

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-18-00245
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:	
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. Long rDNA amplicons scale up DNA barcoding by enabling the accurate recovery of taxonomic and phylogenetic diversity. However, bulk community analyses using long-read approaches may introduce biases and will require further exploration.</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
	Nipam H. Patel
	Rosemary G. Gillespie
	Stefan Prost
Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>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.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>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.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
Availability of data and materials	Yes

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” section of your manuscript.

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



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, Nipam H. Patel^{2,7}, Rosemary G.
7 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) and Stefan Prost
23 (stefan.prost@berkeley.edu)

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. Long rDNA amplicons scale up DNA barcoding by enabling the accurate recovery of taxonomic and phylogenetic diversity. However, bulk community analyses using long-read approaches may introduce biases and will require further exploration.

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

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. This
81 problem is exacerbated by the fast evolutionary rate of mitochondrial DNA, which leads to a
82 quick saturation with mutations, increasing the probability of homoplasy between divergent
83 lineages. The accurate estimation of phylogenetic diversity across wide taxonomic scales,
84 however, is an important component of biodiversity research [9]. Moreover, mitochondrial DNA
85 is not always the best marker to reflect species differentiation, as different factors are known to
86 inflate mitochondrial differentiation in relation to the nuclear genomic background. For example,
87 male biased gene flow [10] or infections with reproductive parasites [11] (e.g. *Wolbachia*) can
88 lead to highly divergent mitochondrial lineages in the absence of nuclear differentiation. In
89 contrast, introgressive hybridization can cause the complete replacement of mitochondrial
90 genomes (see e.g. [12, 13]), resulting in shared mitochondrial 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 evolutionary conservation, which are interrupted by
103 highly variable sequence stretches [16]. While some rDNA gene regions are entirely conserved

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

112
113 Spanning about 8 kb, the ribosomal cluster is fairly large, and current barcoding protocols, e.g.
114 using Sanger sequencing or Illumina amplicon sequencing, can only target short sequence
115 stretches of 150 – 1,000 bp. Such short stretches of 28S and 18S are often too conserved to
116 identify young species pairs [19]. The ITS regions, on the other hand, are too variable to design
117 truly universal primers, leading to a considerable amount of taxon dropout during PCR.
118 Moreover, ITS sequences can show considerable length variation between taxa, and are often
119 too long for short amplicon-based barcoding [20]. Consequently, it would be ideal to amplify and
120 sequence a large part of the ribosomal cluster in one fragment. A solution to sequence the
121 resulting long amplicons is offered by recent developments in third generation sequencing
122 platforms, which now enable researchers to generate ultra-long reads, of up to 800 kb [21].
123 Recently, amplicons of several kilobases of the rDNA cluster were sequenced using Pacific
124 Bioscience (PacBio) technology, to explore fungal community composition [22, 23]. But while
125 PacBio sequencing is well suited for long amplicon sequencing, it is currently not readily
126 available to every laboratory due to the high cost and limited distribution of sequencing
127 machines.

128

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

129 A cost-efficient and readily available alternative is provided by nanopore sequencing
130 technology. The MinION sequencer (Oxford Nanopore Technologies) is small in size,
131 lightweight, allows for sequencing of several Gb's of DNA with average read lengths over 10 kb
132 on a single flow cell [24] and is available starting at \$1,000. Despite a raw read error rate of
133 about 12-22 % [21], highly accurate consensus sequences can be called from nanopore data
134 [25, 26]. The MinION is well suited for amplicon sequencing, and a simple dual indexing
135 strategy can be used to demultiplex amplicon samples [27]. This technology offers tremendous
136 potential for long-read barcoding applications, as recently shown in an analysis in fungi [26].
137 However, current analyses are still exploratory or limited in taxonomic focus and streamlined
138 analysis pipelines to establish the method across the eukaryote tree of life are still missing.

139
140 Considering this background, we explore the feasibility of nanopore sequencing of long rDNA
141 amplicons as a simple, cost efficient and universal eukaryote DNA barcoding approach. We
142 compile a workflow from PCR amplification, to library preparation, to demultiplexing and
143 consensus calling (see Fig. 1 for an overview). We explore the error profile of nanopore
144 consensus sequences and introduce MiniBar, a new software to demultiplex dual indexed
145 nanopore amplicon sequences. We test the utility of the ribosomal cluster for molecular
146 taxonomy and phylogenetics across divergent plant and animal taxa. A particular focus of our
147 analysis are arthropods, the most diverse group in the animal kingdom [28], which are highly
148 threatened by current mass extinctions [29]. Using a dataset of spiders, we compare the
149 taxonomic resolution of the ribosomal cluster with that offered by molecular barcoding using
150 mitochondrial COI, the currently preferred barcode marker for arthropods. Oxford Nanopore
151 Technologies' MinION is a portable sequencer, and Nanopore based DNA barcoding has been
152 applied in remote sites outside of conventional labs (see eg. [25, 30, 31]). Such field-based
153 applications confront researchers with additional complexities and challenges. To highlight the

1
2
3
4 154 simplicity of our approach, we tested it under field conditions and generated long rDNA barcode
5
6 155 sequences using a miniaturized mobile laboratory in a Peruvian rainforest.
7
8
9 156

10 157 We also tested the efficacy of long-read rDNA sequencing for metabarcoding of bulk community
11
12 158 samples. A study of bacterial communities [32] suggests Nanopore long-read sequencing as a
13
14 159 powerful tool for community characterization, but also found pronounced biases in the
15
16 160 recovered taxon abundance. Currently, little is known about the utility of long-read sequencing
17
18 161 for animal community analysis. Metabarcoding protocols for community samples need to be
19
20 162 carefully optimized, as they can suffer from pronounced taxonomic biases, e.g. due to primer
21
22 163 binding or polymerase efficiency [33]. Well established Illumina based short read metabarcoding
23
24 164 protocols can account for these biases and allow for a very high qualitative and even
25
26 165 quantitative recovery of taxa in communities [34]. However, additional, yet unexplored, biases
27
28 166 may affect long-read metabarcoding. We thus also test the utility of long-read rDNA barcoding
29
30 167 to recover taxonomic diversity from arthropod mock communities. We compare the qualitative
31
32 168 (species richness) and quantitative (species abundance) recovery of taxa by long-read
33
34 169 sequencing with that based on short read Illumina amplicon sequencing of the 18SrDNA.
35
36
37
38
39
40 170

41
42 171 Overall, we demonstrate that long rDNA amplification and sequencing on the MinION platform is
43
44 172 a straightforward, cost effective, and universal approach for eukaryote DNA barcoding. It
45
46 173 combines robust taxonomic assignment power with high phylogenetic resolution and will enable
47
48 174 future analyses of taxonomic and phylogenetic diversity across wide taxonomic scales.
49
50

51 175
52
53 176
54

55 177 **Data Description and Analyses**
56

57 178
58
59 179 ***DNA extraction, PCR and library preparation***
60
61
62
63
64
65

1
2
3
4 180 We analyzed 114 specimens of eukaryotes including 17 insect and 42 spider species, two
5
6 181 annelid and nine plant species (Supplementary Table 1). Some feeder insects and the annelids
7
8 182 were purchased at a pet store. The remaining specimens were collected in oak forest on the
9
10 183 University of California Berkeley's campus or in native rainforests of the Hawaiian Archipelago
11
12 184 (under the Hawaii DLNR permit: FHM14-349). We particularly focused our arthropod sampling
13
14 185 on spiders, which are ubiquitous and essential predators in all terrestrial ecosystems. Recent
15
16 186 phylogenomic work [35] provided us with a solid baseline to test the efficiency of rDNA
17
18 187 amplicons for phylogenetic and taxonomic purposes. We included a taxonomically diverse
19
20 188 collection of 16 spider families from the Araneoidea, the RTA clade and a haplogyne outgroup
21
22 189 species. Within spiders, we additionally focused on the genus *Tetragnatha*, which has
23
24 190 undergone a striking adaptive radiation on Hawaii.
25
26 191
27
28
29
30
31 192 DNA was extracted from each sample using the Qiagen Archivepure kit (Qiagen, Valencia, CA,
32
33 193 USA) according to the manufacturer's protocols. The DNA integrity was checked on an agarose
34
35 194 gel. Only samples with high DNA integrity were used for the following PCRs. All DNA extracts
36
37 195 were quantified using a Qubit fluorometer using the high sensitivity dsDNA assay (Thermo
38
39 196 Fisher, Waltham, MA, USA) and diluted to concentrations of 20 ng/μl. We designed a primer
40
41 197 pair of each 27 bases to amplify a ~4,000 bp fragment of the ribosomal DNA, including partial
42
43 198 18S and 28S as well as full ITS1, 5.8S and ITS2 sequences (18S_F4
44
45 199 GGCTACCACATCYAARGAAGGCAGCAG and 28S_R8
46
47 200 TCGGCAGGTGAGTYGTTRCACAYTCCT). The primers were designed using alignments of
48
49 201 partial 18S and 28S sequences of ~1,000 species of eukaryotes, with a focus on animals. The
50
51 202 primers targeted highly conserved regions across all analyzed taxa. Degenerate sites were
52
53 203 incorporated to account for variation. We aimed for high annealing temperatures (65-70°C) to
54
55 204 impose stringent amplification. These were calculated using the NEB Tm Calculator
56
57
58
59 205 (<https://tmcalsculator.neb.com/#!/main>).
60
61
62
63
64
65

1
2
3
4 206
5
6 207 To index every PCR amplicon separately, we used a dual indexing strategy with each primer
7
8 208 carrying a unique 15 bp index sequence at its 5'-tail. Index sequences were designed using
9
10 209 Barcode Generator (http://comailab.genomecenter.ucdavis.edu/index.php/Barcode_generator)
11
12 210 with a minimum distance of 10 bases between each index. A total of 15 forward and 16 reverse
13
14 211 indexes were designed. Every sample was amplified separately using the Q5 Hot Start High-
15
16 212 Fidelity 2X Master Mix (NEB, Ipswich, MA, USA) in 15 µl reactions, at 68°C annealing
17
18 213 temperature, with 35 PCR cycles and using 50 ng of template DNA per PCR. All PCR products
19
20 214 were checked and quantified on an agarose gel and then pooled. The final pool was cleaned
21
22 215 from residual primers by 0.75 X AMPure Beads XP (Beckman Coulter, Brea, CA, USA). DNA
23
24 216 library preparation was carried out according to the 1D PCR barcoding amplicons SQK- LSK108
25
26 217 protocol (Oxford Nanopore Technologies, Oxford, UK). Barcoded DNA products were pooled
27
28 218 with 5 µl of DNA CS (a positive control provided by ONT) and an end-repair was performed
29
30 219 (NEB-Next Ultra II End-prep reaction buffer and enzyme mix), then purified using AMPure XP
31
32 220 beads. Adapter ligation and tethering was carried out with 20 µl Adapter Mix and 50 µl of NEB
33
34 221 Blunt/TA ligation Master Mix. The adapter-ligated DNA library was then purified with AMPure
35
36 222 beads XP, followed by the addition of Adapter Bead binding buffer, and finally eluted in 15 µl of
37
38 223 Elution Buffer. Each R9 flow cell was primed with 1000 µl of a mixture of Fuel Mix and nuclease-
39
40 224 free water. Twelve µl of the amplicon library were diluted in 75 µL of running buffer with 35 µL
41
42 225 RBF, 25.5 uL LLB, and 2.5 µL nuclease-free water and then added to the flow cell via the
43
44 226 SpotON sample port. The "NC_48Hr_sequencing_FLO-MIN107_SQK-
45
46 227 LSK108_plus_Basecaller.py" protocol was initiated using the MinION control software,
47
48 228 MinKNOW.
49
50
51
52
53
54
55
56
57

58 230 ***Field trial in the Amazon rainforest***
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

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

245 **Bioinformatics**

247 ***Raw data processing and consensus calling***

248 The fastq files generated by the ONT software MinKNOW were de-multiplexed using MiniBar
249 (see description below), with index edit distances of 2, 3, and 4 and a primer edit distance of 11.
250 Next, the reads were filtered for quality (>13) and size (>3kb) using Nanofilt [36](
251 <https://github.com/wdecoster/nanofilt>). Individual consensus sequences were created using
252 Allele Wrangler (<https://github.com/transplantation-immunology/allele-wrangler/>) for
253 demultiplexed fastq files with a minimum coverage of 30. Error correction was performed using
254 RACON [37] (<https://github.com/isovic/racon>). To do so, we first mapped all the reads back to
255 the consensus using minimap (<https://github.com/lh3/minimap2>). We performed two cycles of

1
2
3
4 256 running minimap and RACON. Final consensus sequences were compared against the NCBI
5
6 257 database using BLASTn to check if the taxonomic assignment was correct.
7
8
9 258

10
11 259 We performed multiple tests to validate and optimize the consensus accuracy of long-read
12
13 260 barcode sequences. To comparatively assess the accuracy, we used consensus sequences of
14
15 261 short 18S and 28SrDNA amplicons, which were previously generated using Illumina amplicon
16
17 262 sequencing for the 47 analyzed Hawaiian *Tetraglatha* specimens (Kennedy unpublished data).
18
19
20 263 These sequences were aligned with the respective stretches of our nanopore consensus
21
22 264 sequences using ClustalW in MEGA [38]. All alignments were then visually inspected and edited
23
24 265 manually, where necessary. Pairwise distances between Illumina and nanopore consensus
25
26 266 were calculated in MEGA.
27
28
29 267

30
31 268 To measure consensus accuracy over the whole ribosomal amplicon, we utilized genome
32
33 269 skimming data [39] for six Hawaiian *Peperomia* plant species (Lim et al unpublished data). 150
34
35 270 bp paired-end TruSeq gDNA shotgun libraries for the six *Peperomia* samples were sequenced
36
37
38 271 on a single HiSeq v4000 lane (Illumina, San Diego, CA, USA). The resulting paired-end reads
39
40 272 were trimmed and filtered using Trimmomatic v0.36 [40] and mapped to their respective
41
42 273 nanopore consensus sequences using bowtie2 [41] under default parameter values and
43
44 274 allowing for minimum and maximum fragment size of 200 and 700 bases respectively. Mapping
45
46
47 275 coverage of Illumina reads to nanopore consensus sequences ranged between 150 - 600 X with
48
49 276 a mean of ~ 300 X across all six samples. We called Illumina read based consensus sequences
50
51 277 for each *Peperomia* species using bcftools [42], and aligned them with the previously generated
52
53 278 nanopore consensus sequences. Pairwise genetic distances were then calculated in MEGA as
54
55
56 279 described above. We performed two independent distance calculations: 1) excluding indels, i.e.
57
58 280 only using nucleotide substitutions to estimate genetic distances, and 2) including indels as
59
60 281 additional characters.
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

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

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

A recent study [25] showed that accurate consensus sequences from nanopore data can be generated using only 30x coverage. We tested 18 different assembly coverages from 10 to 800 sequences for a *Peperomia* species, to explore optimal assembly coverage. We randomly subsampled the quality filtered and demultiplexed fastq file for the according specimen each 10 times for each tested assembly coverage. Consensus sequences were then assembled and genetic distances to the Illumina consensus calculated as described above.

Phylogenetic and taxonomic analysis

We carried out phylogenetic analyses on two hierarchical levels. First, we built a phylogeny for all higher eukaryote taxa in our dataset, which included plants, animals and fungi. Second, we took a closer look into the phylogeny of spiders. The resulting quality checked consensus sequences of all taxa were aligned using ClustalW in MEGA. The alignments were visually inspected and manually edited. Due to the deep divergence in the eukaryote data set, the highly variable ITS sequences could not be aligned and were excluded. For the analyses of spiders,

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

308 we retained both ITS sequences and aligned the whole rDNA amplicon. Appropriate models of
309 sequence evolution for each gene fragment of the rDNA cluster were identified using
310 PartitionFinder [43]. Phylogenies were built using MrBayes [44], with 4 heated chains, a chain
311 length of 1,100,000, subsampling every 200 generations and a burnin length of 100,000.

312
313 Focusing on the endemic Hawaiian *Tetragnatha* species, we also tested the utility of the
314 ribosomal cluster for taxonomic identification, as we also had COI barcodes available for these
315 species. Our dataset contained ribosomal DNA sequences for 47 specimens in 16 species. We
316 calculated pairwise genetic distances between and within all species for the whole ribosomal
317 cluster and for each separate gene region of the rDNA cluster using MEGA. As the 18S and
318 5.8S did not yield any species level resolution within Hawaiian *Tetragnatha*, they were not
319 analyzed separately. To compare the taxonomic resolution of the ribosomal cluster with that of
320 the commonly used mitochondrial COI, we calculated inter- and intraspecific distances for an
321 alignment of 418 bp of the COI barcode region for the same spider specimens (Kennedy et al.
322 unpublished data). We performed a Mantel test using the R package ade4 [45] to test for a
323 significant correlation between COI and ribosomal DNA based distances. A comparison of
324 intraspecific and interspecific distances for mitochondrial COI and ribosomal DNA also allowed
325 us to test for the presence of a barcode gap.

326

327 ***Nanopore based arthropod metabarcoding***

328 To test for the possibility of estimating arthropod community composition from Nanopore
329 sequencing, we prepared four mock communities of different amounts of DNA extracts from 9
330 species of arthropods from different orders (see Supplementary Table 2). The samples were
331 amplified using the Q5 High Fidelity Mastermix as described above at 68 °C annealing
332 temperature and 35 PCR cycles. We additionally tested two variations of PCR conditions. We
333 either reduced the annealing temperature to 63 °C or reduced the PCR cycle number to 25.

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

334 In order to compare our results with those from an optimized Illumina short read protocol, we
335 amplified all samples for a ~300 bp fragment of the 18SrDNA using the primer pair
336 18S2F/18S4R [46]. Amplification and library preparation were performed as described in [47]
337 using the Qiagen Multiplex PCR kit. The 18S amplicon pools were sequenced on an Illumina
338 MiSeq using V3 chemistry and 2 x 300 bp reads. Sequence quality filtering, read merging and
339 primer trimming were performed as described in [34].

340
341 A library of 18S sequences for all included arthropod species (from [34]) was used as a
342 reference database to identify the recovered sequences using BLASTn [48], with a minimum e-
343 value of 10^{-4} and a minimum overlap of 95 %. Despite the high raw error rate of nanopore reads,
344 taxonomic status of sequences could be assigned using BLAST, as our pools contained
345 members of highly divergent orders. We compared the qualitative (number of species) and
346 quantitative (abundance of species) recovery of taxa from the communities by nanopore long-
347 read and Illumina short read data. To estimate the recovery of taxon abundances, we calculated
348 a fold change between input DNA amount and recovered reads for each taxon and mock
349 community. A fold change of zero corresponded to a 1:1 association of taxon abundance and
350 read count, while positive or negative values indicated higher or lower read counts than the
351 taxon's actual abundance.

352
353 MiniBar
354 We created a de-multiplexing software, called MiniBar. It allows customization of search
355 parameters to account for the high read error rates and has built-in awareness of the dual
356 barcode and primer pairs flanking the sequences. MiniBar takes as input a tab-delimited
357 barcode file and a sequence file in either fasta or fastq format. The barcode file contains, at a
358 minimum, sample name, forward barcode, forward primer, reverse barcode, and reverse primer
359 for each of the samples potentially in the sequence file. The software searches for barcodes 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

360 for a primer, each permitting a user defined number of errors, an error being a mismatch or
361 indel. Error count to determine a match can either be a percentage of each of their lengths or
362 can be separately specified for barcode and primer as a maximum edit distance [49]. Output
363 options permit saving each sample in its own file or all samples in a single file, with the sample
364 names in the fasta or fastq headers. The found barcode primer pairs can be trimmed from the
365 sequence or can remain in the sequence distinguished by case or color. MiniBar, written in
366 Python 2.7, can also run in Python 3 and has the single dependency of the Edlib library module
367 for edit distance measured approximate search [50]. MiniBar can be found at
368 <https://github.com/calacademy-research/minibar> along with test data.

369

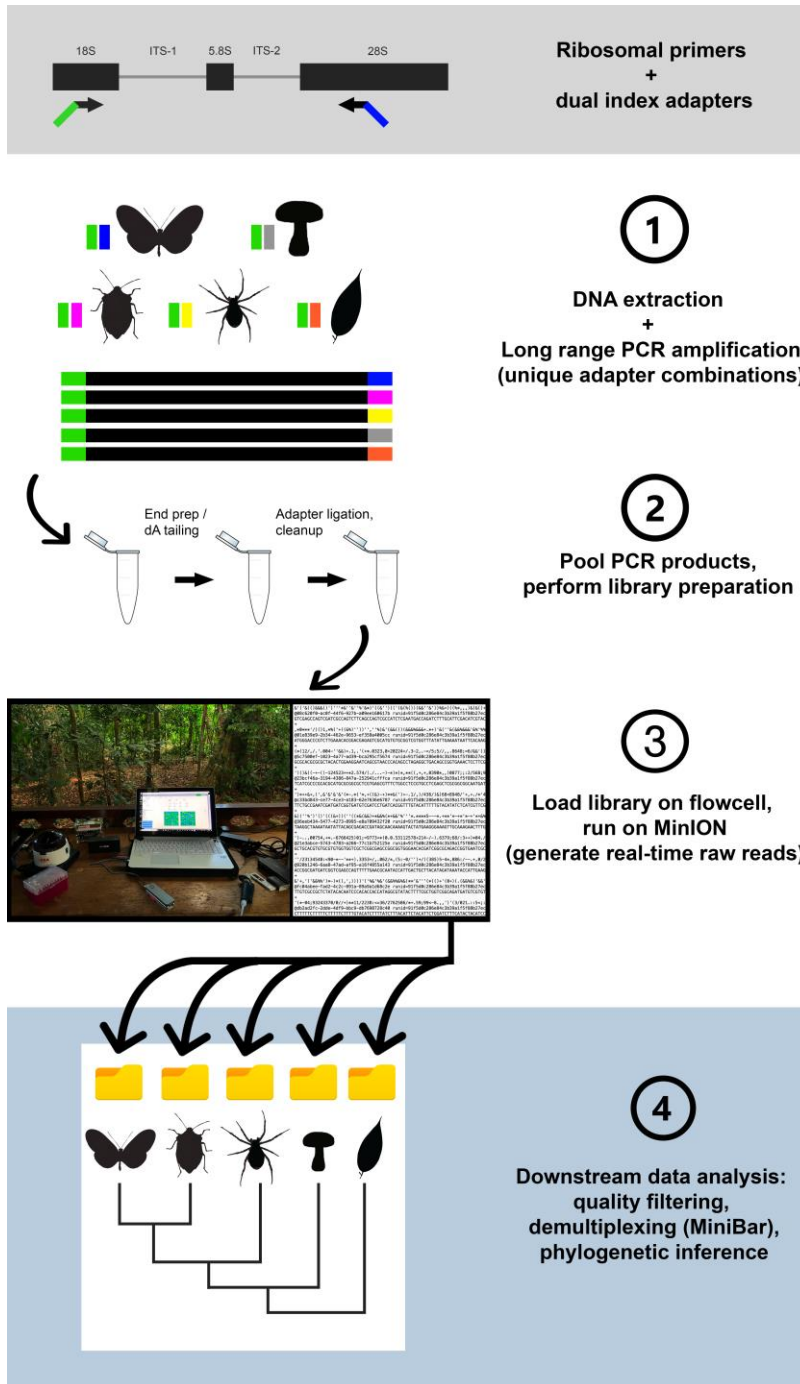


Figure 1. Workflow for the design, amplification, and sequencing of the ribosomal DNA cluster.

Results

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

376

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

378 After quality filtering and trimming, our nanopore run yielded 245,433 reads. We tested edit
379 distances of two, three and four bases in MiniBar to demultiplex samples. Increasing edit
380 distances led to a significant increase in read numbers assigned to index combinations
381 (Pairwise Wilcoxon Test, FDR-corrected P -value < 0.05). On average, we found 355 reads per
382 specimen for an edit distance of two, 647 for a distance of three and 1,051 for a distance of four.
383 However, at an edit distance of four, we found a considerable increase of wrongly assigned
384 samples. Using Illumina shotgun sequencing-derived consensus sequences of rDNA from six
385 *Peperomia* plants, we tested the accuracy of the nanopore consensus assemblies based on the
386 three edit distances (Fig. 2). While a distance of four yielded the highest number of assigned
387 reads (1,785 on average), it also led to slightly more inaccurate consensus assemblies, with an
388 average distance of 2.072 % to Illumina based consensus sequences. We found a significant
389 increase of consensus accuracy (Pairwise Wilcoxon Test, FDR corrected $P < 0.05$) for edit
390 distances of two (0.165 % average distance) and three (0.187 % average distance). Despite
391 significant differences in assigned reads (1,091 vs. 637 reads on average), there was not a
392 significant difference in consensus accuracy of edit distances of two versus three bases
393 (Pairwise Wilcoxon Test, FDR corrected $P > 0.05$).

394

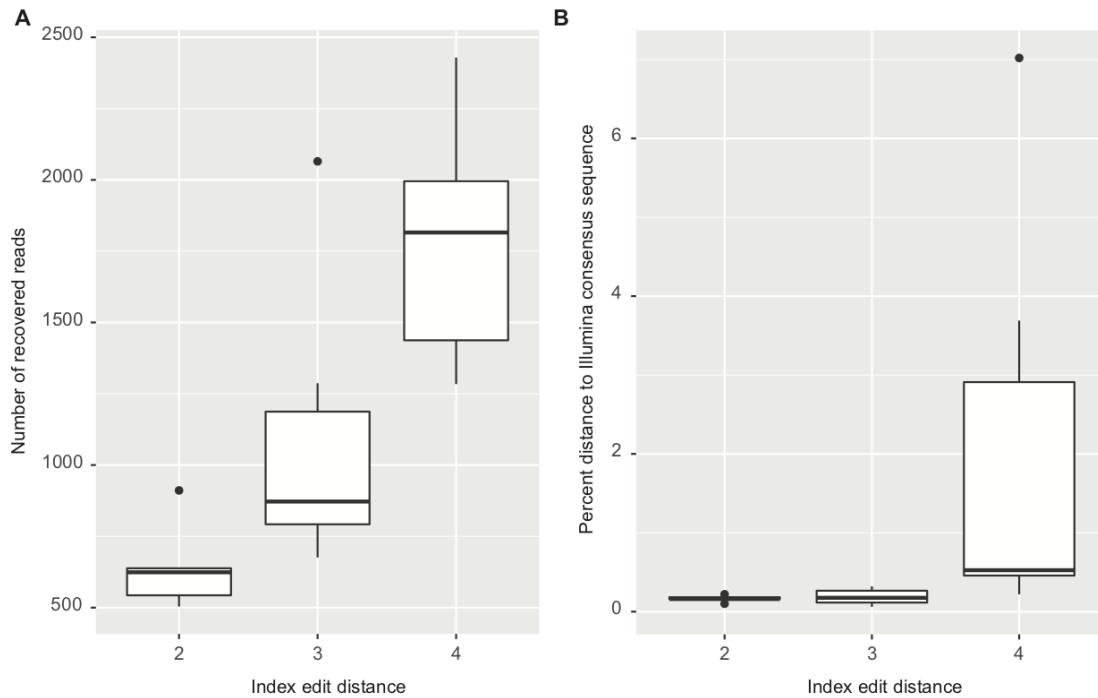


Figure 2: Comparison of recovered sequences and consensus accuracy for different index edit distances in Minibar. A) Number of recovered reads for six *Peperomia* species at index edit distances of two, three and four. B) Pairwise sequence divergence between Illumina and Nanopore based consensus sequences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

1
2
3
4 409 larva of a butterfly and a feeder mealworm (*Zophobas morio*) generated consensus sequences
5
6 410 for plants.
7
8
9 411 A comparison of our consensus sequences for 47 Hawaiian specimens of the spider genus
10
11 412 *Tetragnatha* with short Illumina amplicon sequencing-derived 18S and 28S rDNA sequences
12
13 413 suggests a very high consensus accuracy. Except for a single specimen, with a single
14
15 414 substitution error, all nanopore based consensus sequences were completely identical to the
16
17 415 Illumina based consensus. However, the corresponding 18S and 28S fragments did not contain
18
19
20 416 long stretches of homopolymer sequences, where nanopore raw read errors are known to
21
22 417 accumulate [51]. Despite containing several homopolymers, the nanopore derived *Peperomia*
23
24 418 consensus sequences were highly accurate (Supplementary Fig. 1). Including gaps in the
25
26 419 alignment, an average distance of 0.165 % to Illumina based consensus sequences was found.
27
28
29 420 Errors were clustered in Indel regions. After excluding gaps, the average distance dropped to
30
31 421 0.102 %.
32
33 422
34
35 423 We found only a small effect of sequence coverage on consensus assembly accuracy
36
37 424 (Supplementary Fig. 2). Even at 10-fold coverage, a low average distance of 0.257% to Illumina
38
39 425 consensus sequences was observed. However, at 20-fold coverage, the average distance
40
41 426 significantly decreased to 0.128 % (Pairwise Wilcoxon Test, FDR corrected $P < 0.05$). A slight,
42
43 427 but not significant, decrease of distance was observed with increasing coverage, with optimal
44
45 428 consensus accuracy at 300-fold coverage (0.031 % distance). At coverages larger than 300, the
46
47 429 consensus accuracy slightly decreased (average distance of 0.103 % at 800 X coverage).
48
49
50
51 430
52
53 431 The length of the rDNA amplicon was quite variable between taxa. Compared to animals, plant
54
55 432 specimens showed a significantly shorter amplicon (Pairwise Wilcoxon Test, FDR corrected $P <$
56
57 433 0.05). The length difference was found for the actual gene sequences (18S, 5.8S, 28S: 3,063 vs
58
59 434 2,781 bp on average; Supplementary Fig. 3A) as well as including the ITS sequences (3,741 vs.
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

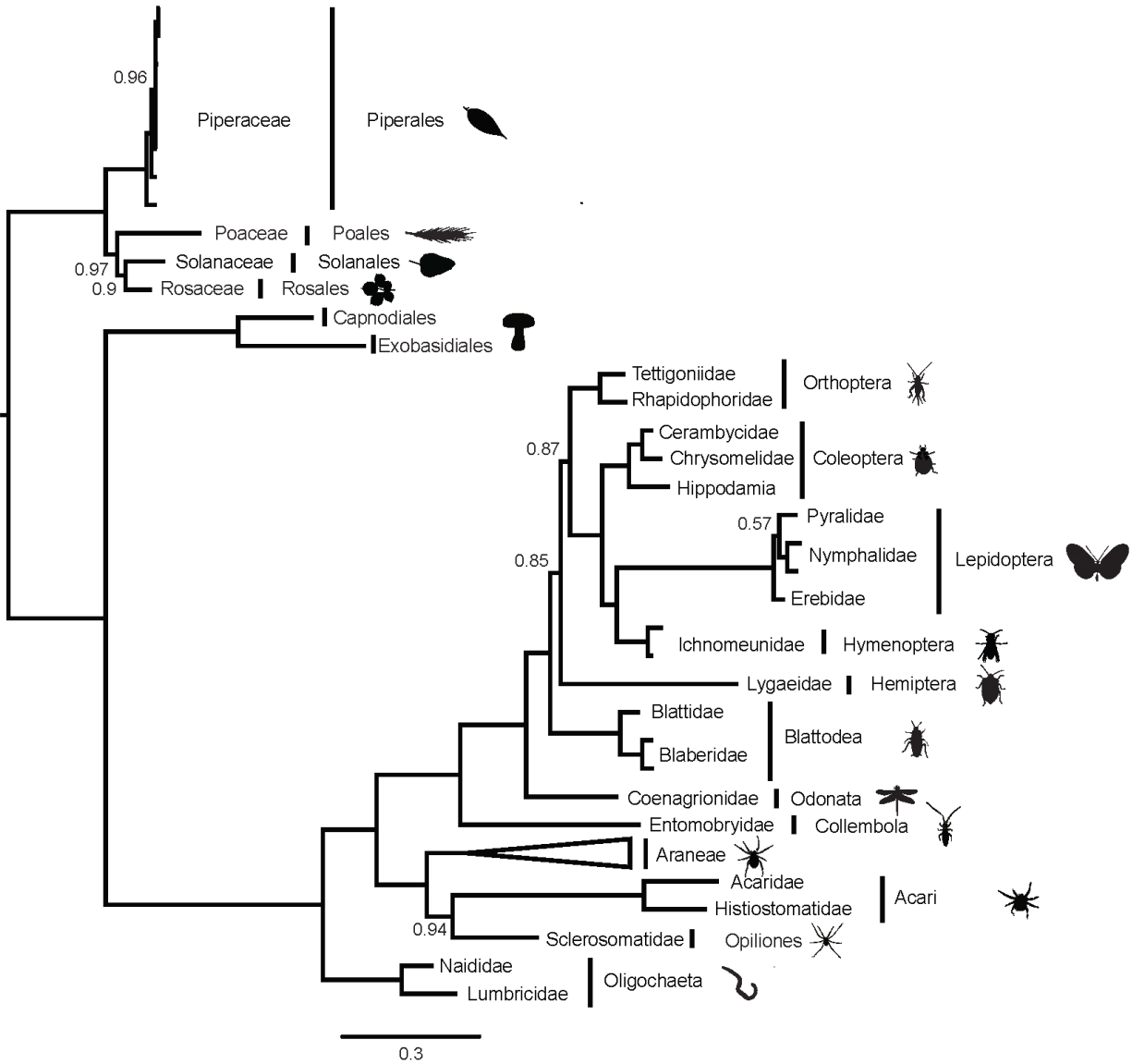
435 3,241 bp on average, Supplementary Fig. 3B). Within arthropods, we found significant length
436 differences between arachnid and insect sequences. On average, insects carried a significantly
437 longer rDNA sequence than arachnids (Supplementary Fig. 4; Pairwise Wilcoxon Test, FDR
438 corrected $P < 0.05$). This holds true for the gene sequences (3,154 vs. 3,047 bp for 18S, 5.8S,
439 28S on average), as well as the whole amplicon, including ITS sequences (4,192 vs. 3,644 bp
440 on average). While most spiders showed very stable length distributions for the rDNA amplicon
441 length (average length \pm standard deviation across all Araneae: 3,629 bp \pm 81), several insect
442 orders had rDNA sequences of more variable length (Coleoptera: 4,488 bp \pm 352; Lepidoptera:
443 4363 bp \pm 603).

444
445 In contrast to the variable length of the rDNA cluster, we found a very stable GC content across
446 the whole taxonomic spectrum (46.75 \pm 2.67 % across all taxa). GC content of plants and
447 animals was highly similar (Supplementary Fig 3c) (plants: 46.01 \pm 1.66 %; animals: 46.82 \pm
448 2.74 %). Highly similar GC content was also found between insects (46.67 \pm 3.73 %) and
449 arachnids (46.93 \pm 2.47 %) (Supplementary Fig 4c).

450
451

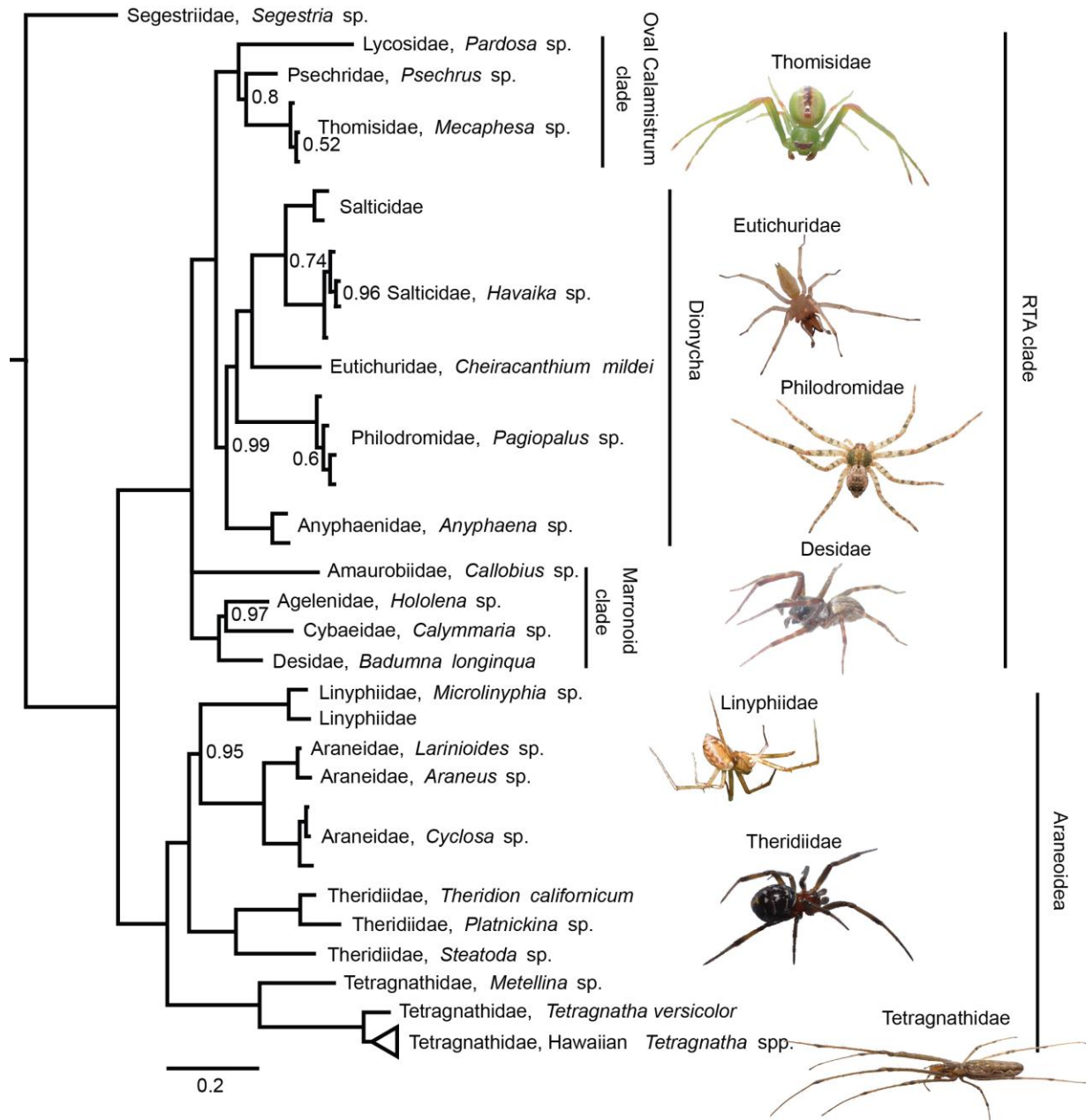
1
2
3
4 452 **Phylogenetic reconstruction**

5
6 453
7
8
9



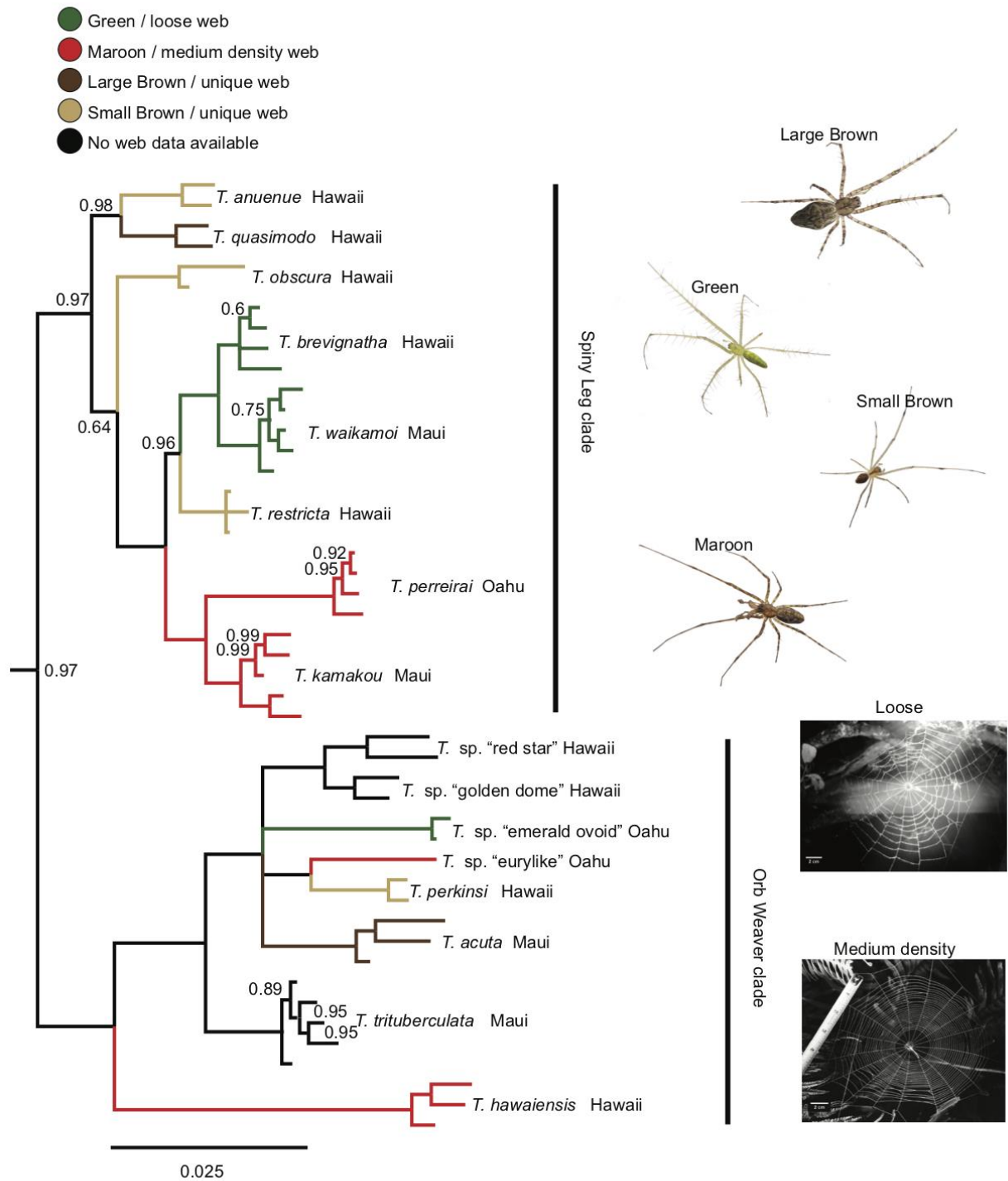
455 **Figure 3 Bayesian consensus phylogeny based on a 3,656 bp alignment of 18S, 5.8S and**
456 **28S sequences of 117 animal, fungal and plant taxa.** The phylogeny is rooted using plants
457 as outgroup. Branches are annotated with family and order level taxonomy. The Araneae clade
458 of 83 specimens is collapsed. 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



460
461 **Figure 4. Bayesian consensus phylogeny of 83 spiders in 16 families, based on a 4,214**
462 **bp alignment of 18S, ITS1, 5.8S, ITS2 and 28S.** The phylogeny is rooted using the basal
463 haplogyne *Segestria* sp. The clade containing Hawaiian members of the genus *Tetragnatha* is
464 collapsed (the uncollapsed clade is shown in Fig. 5). Only posterior probability values below 1
465 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



466
467 **Figure 5. Section of the same phylogeny as Fig. 4, with expansion of the clade of 16**
468 **Hawaiian *Tetragnatha* species.** Different “Spiny Leg” ecomorphs and web architectures are
469 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

470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495

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.

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, 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

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

work: Lycosidae was basal to the clade formed by Psecridae 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.

500

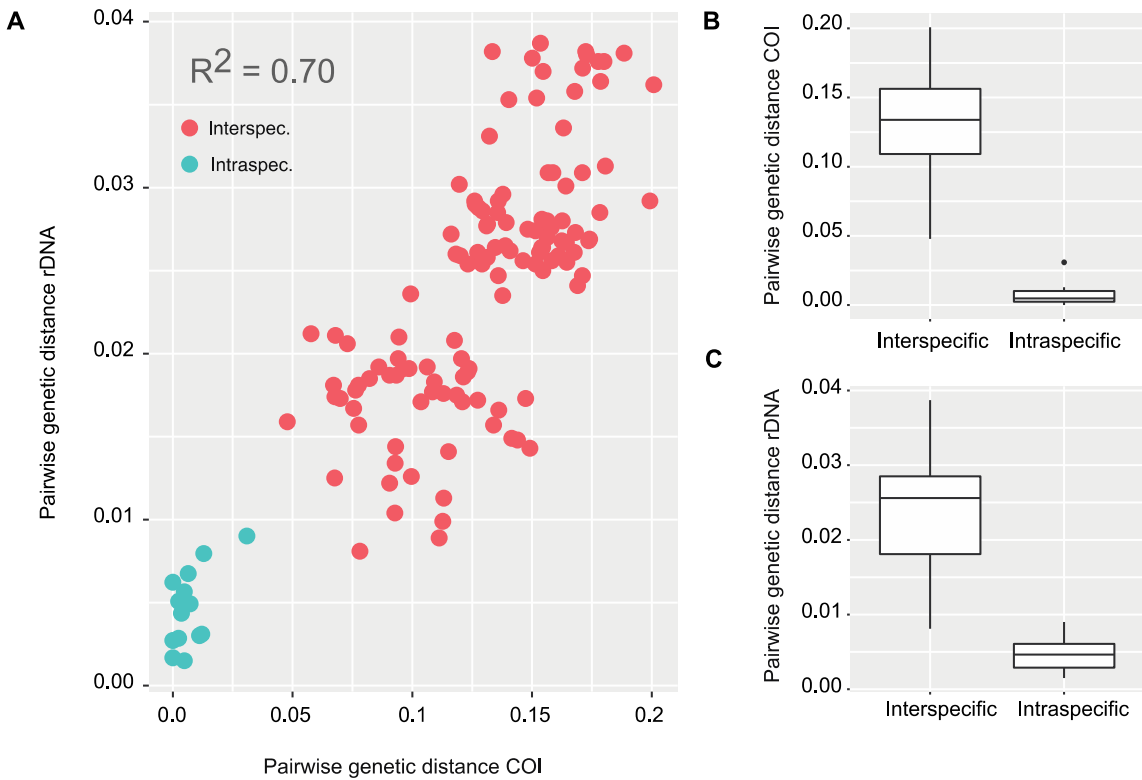
We recovered Hawaiian *Tetragnatha* as a well supported monophyletic clade within the Tetragnathidae. We found two main clades of Hawaiian *Tetragnatha* (Fig. 5), both of which have been supported by earlier work [54-57]: the orb weaving clade and the “Spiny Leg clade” of actively hunting species. All *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 ecotypes, each of which is associated with a particular substrate type: “large brown” (*T. quasimodo*) with tree bark, “small brown” (*T. anuenue*, *T. obscura* and *T. restricta*) with twigs, “green” (*T. brevignatha* and *T. waikamoi*) with green leaves, and “maroon” (*T. perreirai* and *T. kamakou*) with lichen. While green and maroon ecotypes 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 [35, 58]. 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 [58], and are thus classified as “unique”.

520

521

1
2
3
4 **522 Taxonomic resolution**

5
6 523 Our inferred genetic distances for rDNA sequences within and between Hawaiian *Tetragnatha*
7
8 524 species were significantly correlated to those found for COI sequences of the same taxa ($R^2 =$
9
10 525 0.70, $P < 0.001$) (Fig. 6a). A Mantel test also suggested highly significant correlation of
11
12 526 mitochondrial COI and nuclear rDNA based distances (Mantel test, 9999 replicates, $P < 0.001$).
13
14 527 Hence, the rDNA cluster supported a very similar pattern of genetic differentiation to COI.
15
16 528 However, the faster evolutionary rate of COI was reflected in lower distances for the whole
17
18 529 rDNA than for COI. Interspecific distances were significantly higher than intraspecific ones for
19
20 530 COI and rDNA (Fig 6b,c). No overlap of intra and interspecific distances was evident for COI,
21
22 531 suggesting the presence of a barcode gap. A small overlap of intra and interspecific distances
23
24 532 was evident for the rDNA (Supplementary Table 3). Like the combined rDNA cluster, genetic
25
26 533 distances for different parts of the rDNA cluster all showed significant correlation with COI
27
28 534 based distances, when analyzed separately (R^2 28S = 0.57, R^2 ITS1 = 0.68, R^2 ITS2 = 0.56, $P <$
29
30 535 0.001) (Supplementary Fig. 5). While the 28SrDNA showed considerably lower distances than
31
32 536 COI, those for ITS1 and ITS2 were more comparable to COI (Supplementary Fig. 5b-d). Yet,
33
34 537 interspecific and intraspecific distances for COI were significantly different from those for any
35
36 538 part of the rDNA cluster (Pairwise Wilcoxon Test, FDR corrected $P < 0.05$).
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



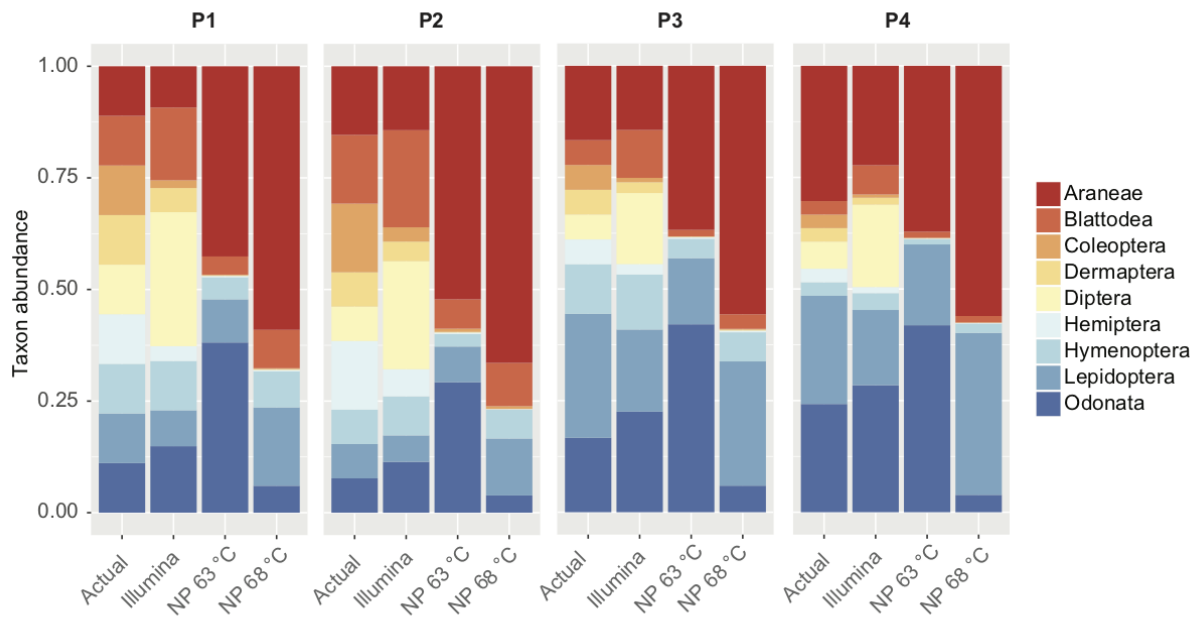
540
 541 **Figure 6 Inter and intraspecific genetic distances for the nuclear rDNA and mitochondrial**
 542 **COI for Hawaiian *Tetragnatha* spiders.** A) Correlation of pairwise genetic distance between
 543 (red) and within (green) 16 Hawaiian *Tetragnatha* species based on COI and the full rDNA
 544 amplicon. B) Interspecific and intraspecific genetic distances for the same spider species based
 545 on mitochondrial COI and C) the whole rDNA amplicon.

546
 547
 548
 549 ***Field trial in the Amazon rainforest***

550 On March 26, 2018 we set out to test this method and a portable laboratory (as described in
 551 Pomerantz, et al. [25]) during an expedition to the Peruvian Amazon at the Refugio Amazonas
 552 Lodge (Supplementary Fig. 6). This field site is a “Terra firme” forest in the sector of
 553 “Condenado”, approximately two and a half hours by boat up river from the native community of

1
 2
 3
 4 554 Infierno on the buffer zone of the Tambopata National Reserve. We collected plant and insect
 5
 6 555 material, extracted DNA, amplified the rDNA cluster, and sequenced material on the MinION
 7
 8
 9 556 platform using the MinKNOW offline software (provided by ONT). The first run generated 17,149
 10
 11 557 reads and the second one 20,167 reads. We generated consensus sequences for five out of the
 12
 13 558 seven analyzed specimens. One plant sample and the grasshopper could not be assembled
 14
 15 559 due to too low read coverage. Moreover, BLAST analysis of the reads assigned to the
 16
 17 560 grasshopper suggested that we had sequenced a mite, instead of the grasshopper DNA. The
 18
 19
 20 561 unidentified insect eggs resulted in a butterfly consensus sequence, possibly a Pierid species.
 21
 22 562
 23
 24 563

25 ***Nanopore based arthropod metabarcoding***



49 564
 50
 51 565 **Figure 7: Relative abundances for nine arthropod species in our four mock communities**
 52
 53 566 **(actual), compared to an Illumina amplicon sequencing protocol, and nanopore protocols**
 54
 55
 56 567 **at 63 °C and 68 °C annealing temperature**
 57
 58 568
 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

569 On average, we recovered 2,645 reads for each Illumina sequenced mock community and
570 1,149 for each nanopore mock community. The optimized Illumina amplicon sequencing based
571 18SrDNA protocol resulted in a very good taxon recovery. All nine taxa were recovered from all
572 four mock communities (Fig. 7). Moreover, the Illumina based protocol allowed for accurate
573 predictions of taxon abundances. The average fold change between input DNA and recovered
574 read count was closely distributed around zero (Supplementary Fig 6). In contrast, the long-read
575 nanopore protocol showed very biased qualitative and quantitative taxon recovery (Fig. 7). On
576 average, only 83.33 % of taxa were recovered per nanopore sequenced mock community.
577 Moreover, the fold change of input DNA and recovered read count were highly biased between
578 taxa. Some taxa were considerably over or underrepresented among the read population. This
579 led to a significantly higher variation of fold change between input DNA and read count
580 compared to the Illumina amplicon based protocol (Levene's test $P < 0.05$; Supplementary Fig.
581 7). A reduction of PCR annealing temperature did result in a considerable increase of Odonata
582 sequences, but overall did not have a strong effect on qualitative (77.78 % of taxa recovered) or
583 quantitative taxon recovery (Fig. 7). The variation of fold change between different PCR
584 annealing temperatures was not significantly different (Levene's test, $P > 0.05$). A reduction of
585 PCR cycle number by 10 also did not yield any significant effect on qualitative (88.89 % of taxa
586 recovered) or quantitative taxon recovery (Supplementary Fig. 7).

587

588 Discussion and Potential implications

589

590 *Phylogenetic and taxonomic utility of long rDNA amplicons*

591 Developments in long-read sequencing hold great promise for molecular taxonomy and
592 phylogenetics across very broad taxonomic scales. We recovered phylogenetic relationships
593 across the eukaryote tree of life, which were mostly consistent with the current state of research
594 (e.g. [52]). Separate orders of arthropods all formed well supported monophyletic groups. 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

595 spider phylogeny was highly congruent with recent work based on whole transcriptomes [35]
596 and multi-amplicon data [53]. Moreover, using the rDNA cluster allowed us to resolve young
597 phylogenetic divergences: the relationships within the recent adaptive radiation of the genus
598 *Tetragnatha* in Hawaii confirmed previous research [59, 60].

599
600 Besides their high phylogenetic utility, long rDNA amplicons showed excellent support for
601 taxonomic hypotheses. All morphologically identified species of Hawaiian *Tetragnatha* were
602 recovered as monophyletic groups. The divergence patterns and taxonomic classifications of
603 spiders based on rDNA were strongly correlated to those based on mitochondrial COI, the most
604 commonly used animal barcode marker [4]. rDNA may thus be ideal to complement
605 mitochondrial barcoding. A universal and variable nuclear marker as a supplement to COI
606 barcoding will be particularly useful in cases of mito-nuclear discordance due to male biased
607 gene flow [10, 61], hybridization [12] or infections with reproductive parasites [11].

608
609 Their high phylogenetic utility across very broad taxonomic categories also provide long rDNA
610 amplicons with a distinct advantage over short read barcoding protocols, which are not well
611 suited to support broad scale phylogenetic hypotheses [62]. The inclusion of long amplicons
612 would make it possible to scale up barcoding from simple taxon assignment to community wide
613 phylogenetic inferences [9]. Recently, the amplification of whole mitochondrial genomes was
614 suggested for animal barcoding [63]. This would increase taxonomic and phylogenetic
615 resolution and thus alleviate some disadvantages of short COI amplicons. However, it is
616 challenging to develop truly universal primers to target mitochondrial genomes across a wide
617 range of taxonomic groups [64]. Moreover, mitochondrial genomes will not allow cases of mito-
618 nuclear discordance to be identified. A straightforward way to achieve highly resolved
619 phylogenies may be the combination of long rDNA amplicon sequencing with multiplex PCRs of
620 short mitochondrial amplicons, to amplify multiple mitochondrial DNA fragments [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

621

622 ***Simple, accurate, universal and cost efficient long-read DNA barcoding***

623

Despite the high raw read error of nanopore data, consensus sequences were highly accurate, and library preparation and sequencing for our protocol are simple and cost efficient. Using a single pair of universal primers, long rDNA amplicons can be amplified across diverse eukaryote taxa. A simple dual indexing approach during PCR allows large numbers of samples to be pooled before library preparation [27]. Only a single PCR is required per specimen, while subsequent cleanup and library preparation can be performed on pooled samples. The simplicity of our approach is additionally highlighted by its effectiveness even under field conditions in a remote rainforest site. Nanopore sequencing technology is affordable and universally available to any laboratory. Our ONT MinION generated about 250,000 reads per run. Aiming for about 1,000 reads per amplified specimen, 250 long rDNA barcodes could be generated in single MinION run. Input DNA amounts for different specimens will have to be carefully balanced to maximize the recovery. The total reagent costs, including PCR, library preparation and sequencing, then amount to less than \$4 for each long barcode sequence generated.

637

638 ***Pitfalls of nanopore based long-read barcoding***

639

While our protocol was generally straightforward and reliable, we found several drawbacks, which require further considerations and optimization. First, it needs to be noted that long rDNA amplification will not be possible with highly degraded DNA molecules, e.g. from historical specimens [66]. Moreover, amplification success of long range PCRs proved less consistent than that for amplification of short amplicons. We observed a complete failure of some PCRs when too high template DNA concentrations were loaded. The long range polymerase may be more sensitive to PCR inhibitors present in some arthropod DNA extractions [67]. PCR conditions will have to be carefully optimized for reliable and consistent amplification. We also

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

647 found that highly universal eukaryote primers may result in undesired amplification, for example
648 plants from beetle and butterfly larval guts, phoretic mites, or fungal sequences. However, as
649 long as the DNA of the target taxon is still dominating the resulting amplicon mixture, this
650 undesired amplification will not affect consensus calling. It may be advisable to check the
651 taxonomic composition of amplicon samples before assembly, e.g. by blasting against a
652 reference library. To avoid unspecific amplification, PCR primers could also be redesigned to
653 exclude certain lineages from amplification. It should also be noted that our approach results in
654 only a single consensus sequence for each processed specimen. As a diploid marker, the rDNA
655 cluster can contain heterozygous positions in some specimens, in particular within the ITS
656 regions. This information is currently lost, and a different assembly approach may be necessary
657 to recover heterozygosity as well. Furthermore, index length and edit distance are also
658 important considerations. We used indexes of 15 bp and with a minimum distance of 10 bp to
659 index both sides of our amplicons. Index edit distance of only 4 bp between samples already led
660 to considerable cross-specimen index bleeding. It may thus be better to increase the length and
661 edit distances of indexes. Indexes of 20 or 30 bp could be easily attached to the 5'-tails of PCR
662 primers without strongly affecting PCR efficiency.

663

664 ***Nanopore based arthropod metabarcoding***

665 It is well known that Illumina amplicon sequencing of short 18SrDNA fragments can yield very
666 accurate qualitative and quantitative taxon recovery in metabarcoding experiments [48], a
667 finding that is confirmed by our results. In contrast, little is known on the performance of long-
668 read nanopore sequencing for community diversity assessments [32]. Our long barcode based
669 approach resulted in the dropout of several taxa and highly skewed relative taxon abundances.
670 Skewed abundances were already found in microbial community analysis using nanopore [32].
671 In the simplest case, primer mismatches may be responsible for biased amplification [32, 68].
672 However, the targeted priming sites in our study were extremely conserved. Also, a change 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

673 PCR cycle number and annealing temperature did not have a strong effect on taxon
674 abundances, as would be expected in the case of PCR priming bias [69]. Another possibility is
675 the preferential amplification of template molecules with a certain GC content by the DNA
676 polymerase [33]. However, we found the GC content of the rDNA cluster to be very stable
677 across taxa. Yet another potential explanation for the differential recovery of taxa in community
678 samples is taxonomic bias in DNA degradation [70], but we do not expect DNA degradation to
679 have played a role in our experiment because we used only high quality DNA extractions
680 (verified by gel electrophoresis) from fresh specimens. The most plausible explanation appears
681 to be that variable rDNA lengths are found between different taxa. It is well known that shorter
682 sequences are amplified preferentially in a PCR, especially after it reaches the plateau stage
683 [71]. Such dominance of shorter amplicons could explain the observed biases very well. In fact,
684 the most abundant taxon in our pools was a spider, which also had the shortest amplicon
685 length. The dominant amplification of shorter sequences may also explain the amplification of
686 plant DNA from a butterfly and a flour beetle larva, as plants showed considerably shorter rDNA
687 amplicons than insects. We found a very high variation of rDNA amplicon length within many
688 taxonomic groups, this could be a considerable problem for long read metabarcoding
689 applications. More research into the causes and possible mitigation of these biases will be
690 required before long-read sequencing can be routinely utilized for metabarcoding applications.

691
692 **Conclusion**

693 Sequencing long dual indexed rDNA amplicons on Oxford Nanopore Technologies' MinION is a
694 simple, cost effective, accurate and universal approach for eukaryote DNA barcoding. Long
695 rDNA amplicons offer high phylogenetic and taxonomic resolution across broad taxonomic
696 scales from kingdom down to species. They also prove to be an excellent complement to
697 mitochondrial COI based barcoding in arthropods. However, despite the long-read advantages
698 in the analysis of separate specimens, we found considerable biases associated with

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

699 sequencing bulk community samples. The observed taxonomic bias is possibly a result of
700 taxon-specific length variation of the rDNA cluster and preferential amplification of species with
701 shorter rDNA. Further research into the sources of the observed bias is required before long
702 rDNA amplicon sequencing can be utilized as a reliable resource for the analysis of bulk
703 samples.

704

705 Availability of source code and requirements

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

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

708 Operating systems: MacOS, Linux and Windows

709

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

711

712

713 Availbity of supporting data

714 The following data supporting the results of this article are available in the [**will be submitted to**
715 **the GlgaScience database**] repository.

716

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

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

721 3. Newick tree files

722 4. Analysis tables for the mock community sequencing experiment, the comparison of genetic
723 distances within and between Hawaiian *Tetragnatha* species for COI and rDNA and the
724 distance between Nanopore based and Illumina based 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

725

726 **Author contributions**

727 HK and SP designed the study. HK, AP, SRK, JYL, VS and JDS collected the specimens.

728 Laboratory work was carried out by HK, AP and SRK and the data were subsequently analyzed

729 by HK, AP, JBH, SRK and SP. The paper was written by HK, AP, JBH, SRK, JYL, VS, JDS,

730 NHP, RGG, SP.

731

732

733 **Abbreviations**

734 ONT: Oxford Nanopore Technologies; PCR: polymerase chain reaction; rDNA: ribosomal DNA;

735 COI: Cytochrome c oxidase subunit I; RTA: retrolateral tibial apophysis

736

737

738 **Competing interests**

739 The authors declare that they have no competing interests.

740

741

742 **Acknowledgements**

743 We thank Taylor Liu for help during laboratory work, and Natalie Graham and Tara Gallant for

744 help during specimen collection. Hitomi Asahara graciously provided access to a laboratory

745 facility and the necessary software for our MinION sequencing run. We thank the State of

746 Hawaii Department of Land and Natural Resources and the Servicio Nacional Forestal y de

747 Fauna Silvestre, who provided collection permits, rainforest Expeditions and Gabriela Orihuela

748 for providing assistance and support with fieldwork in Peru. We thank Anna Holmquist for

749 providing the *Psechrus* sp. Specimen. The specimen was collected with permits from the

750 Indonesian Ministry of Research and Technology (KEMENRISTEK) and in collaboration with

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

751 Pungki Lupiyaningdyah and Anang Achmadi from the Museum Zoologicum Bogoriense and
752 funded by an NSF grant DEB 1457845.

753
754

References

755
756
757 1. Sala, O.E., Chapin, F.S., Armesto, J.J., Berlow, E., Bloomfield, J., Dirzo, R., Huber-
758 Sanwald, E., Huenneke, L.F., Jackson, R.B., and Kinzig, A. (2000). Global biodiversity
759 scenarios for the year 2100. *Science* 287, 1770-1774.
760 2. Pimm, S.L., Jenkins, C.N., Abell, R., Brooks, T.M., Gittleman, J.L., Joppa, L.N., Raven,
761 P.H., Roberts, C.M., and Sexton, J.O. (2014). The biodiversity of species and their rates
762 of extinction, distribution, and protection. *Science* 344, 1246752.
763 3. Rominger, A., Goodman, K., Lim, J., Armstrong, E., Becking, L., Bennett, G., Brewer, M.,
764 Cotoras, D., Ewing, C., and Harte, J. (2016). Community assembly on isolated islands:
765 macroecology meets evolution. *Global ecology and biogeography* 25, 769-780.
766 4. Hebert, P.D., Ratnasingham, S., and de Waard, J.R. (2003). Barcoding animal life:
767 cytochrome c oxidase subunit 1 divergences among closely related species.
768 *Proceedings of the Royal Society of London B: Biological Sciences* 270, S96-S99.
769 5. Schoch, C.L., Seifert, K.A., Huhndorf, S., Robert, V., Spouge, J.L., Levesque, C.A.,
770 Chen, W., Bolchacova, E., Voigt, K., and Crous, P.W. (2012). Nuclear ribosomal internal
771 transcribed spacer (ITS) region as a universal DNA barcode marker for Fungi.
772 *Proceedings of the National Academy of Sciences* 109, 6241-6246.
773 6. China Plant BOL Group, Li, D.-Z., Gao, L.-M., Li, H.-T., Wang, H., Ge, X.-J., Liu, J.-Q.,
774 Chen, Z.-D., Zhou, S.-L., and Chen, S.-L. (2011). Comparative analysis of a large
775 dataset indicates that internal transcribed spacer (ITS) should be incorporated into the

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

776 core barcode for seed plants. *Proceedings of the National Academy of Sciences* 108,
777 19641-19646.

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

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

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

786 10. Krehenwinkel, H., Graze, M., Rödder, D., Tanaka, K., Baba, Y.G., Muster, C., and Uhl,
787 G. (2016). A phylogeographical survey of a highly dispersive spider reveals eastern Asia
788 as a major glacial refugium for Palaearctic fauna. *Journal of Biogeography* 43, 1583-
789 1594.

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

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

797 13. Melo-Ferreira, J., Boursot, P., Suchentrunk, F., Ferrand, N., and Alves, P. (2005).
798 Invasion from the cold past: extensive introgression of mountain hare (*Lepus timidus*)
799 mitochondrial DNA into three other hare species in northern Iberia. *Molecular Ecology*
800 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

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

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

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

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

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

814 19. Tang, C.Q., Leasi, F., Obertegger, U., Kieneker, A., Barraclough, T.G., and Fontaneto, D.
815 (2012). The widely used small subunit 18S rDNA molecule greatly underestimates true
816 diversity in biodiversity surveys of the meiofauna. *Proceedings of the National Academy*
817 *of Sciences* 109, 16208-16212.

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

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

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

825 22. Heeger, F., Bourne, E.C., Baschien, C., Yurkov, A., Bunk, B., Spröer, C., Overmann, J.,
826 Mazzoni, C.J., and Monaghan, M.T. (2018). Long-read DNA metabarcoding of ribosomal
827 rRNA in the analysis of fungi from aquatic environments. *bioRxiv*, 283127.

828 23. Tedersoo, L., Tooming-Klunderud, A., and Anslan, S. (2018). PacBio metabarcoding of
829 Fungi and other eukaryotes: errors, biases and perspectives. *New Phytologist* 217,
830 1370-1385.

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

834 25. Pomerantz, A., Peñafiel, N., Arteaga, A., Bustamante, L., Pichardo, F., Coloma, L.A.,
835 Barrio-Amorós, C.L., Salazar-Valenzuela, D., and Prost, S. (2018). Real-time DNA
836 barcoding in a rainforest using nanopore sequencing: opportunities for rapid biodiversity
837 assessments and local capacity building. *GigaScience* 7, giy033.

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

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

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

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

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

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

850 31. Edwards, A., Debonnaire, A.R., Sattler, B., Mur, L.A., and Hodson, A.J. (2016). Extreme
851 metagenomics using nanopore DNA sequencing: a field report from Svalbard, 78 N.
852 bioRxiv, 073965.

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

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

859 34. Krehenwinkel, H., Wolf, M., Lim, J.Y., Rominger, A.J., Simison, W.B., and Gillespie, R.G.
860 (2017). Estimating and mitigating amplification bias in qualitative and quantitative
861 arthropod metabarcoding. *Scientific reports* 7, 17668.

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

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

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

869 38. Tamura, K., Stecher, G., Peterson, D., Filipowski, A., and Kumar, S. (2013). MEGA6:
870 molecular evolutionary genetics analysis version 6.0. *Molecular Biology and Evolution*
871 30, 2725-2729.

872 39. Straub, S.C., Parks, M., Weitemier, K., Fishbein, M., Cronn, R.C., and Liston, A. (2012).
873 Navigating the tip of the genomic iceberg: Next-generation sequencing for plant
874 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

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

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

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

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

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

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

889 46. Machida, R.J., and Knowlton, N. (2012). PCR primers for metazoan nuclear 18S and
890 28S ribosomal DNA sequences. *PLoS one* 7, e46180.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

940 65. Krehenwinkel, H., Kennedy, S., Rueda, M., and Gillespie, R. (in press). Low cost
941 molecular systematics of entire arthropod communities: Primer sets for rapid multi locus
942 analyses by multiplex PCRs and Illumina amplicon sequencing. . *Methods in Ecology*
943 *and Evolution*.

944 66. Krehenwinkel, H., and Pekar, S. (2015). An analysis of factors affecting genotyping
945 success from museum specimens reveals an increase of genetic and morphological

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 variation during a historical range expansion of a European spider. PLoS one 10,
947 e0136337.

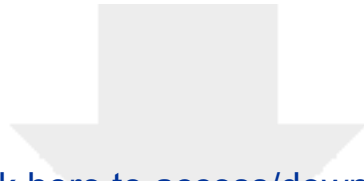
948 67. Margam, V.M., Gachomo, E.W., Shukle, J.H., Ariyo, O.O., Seufferheld, M.J., and
949 Kotchoni, S.O. (2010). A simplified arthropod genomic-DNA extraction protocol for
950 polymerase chain reaction (PCR)-based specimen identification through barcoding.
951 Molecular biology reports 37, 3631-3635.

952 68. Sipos, R., Székely, A.J., Palatinszky, M., Révész, S., Márialigeti, K., and Nikolausz, M.
953 (2007). Effect of primer mismatch, annealing temperature and PCR cycle number on
954 16S rRNA gene-targetting bacterial community analysis. FEMS Microbiology Ecology
955 60, 341-350.

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

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

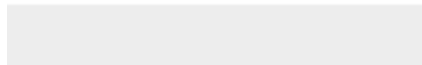
963 71. Wattier, R., Engel, C., Saumitou-Laprade, P., and Valero, M. (1998). Short allele
964 dominance as a source of heterozygote deficiency at microsatellite loci: experimental
965 evidence at the dinucleotide locus Gv1CT in *Gracilaria gracilis* (Rhodophyta). Molecular
966 Ecology 7, 1569-1573.



Click here to access/download

Supplementary Material

SupplementaryTable1_SampleList.xlsx





[Click here to access/download](#)

Supplementary Material

[Ir-barcoding_publ_SUPPL_PRO_final.docx](#)

