular Ecology* 14, 2459-2464.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

870 14. Soltis, P.S., and Soltis, D.E. (1998). Molecular evolution of 18S rDNA in angiosperms:
871 implications for character weighting in phylogenetic analysis. In *Molecular systematics of*
872 *plants II*. (Springer), pp. 188-210.

873 15. Hillis, D.M., and Dixon, M.T. (1991). Ribosomal DNA: molecular evolution and
874 phylogenetic inference. *The Quarterly review of biology* 66, 411-453.

875 16. Black IV, W.C., Klompen, J., and Keirans, J.E. (1997). Phylogenetic relationships among
876 tick subfamilies (Ixodida: Ixodidae: Argasidae) based on the 18S nuclear rDNA gene.
877 *Molecular Phylogenetics and Evolution* 7, 129-144.

878 17. Powers, T.O., Todd, T., Burnell, A., Murray, P., Fleming, C., Szalanski, A.L., Adams, B.,
879 and Harris, T. (1997). The rDNA internal transcribed spacer region as a taxonomic
880 marker for nematodes. *Journal of Nematology* 29, 441.

881 18. Sonnenberg, R., Nolte, A.W., and Tautz, D. (2007). An evaluation of LSU rDNA D1-D2
882 sequences for their use in species identification. *Frontiers in zoology* 4, 6.

883 19. Tang, C.Q., Leasi, F., Obertegger, U., Kieneker, A., Barraclough, T.G., and Fontaneto, D.
884 (2012). The widely used small subunit 18S rDNA molecule greatly underestimates true
885 diversity in biodiversity surveys of the meiofauna. *Proceedings of the National Academy*
886 *of Sciences* 109, 16208-16212.

887 20. von der Schulenburg, J.H.G., Hancock, J.M., Pagnamenta, A., Sloggett, J.J., Majerus,
888 M.E., and Hurst, G.D. (2001). Extreme length and length variation in the first ribosomal
889 internal transcribed spacer of ladybird beetles (Coleoptera: Coccinellidae). *Molecular*
890 *Biology and Evolution* 18, 648-660.

891 21. Jain, M., Koren, S., Miga, K.H., Quick, J., Rand, A.C., Sasani, T.A., Tyson, J.R., Beggs,
892 A.D., Dilthey, A.T., and Fiddes, I.T. (2018). Nanopore sequencing and assembly of a
893 human genome with ultra-long reads. *Nature biotechnology* 36, 338.

894 22. Heeger, F., Bourne, E.C., Baschien, C., Yurkov, A., Bunk, B., Spröer, C., Overmann, J.,
895 Mazzoni, C.J., and Monaghan, M.T. (2018). Long-amplicon DNA metabarcoding of

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

896 ribosomal rRNA in the analysis of fungi from aquatic environments. *Molecular Ecology*
897 *Resources*.

898 23. Tedersoo, L., Tooming-Klunderud, A., and Anslan, S. (2018). PacBio metabarcoding of
899 Fungi and other eukaryotes: errors, biases and perspectives. *New Phytologist* 217,
900 1370-1385.

901 24. Giordano, F., Aigrain, L., Quail, M.A., Coupland, P., Bonfield, J.K., Davies, R.M.,
902 Tischler, G., Jackson, D.K., Keane, T.M., and Li, J. (2017). De novo yeast genome
903 assemblies from MinION, PacBio and MiSeq platforms. *Scientific Reports* 7, 3935.

904 25. Pomerantz, A., Peñafiel, N., Arteaga, A., Bustamante, L., Pichardo, F., Coloma, L.A.,
905 Barrio-Amorós, C.L., Salazar-Valenzuela, D., and Prost, S. (2018). Real-time DNA
906 barcoding in a rainforest using nanopore sequencing: opportunities for rapid biodiversity
907 assessments and local capacity building. *Gigascience*. 2018 Apr 1;7(4). doi:
908 10.1093/gigascience/giy033.

909 26. Wurzbacher, C., Larsson, E., Bengtsson-Palme, J., Van den Wyngaert, S., Svantesson,
910 S., Kristiansson, E., Kagami, M., and Nilsson, R.H. (2018). Introducing ribosomal
911 tandem repeat barcoding for fungi. *bioRxiv*, 310540.

912 27. Srivathsan, A., Baloğlu, B., Wang, W., Tan, W.X., Bertrand, D., Ng, A.H., Boey, E.J.,
913 Koh, J.J., Nagarajan, N., and Meier, R. (2018). A Min ION™- based pipeline for fast and
914 cost- effective DNA barcoding. *Molecular ecology resources*.

915 28. Quick, J., Loman, N.J., Duraffour, S., Simpson, J.T., Severi, E., Cowley, L., Bore, J.A.,
916 Koundouno, R., Dudas, G., and Mikhail, A. (2016). Real-time, portable genome
917 sequencing for Ebola surveillance. *Nature* 530, 228.

918 29. Edwards, A., Debonnaire, A.R., Sattler, B., Mur, L.A., and Hodson, A.J. (2016). Extreme
919 metagenomics using nanopore DNA sequencing: a field report from Svalbard, 78 N.
920 *bioRxiv*, 073965.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

921 30. Giribet, G., and Edgecombe, G.D. (2012). Reevaluating the arthropod tree of life. *Annual*
922 *review of entomology* 57, 167-186.

923 31. Hochkirch, A. (2016). The insect crisis we can't ignore. *Nature News* 539, 141.

924 32. Benítez-Páez, A., Portune, K.J., and Sanz, Y. (2016). Species-level resolution of 16S
925 rRNA gene amplicons sequenced through the MinION™ portable nanopore sequencer.
926 *GigaScience* 5, 4.

927 33. Nichols, R.V., Vollmers, C., Newsom, L.A., Wang, Y., Heintzman, P.D., Leighton, M.,
928 Green, R.E., and Shapiro, B. (2018). Minimizing polymerase biases in metabarcoding.
929 *Molecular ecology resources*.

930 34. Krehenwinkel, H., Wolf, M., Lim, J.Y., Rominger, A.J., Simison, W.B., and Gillespie, R.G.
931 (2017). Estimating and mitigating amplification bias in qualitative and quantitative
932 arthropod metabarcoding. *Scientific reports* 7, 17668.

933 35. Fernández, R., Kallal, R.J., Dimitrov, D., Ballesteros, J.A., Arnedo, M.A., Giribet, G., and
934 Hormiga, G. (2018). Phylogenomics, Diversification Dynamics, and Comparative
935 Transcriptomics across the Spider Tree of Life. *Current Biology* 28, 1489-1497. e1485.

936 36. De Coster, W., D'Hert, S., Schultz, D.T., Cruts, M., and Van Broeckhoven, C. (2018).
937 NanoPack: visualizing and processing long read sequencing data. *bioRxiv*, 237180.

938 37. Vaser, R., Sović, I., Nagarajan, N., and Šikić, M. (2017). Fast and accurate de novo
939 genome assembly from long uncorrected reads. *Genome Research* 27, 737-746.

940 38. Tamura, K., Stecher, G., Peterson, D., Filipowski, A., and Kumar, S. (2013). MEGA6:
941 molecular evolutionary genetics analysis version 6.0. *Molecular Biology and Evolution*
942 30, 2725-2729.

943 39. Straub, S.C., Parks, M., Weitemier, K., Fishbein, M., Cronn, R.C., and Liston, A. (2012).
944 Navigating the tip of the genomic iceberg: Next-generation sequencing for plant
945 systematics. *American Journal of Botany* 99, 349-364.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

946 40. Bolger, A., and Giorgi, F. Trimmomatic: A Flexible Read Trimming Tool for Illumina NGS
947 Data. URL <http://www.usadellab.org/cms/index.php>.

948 41. Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows–
949 Wheeler transform. *Bioinformatics* 25, 1754-1760.

950 42. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,
951 Abecasis, G., and Durbin, R. (2009). The sequence alignment/map format and
952 SAMtools. *Bioinformatics* 25, 2078-2079.

953 43. Lanfear, R., Calcott, B., Ho, S.Y., and Guindon, S. (2012). PartitionFinder: combined
954 selection of partitioning schemes and substitution models for phylogenetic analyses.
955 *Molecular Biology and Evolution* 29, 1695-1701.

956 44. Huelsenbeck, J.P., and Ronquist, F. (2001). MRBAYES: Bayesian inference of
957 phylogenetic trees. *Bioinformatics* 17, 754-755.

958 45. Dray, S., and Dufour, A.-B. (2007). The ade4 package: implementing the duality diagram
959 for ecologists. *Journal of statistical software* 22, 1-20.

960 46. Machida, R.J., and Knowlton, N. (2012). PCR primers for metazoan nuclear 18S and
961 28S ribosomal DNA sequences. *PLoS one* 7, e46180.

962 47. Krehenwinkel, H., Kennedy, S., Pekár, S., and Gillespie, R.G. (2017). A cost- efficient
963 and simple protocol to enrich prey DNA from extractions of predatory arthropods for
964 large- scale gut content analysis by Illumina sequencing. *Methods in Ecology and*
965 *Evolution* 8, 126-134.

966 48. Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local
967 alignment search tool. *Journal of molecular biology* 215, 403-410.

968 49. Levenshtein, V.I. (1966). Binary codes capable of correcting deletions, insertions, and
969 reversals. In *Soviet physics doklady*, Volume 10. pp. 707-710.

970 50. Šošić, M., and Šikić, M. (2017). Edlib: a C/C++ library for fast, exact sequence alignment
971 using edit distance. *Bioinformatics* 33, 1394-1395.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

51. Loman, N.J., Quick, J., and Simpson, J.T. (2015). A complete bacterial genome assembled de novo using only nanopore sequencing data. *bioRxiv*, 015552.

52. Misof, B., Liu, S., Meusemann, K., Peters, R.S., Donath, A., Mayer, C., Frandsen, P.B., Ware, J., Flouri, T., and Beutel, R.G. (2014). Phylogenomics resolves the timing and pattern of insect evolution. *Science* *346*, 763-767.

53. Wheeler, W.C., Coddington, J.A., Crowley, L.M., Dimitrov, D., Goloboff, P.A., Griswold, C.E., Hormiga, G., Prendini, L., Ramírez, M.J., and Sierwald, P. (2017). The spider tree of life: phylogeny of Araneae based on target- gene analyses from an extensive taxon sampling. *Cladistics* *33*, 574-616.

54. Gillespie, R.G. (1991). Hawaiian spiders of the genus *Tetragnatha*: I. Spiny leg clade. *Journal of Arachnology*, 174-209.

55. Gillespie, R.G. (1999). Comparison of rates of speciation in web-building and non-web-building groups within a Hawaiian spider radiation. *Journal of Arachnology*, 79-85.

56. Gillespie, R.G. (2016). Island time and the interplay between ecology and evolution in species diversification. *Evolutionary applications* *9*, 53-73.

57. Gillespie, R.G., Croom, H.B., and Hasty, G.L. (1997). Phylogenetic relationships and adaptive shifts among major clades of *Tetragnatha* spiders (Araneae: Tetragnathidae) in Hawai'i.

58. Gillespie, R. (2004). Community assembly through adaptive radiation in Hawaiian spiders. *Science* *303*, 356-359.

59. Blackledge, T.A., Binford, G.J., and Gillespie, R.G. (2003). Resource use within a community of Hawaiian spiders (Araneae: Tetragnathidae). In *Annales Zoologici Fennici*. (JSTOR), pp. 293-303.

60. Blackledge, T.A., and Gillespie, R.G. (2004). Convergent evolution of behavior in an adaptive radiation of Hawaiian web-building spiders. *Proceedings of the National Academy of Sciences of the United States of America* *101*, 16228-16233.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

61. Wilmer, J.W., Hall, L., Barratt, E., and Moritz, C. (1999). Genetic Structure and Male-Mediated Gene Flow in the Ghost Bat (*Macroderma gigas*). *Evolution*, 1582-1591.

62. Kjer, K.M., Zhou, X., Frandsen, P.B., Thomas, J.A., and Blahnik, R.J. (2014). Moving toward species-level phylogeny using ribosomal DNA and COI barcodes: an example from the diverse caddisfly genus *Chimarra* (Trichoptera: Philopotamidae). *Arthropod Systematics & Phylogeny* 72, 345-354.

63. Deiner, K., Renshaw, M.A., Li, Y., Olds, B.P., Lodge, D.M., and Pfrender, M.E. (2017). Long-range PCR allows sequencing of mitochondrial genomes from environmental DNA. *Methods in Ecology and Evolution* 8, 1888-1898.

64. Briscoe, A.G., Goodacre, S., Masta, S.E., Taylor, M.I., Arnedo, M.A., Penney, D., Kenny, J., and Creer, S. (2013). Can long-range PCR be used to amplify genetically divergent mitochondrial genomes for comparative phylogenetics? A case study within spiders (Arthropoda: Araneae). *PLoS one* 8, e62404.

65. Krehenwinkel, H., Kennedy, S. R., Rueda, A., Lam, A., & Gillespie, R. G. (2018). Scaling up DNA barcoding—Primer sets for simple and cost-efficient arthropod systematics by multiplex PCR and Illumina amplicon sequencing. *Methods in Ecology and Evolution*, 9(11), 2181-2193.

66. Krehenwinkel, H., and Pekar, S. (2015). An analysis of factors affecting genotyping success from museum specimens reveals an increase of genetic and morphological variation during a historical range expansion of a European spider. *PLoS one* 10, e0136337.

67. Margam, V.M., Gachomo, E.W., Shukle, J.H., Ariyo, O.O., Seufferheld, M.J., and Kotchoni, S.O. (2010). A simplified arthropod genomic-DNA extraction protocol for polymerase chain reaction (PCR)-based specimen identification through barcoding. *Molecular biology reports* 37, 3631-3635.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

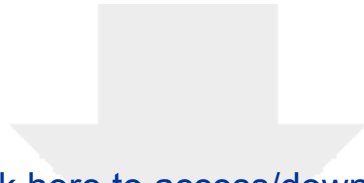
68. Sipos, R., Székely, A.J., Palatinszky, M., Révész, S., Márialigeti, K., and Nikolausz, M. (2007). Effect of primer mismatch, annealing temperature and PCR cycle number on 16S rRNA gene-targeting bacterial community analysis. *FEMS Microbiology Ecology* 60, 341-350.

69. Suzuki, M.T., and Giovannoni, S.J. (1996). Bias caused by template annealing in the amplification of mixtures of 16S rRNA genes by PCR. *Applied and environmental microbiology* 62, 625-630.

70. Krehenwinkel, H., Fong, M., Kennedy, S., Huang, E.G., Noriyuki, S., Cayetano, L., and Gillespie, R. (2018). The effect of DNA degradation bias in passive sampling devices on metabarcoding studies of arthropod communities and their associated microbiota. *PLoS one* 13, e0189188.

71. Wattier, R., Engel, C., Saumitou- Laprade, P., and Valero, M. (1998). Short allele dominance as a source of heterozygote deficiency at microsatellite loci: experimental evidence at the dinucleotide locus Gv1CT in *Gracilaria gracilis* (Rhodophyta). *Molecular Ecology* 7, 1569-1573.

72. Krehenwinke I H; Pomerantz A; Henderson JB; Kennedy SR; Lim JY; Swamy V; Shoobridge JD; Graham N; Patel NH; Gillespie RG; Prost S (2019): Supporting data for: "Nanopore sequencing of long ribosomal DNA amplicons enables portable and simple biodiversity assessments with high phylogenetic resolution across broad taxonomic scale".. GigaScience Database. <http://dx.doi.org/10.5524/100552>



Click here to access/download

Supplementary Material

SupplementaryTable1_SampleList.xlsx

