






## RESEARCH

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

Henrik Krehenwinkel <sup>1,2,5,\*</sup>, Aaron Pomerantz <sup>3</sup>, James B. Henderson<sup>4,5</sup>, Susan R. Kennedy<sup>2</sup>, Jun Ying Lim<sup>2,3</sup>, Varun Swamy<sup>6</sup>, Juan Diego Shoobridge<sup>7</sup>, Natalie Graham<sup>2</sup>, Nipam H. Patel<sup>3,8</sup>, Rosemary G. Gillespie<sup>2</sup> and Stefan Prost <sup>3,9,10,\*</sup>

<sup>1</sup>Department of Biogeography, Trier University, Faculty of Regional and Environmental Sciences, Trier 54286, Germany, <sup>2</sup>Department of Environmental Science, Policy and Management, University of California, Berkeley, California, 94720, USA, <sup>3</sup>Department of Integrative Biology, University of California, Berkeley, California, 94720, USA, <sup>4</sup>Institute for Biodiversity Science and Sustainability, California Academy of Sciences, 55 Music Concourse Drive, San Francisco, California, 94118, USA, <sup>5</sup>Center for Comparative Genomics, California Academy of Sciences, 55 Music Concourse Drive, San Francisco, California, 94118, USA, <sup>6</sup>San Diego Zoo Institute for Conservation Research, 15600 San Pasqual Valley Road, Escondido, California, 92027, USA, <sup>7</sup>Applied Botany Laboratory, Research and development Laboratories, Cayetano Heredia University, Av. Honorio Delgado 430, Urb Ingenieria, Lima, Perú, <sup>8</sup>Department of Molecular and Cell Biology, University of California, Berkeley, California, 94720, USA, <sup>9</sup>Research Institute of Wildlife Ecology, Department of Integrative Biology and Evolution, University of Veterinary Medicine, Vienna, Austria and <sup>10</sup>South African National Biodiversity Institute, National Zoological Garden, Pretoria, 0184, South Africa

\*Correspondence address. Henrik Krehenwinkel, E-mail: [krehenwinkel@uni-trier.de](mailto:krehenwinkel@uni-trier.de)  <http://orcid.org/0000-0001-5069-8601>; Mailing address: Department of Biogeography, Trier University, Faculty of Regional and Environmental Sciences, Trier 54286, Germany; Stefan Prost, E-mail: [stefan.prost@protonmail.com](mailto:stefan.prost@protonmail.com)  <http://orcid.org/0000-0002-6229-3596>; Mailing address: LOEWE - Center for Translational Biodiversity Genomics, Senckenberg, Frankfurt 60325, Germany

## Abstract

**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

Received: 4 May 2018; Revised: 30 October 2018; Accepted: 10 January 2019

© The Author(s) 2019. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

solution is offered by sequencing long ribosomal DNA (rDNA) amplicons on the MinION platform (Oxford Nanopore Technologies). **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. **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.

**Keywords:** biodiversity; ribosomal; eukaryotes; long DNA barcodes; Oxford Nanopore Technologies; MinION; metabarcoding

## Background

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

Comparative ecological studies across multiple taxa have been greatly simplified by molecular barcoding [4], where species identifications are based on short polymerase chain reaction (PCR) amplicon “barcode” sequences. Different barcode marker genes have been established across the tree of life [5, 6], with mitochondrial cytochrome oxidase subunit I (COI) commonly used for animal barcoding [4]. The availability of large sequence reference databases and universal primers, together with its uniparental inheritance and fast evolutionary rate, make COI a useful marker to distinguish even recently diverged taxa. In recent years, DNA barcoding has greatly profited from the emergence of next-generation sequencing (NGS) technology. Current NGS platforms enable the parallel generation of barcodes for hundreds of specimens at a fraction of the cost of Sanger sequencing [7]. Furthermore, NGS technology has enabled metabarcoding, the sequencing of bulk community samples, which allows scoring the diversity of entire ecosystems [8].

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

placement of mitochondrial genomes (see, e.g., [12, 13]), resulting in shared mitochondrial variation between species.

Considering this background, it would be desirable to complement mitochondrial DNA based barcoding with additional information from the nuclear genome. An ideal nuclear barcoding marker should possess sufficient variation to distinguish young species pairs but also provide support for phylogenetic hypotheses between divergent lineages. Moreover, the marker should be present across a wide range of taxa, and amplification should be possible using universal primers. A marker that fulfils all the above requirements is the nuclear ribosomal DNA (rDNA). As an essential component of the ribosomal machinery, rDNA is a common feature across the tree of life from microbes to higher eukaryotes [14]. All eukaryotes share homologous transcription units of the 18S, 5.8S, and 28S-rDNA genes, which include two internal transcribed spacers (ITS1 and ITS2) [15]. Due to varying evolutionary constraints acting on different parts of the rDNA, it consists of regions of extreme sequence conservation, which are interrupted by highly variable stretches [16]. While some rDNA gene regions are entirely conserved across all eukaryotes, the two ITS sequences are distinguished by such rapid evolutionary change that they separate even lineages within species [5, 17]. rDNA markers thus offer taxonomic and phylogenetic resolution at a very broad taxonomic scale. As an essential component of the translation machinery, nuclear rDNA is required in large quantities in each cell. It is thus present in multiple copies across the genome [15] and is readily accessible for PCR amplification. Due to the above advantages, rDNA already is a popular and widely used marker for molecular taxonomy and phylogenetics in many groups of organisms [5, 6, 15, 17, 18]. However, its presence in multiple copies across the genome may also make rDNA susceptible to the emergence of paralogs and pseudogenization, which could affect taxonomic and phylogenetic utility.

Spanning about 8 kb, the ribosomal cluster is fairly large, and current barcoding protocols, e.g., using Sanger sequencing or Illumina amplicon sequencing, can only target short sequence stretches of 150–1,000 bp. Such short stretches of 28S and 18S are often too conserved to identify young species pairs [19]. The ITS regions, on the other hand, are so variable that they cannot be properly aligned across deeply divergent lineages. Moreover, ITS sequences can show considerable length variation between taxa. This holds particularly true for the ITS1 region whose length can vary between a few 100 and more than 1,000 bp [20]. Considerable, but less pronounced, length variation can also be observed in ITS2. Short amplicon-based sequencing approaches are limited to a maximum fragment length of about 500 bp. As ITS priming sites have to rest in the conserved flanking rDNA gene sequences, the resulting amplicon often exceeds

this length and thus cannot be used for short amplicon-based barcoding in some taxa.

Consequently, it would be ideal to amplify and sequence a large part of the ribosomal cluster in one fragment. A solution to sequence the resulting long amplicons is offered by recent developments in third-generation sequencing platforms, which now enable researchers to generate ultra-long reads, of up to 800 kb [21]. Recently, amplicons of several kilobases of the rDNA cluster were sequenced using Pacific Bioscience (PacBio) technology to explore fungal community composition [22, 23]. With its circular consensus sequencing technology, PacBio allows the generation of very accurate consensus reads. However, while PacBio sequencing is well suited for long amplicon sequencing, it is currently not readily available to every laboratory due to the high cost and limited distribution of sequencing machines. PacBio sequencers are also bulky and cannot be used outside of conventional laboratory settings.

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

Considering this background, we explore the feasibility of Nanopore sequencing of long rDNA amplicons as a simple, cost-efficient DNA barcoding approach for animals and other eukaryote taxa. We compile a workflow from PCR amplification, to library preparation, to demultiplexing and consensus calling (see Fig. 1 for an overview). We explore the error profile of Nanopore consensus sequences and introduce MiniBar, a new software to demultiplex dual indexed Nanopore amplicon sequences. We test the utility of the ribosomal cluster for molecular taxonomy and phylogenetics across divergent plant and animal taxa. A particular focus of our analysis is arthropods, the most diverse group in the animal kingdom [30], which are highly threatened by current mass extinctions [31]. Using a dataset of spiders, we compare the taxonomic resolution of the ribosomal cluster with that offered by molecular barcoding using mitochondrial COI, the currently preferred barcode marker for arthropods. ONT's MinION is a portable sequencer, and Nanopore-based DNA barcoding has been applied in remote sites outside of conventional labs (see, e.g., [25, 30, 31]). As mentioned above, the MinION is portable and can be used for DNA barcoding in field settings. Such field-based applications confront researchers with additional complexities and challenges. To highlight the simplicity of our approach, we tested it under field conditions and generated long rDNA barcode sequences using a miniaturized mobile laboratory in a Peruvian rainforest.

We also tested the efficacy of long-amplicon rDNA sequencing for metabarcoding of bulk community samples. A study of bacterial communities [32] suggests Nanopore long-amplicon

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

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

## Data Description and Analyses

### DNA extraction, PCR, and library preparation

We analyzed 114 specimens of eukaryotes including 17 insect and 42 spider species, 2 annelid and 9 plant species (Supplementary Table S1). Some feeder insects and the annelids were purchased at a pet store. The remaining specimens were collected in oak forest on the University of California Berkeley's campus or in native rainforests of the Hawaiian Archipelago (under the Hawaii DLNR permit: FHM14-349). We focused our arthropod sampling on spiders, which are ubiquitous and essential predators in all terrestrial ecosystems. Recent phylogenomic work [35] provided us with a solid baseline to test the efficiency of rDNA amplicons for phylogenetic and taxonomic purposes. We included a taxonomically diverse collection of 16 spider families from the Araneoidea, the retrolateral tibial apophysis (RTA) clade, and a haplogyne outgroup species. Within spiders, we additionally focused on the genus *Tetragnatha*, which has undergone a striking adaptive radiation on Hawaii.

DNA was extracted from each sample using the Qiagen Archivepure kit (Qiagen, Valencia, CA) according to the manufacturer's protocols. The DNA integrity was checked on an agarose gel. Only samples with high DNA integrity were used for the following PCRs. All DNA extracts were quantified using a Qubit fluorometer using the high-sensitivity dsDNA assay (Thermo Fisher, Waltham, MA) and diluted to concentrations of 20 ng/ $\mu$ L. We designed a primer pair of each 27 bases to amplify a ~4000 bp fragment of the ribosomal DNA, including partial 18S and 28S as well as full ITS1, 5.8S, and ITS2 sequences (18S.F4 GGCTACCA-CATCYAARGAAGGCAGCAG and 28S.R8 TCGGCAGGTGAGTYGT-TRCACAYTCCT). The primers were designed using alignments of partial 18S and 28S sequences of ~1,000 species of eukaryotes, with a focus on animals (Supplementary Fig. S1). The primers targeted highly conserved regions across all analyzed taxa. Degenerate sites were incorporated to account for variation. We aimed for high annealing temperatures (65°C–70°C) to impose stringent amplification. These were calculated using the NEB Tm Calculator [36].

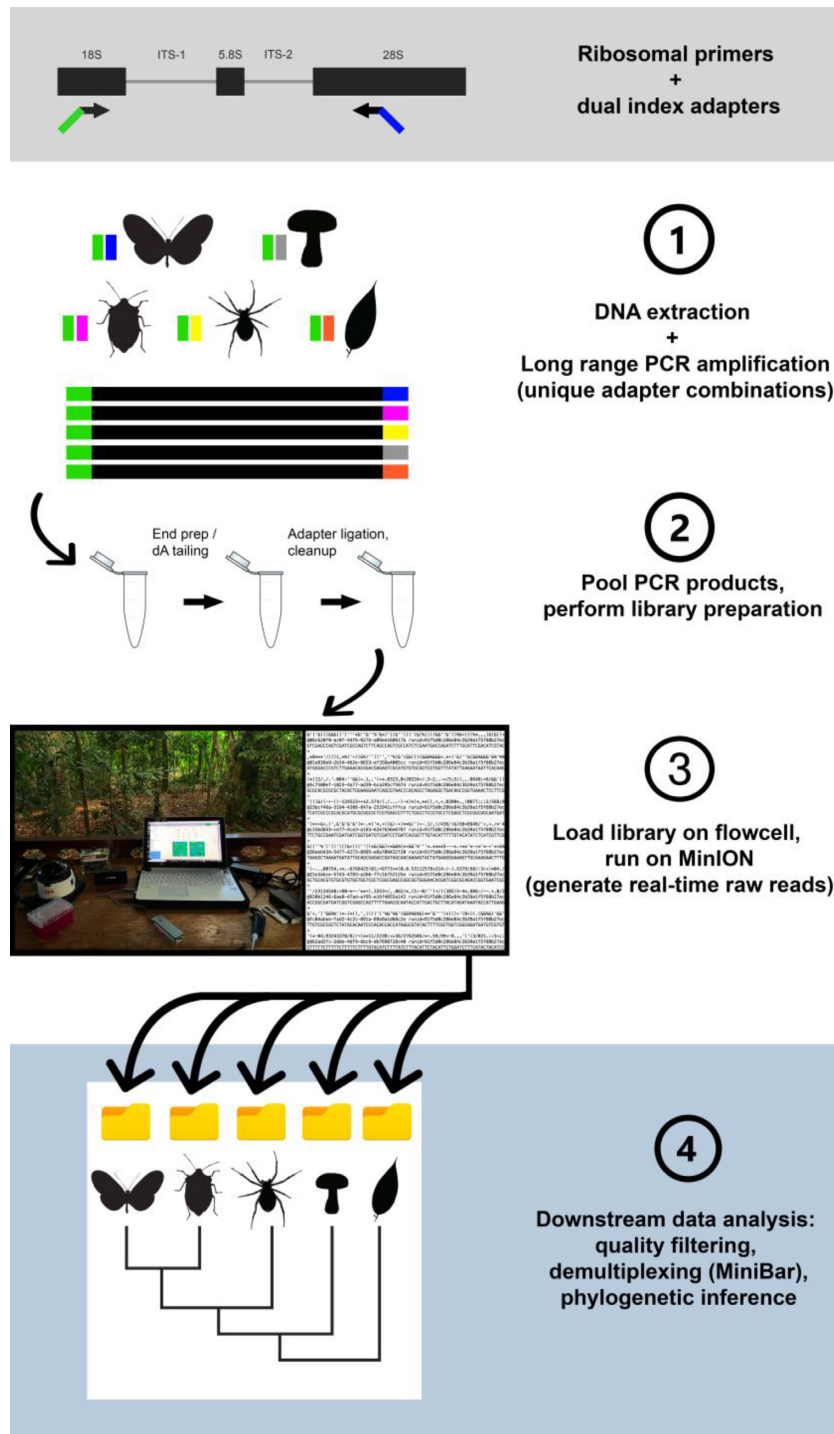


Figure 1: Workflow for the design, amplification, and sequencing of the ribosomal DNA cluster.

To index every PCR amplicon separately, we used a dual indexing strategy with each primer carrying a unique 15 bp index sequence at its 5'-tail. Index sequences were designed using Barcode Generator [37] with a minimum distance of 10 bases between each index. A total of 15 forward and 16 reverse indexes were designed. Every sample was amplified separately using the Q5 Hot Start High-Fidelity 2X Master Mix (NEB, Ipswich, MA) in 15  $\mu$ L reactions, at 68°C annealing temperature, with 35 PCR cycles and using 50 ng of template DNA per PCR. All PCR products

were quantified on an agarose gel, based on band intensity on the gel, using the Gel Doc XR System with the Quantity One software (Bio-Rad, CA) and then pooled.

Then, 100  $\mu$ L of the final pool were cleaned from residual primers by 0.75 X AMPure Beads XP (Beckman Coulter, Brea, CA). DNA library preparation was carried out according to the 1D PCR barcoding amplicons SQK-LSK108 protocol (Oxford Nanopore Technologies, Oxford, UK). Barcoded DNA products were pooled with 5  $\mu$ L of DNA GS (a positive control provided by ONT) and

an end-repair was performed (NEB-Next Ultra II End-prep reaction buffer and enzyme mix), then purified using AMPure XP beads. Adapter ligation and tethering was carried out with 20  $\mu\text{L}$  Adapter Mix and 50  $\mu\text{L}$  of NEB Blunt/TA ligation Master Mix. The adapter-ligated DNA library was then purified with AMPure beads XP, followed by the addition of Adapter Bead binding buffer, and finally eluted in 15  $\mu\text{L}$  of Elution Buffer. Each R9 flow cell was primed with 1,000  $\mu\text{L}$  of a mixture of Fuel Mix and nuclease-free water. Twelve  $\mu\text{L}$  of the amplicon library were diluted in 75  $\mu\text{L}$  of running buffer with 35  $\mu\text{L}$  RBF, 25.5  $\mu\text{L}$  LLB, and 2.5  $\mu\text{L}$  nuclease-free water and then added to the flow cell via the SpotON sample port. The “NC.48Hr.sequencing.FLO-MIN107.SQK-LSK108.plus.Basecaller.py” protocol was initiated using the MinION control software, MinKNOW.

### Field trial in the Amazon rainforest

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

### Bioinformatics

#### Raw data processing and consensus calling

The fastq files generated by the ONT software MinKNOW were de-multiplexed using MiniBar (see description below), with index edit distances of 2, 3, and 4 and a primer edit distance of 11. Next, the reads were filtered for quality (>13) and size (>3 kb) using Nanofilt [38, 39]. Individual consensus sequences were created using Allele Wrangler [40] for demultiplexed fastq files with a minimum coverage of 30. Error correction was performed using RACON [41, 42]. To do so, we first mapped all the reads back to the consensus using minimap [43]. We performed two cycles of running minimap and RACON. Final consensus sequences were compared against the National Center for Biotechnology Information database using Basic Local Alignment Search Tool n (BLASTn) to check if the taxonomic assignment was correct.

We performed multiple tests to validate and optimize the consensus accuracy of long-amplicon barcode sequences. To comparatively assess the accuracy, we used consensus sequences of short 18S and 28S rDNA amplicons, which were previously generated using Illumina amplicon sequencing for the 47 analyzed Hawaiian *Tetragnatha* specimens (Kennedy, unpublished data). These sequences were aligned with the respective stretches of our Nanopore consensus sequences using ClustalW in MEGA (MEGA Software, RRID:SCR.000667) [44]. All alignments were then visually inspected and edited manually, where necessary. Pairwise distances between Illumina and Nanopore consensus were calculated in MEGA.

To measure consensus accuracy over the whole ribosomal amplicon, we utilized genome skimming data [45] for six Hawaiian *Peperomia* plant species (Lim et al, unpublished data). A total of 150 bp paired-end TruSeq gDNA shotgun libraries for the six *Peperomia* samples were sequenced on a single HiSeq v4000 lane (Illumina, San Diego, CA). The resulting paired-end reads were trimmed and filtered using Trimmomatic v0.36 (Trimmomatic, RRID:SCR.011848) [46] and mapped to their respective Nanopore consensus sequences using bowtie2 (Bowtie, RRID:SCR.005476) [47] under default parameter values and allowing for minimum and maximum fragment size of 200 and 700 bases, respectively. Mapping coverage of Illumina reads to Nanopore consensus sequences ranged from 150 to 600 X with a mean of ~300 X across all six samples. We called Illumina read-based consensus sequences for each *Peperomia* species using bcftools [48] and aligned them with the previously generated Nanopore consensus sequences. Pairwise genetic distances were then calculated in MEGA as described above. We performed two independent distance calculations: excluding indels, i.e., only using nucleotide substitutions to estimate genetic distances, and including indels as additional characters.

Our demultiplexing software allows flexible edit distances to identify forward and reverse indexes from Nanopore reads. Due to the high raw read error rate, edit distances that are too large could lead to crossover between samples during demultiplexing. This crossover could possibly affect the accuracy of the called consensus sequence. On the other hand, edit distances that are too stringent 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 relevant 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 relevant specimen 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. The exact position of gene sequences was identified by downloading full length 18S, 5.8S, and 28S sequences from GenBank and then aligning them against the amplicons. Due to the deep divergence in the eukaryote dataset, the highly variable ITS sequences could not be aligned and were excluded. For the analyses of spiders, we retained both ITS sequences and aligned the whole rDNA amplicon. Appropriate models of sequence evolution for each gene fragment of the rDNA cluster were identified using PartitionFinder [49]. Phylogenies were built using MrBayes [50], with four heated chains, a chain length of 1,100,000, subsampling every 200 generations, and a burnin length of 100,000.

Focusing on the endemic Hawaiian *Tetragnatha* species, we also tested the utility of the ribosomal cluster for taxonomic identification, as we also had COI barcodes available for these species. Our dataset contained ribosomal DNA sequences for 47 specimens in 16 species, which had been identified morphologically before barcoding. We calculated pairwise genetic distances between and within all species for the whole ribosomal cluster and for each separate gene region of the rDNA cluster using MEGA. As the 18S and 5.8S did not yield any species-level resolution within Hawaiian *Tetragnatha*, they were not analyzed separately. To compare the taxonomic resolution of the ribosomal cluster with that of the commonly used mitochondrial COI, we calculated inter- and intraspecific distances for an alignment of 418 bp of the COI barcode region for the same spider specimens (Kennedy et al., unpublished data). We performed a Mantel test using the R package *ade4* [51] to test for a significant correlation between COI and ribosomal DNA-based distances. A comparison of intraspecific and interspecific distances for mitochondrial COI and ribosomal DNA also allowed us to test for the presence of a barcode gap.

#### Nanopore-based arthropod metabarcoding

To test for the possibility of estimating arthropod community composition from Nanopore sequencing, we prepared four mock communities of different amounts of DNA extracts from nine species of arthropods from different orders (see Supplementary Table S2). It should be noted that with representatives of nine different orders, these community samples were highly simplified and are not necessarily representative of a natural arthropod community. Due to the high error rate of individual reads, we did not know if, and how, the MinION's high error rate would affect taxonomic assignment, hence we decided to limit our current analysis to these simplified communities.

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

In order to compare our results with those from an optimized Illumina short read protocol, we amplified all samples for a ~300 bp fragment of the 18S rDNA using the primer pair 18S2F/18S4R [52]. Amplification and library preparation were performed as described in [53] using Qiagen Multiplex PCR kits. The 18S amplicon pools were sequenced on an Illumina MiSeq using V3 chemistry and 2 × 300 bp reads. Sequence quality filtering, read merging, and primer trimming were performed as described in [34].

A library of 18S sequences for all included arthropod species (from [34]) was used as a reference database to identify the recovered sequences using BLASTn [54], with a minimum e-value of 10<sup>-4</sup> and a minimum overlap of 95%. Despite the high raw error rate of Nanopore reads, taxonomic status of sequences could be assigned using BLAST, as our pools contained members of highly divergent orders. We compared the qualitative (number of species) and quantitative (abundance of species) recovery of taxa from the communities by Nanopore long-amplicon and Illumina short read data. To estimate the recovery of taxon abundances, we calculated a fold change between input DNA amount and recovered reads for each taxon and mock community. A fold change of zero corresponded to a 1:1 association of taxon abundance and read count, while positive or negative values indicated higher or lower read counts than the taxon's actual abundance.

#### MiniBar

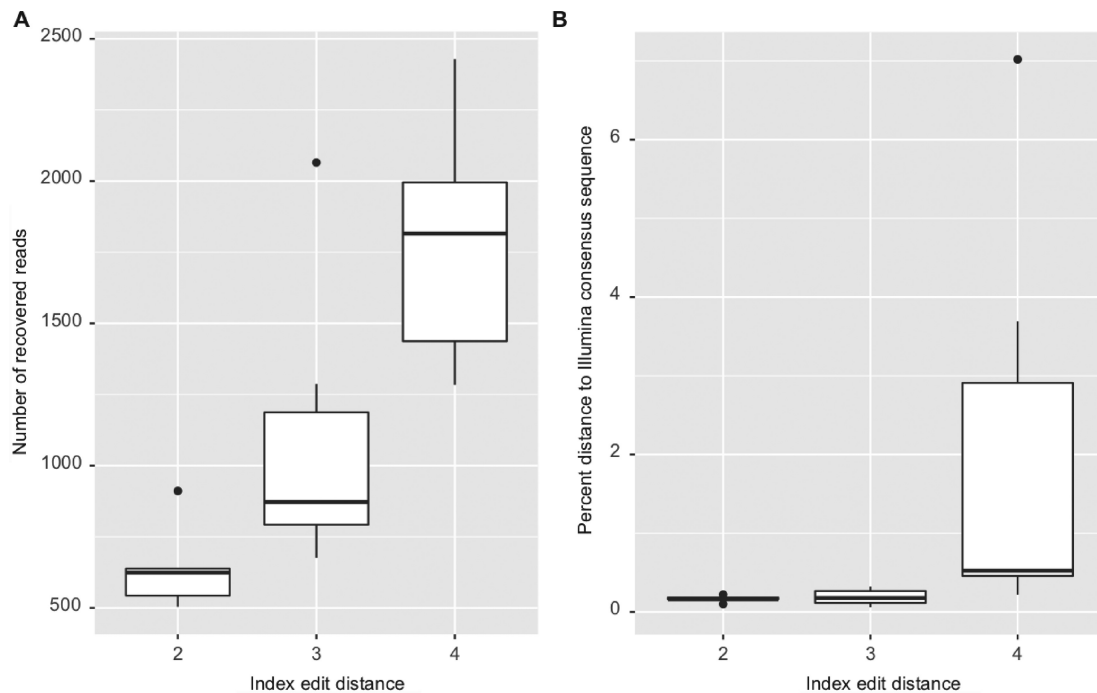
We created a de-multiplexing software, called MiniBar. It allows customization of search parameters to account for the high read error rates and has built-in awareness of the dual barcode and primer pairs flanking the sequences. MiniBar takes as input a tab-delimited barcode file and a sequence file in either fasta or fastq format. The barcode file contains, at a minimum, sample name, forward barcode, forward primer, reverse barcode, and reverse primer for each of the samples potentially in the sequence file. The software searches for barcodes and for a primer, each permitting a user-defined number of errors, an error being a mismatch or indel. Error count to determine a match can either be a percentage of each of their lengths or can be separately specified for barcode and primer as a maximum edit distance [55]. Output options permit saving each sample in its own file or all samples in a single file, with the sample names in the fasta or fastq headers. The found barcode primer pairs can be trimmed from the sequence or can remain in the sequence distinguished by case or color. MiniBar, written in 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 [56]. MiniBar can be found at [57] along with test data.

## Results

### Sequencing, specimen recovery, and consensus quality

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

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



**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.

mealworm (*Zophobas morio*) generated consensus sequences for plants. In most of these samples, the targeted arthropod species was either extremely underrepresented among the read populations or completely absent.

A comparison of our consensus sequences for 47 Hawaiian specimens of the spider genus *Tetragnatha* with short Illumina amplicon sequencing-derived 18S and 28S rDNA sequences suggests a very high consensus accuracy. Except for a single specimen, with a single substitution error, all Nanopore-based consensus sequences were completely identical to the Illumina-based consensus. However, the corresponding 18S and 28S fragments did not contain long stretches of homopolymer sequences, where Nanopore raw read errors are known to accumulate [58]. Despite containing several homopolymers, the Nanopore derived *Peperomia* consensus sequences were highly accurate (Supplementary Fig. S2). Including gaps in the alignment, an average distance of 0.165% to Illumina-based consensus sequences was found. Errors were clustered in indel regions (Supplementary Fig. S3). After excluding gaps, the average distance dropped to 0.102%.

We found only a small effect of sequence coverage on consensus assembly accuracy (Supplementary Fig. S4). Even at 10-fold coverage, a low average distance of 0.257% to Illumina consensus sequences was observed. However, at 20-fold coverage, the average distance significantly decreased to 0.128% (pairwise Wilcoxon test, FDR-corrected  $P < 0.05$ ). A slight, but not significant, decrease of distance was observed with increasing coverage, with optimal consensus accuracy at 300-fold coverage (0.031% distance). At coverages larger than 300, the consensus accuracy slightly decreased (average distance of 0.103% at 800 X coverage), but this change was not significant.

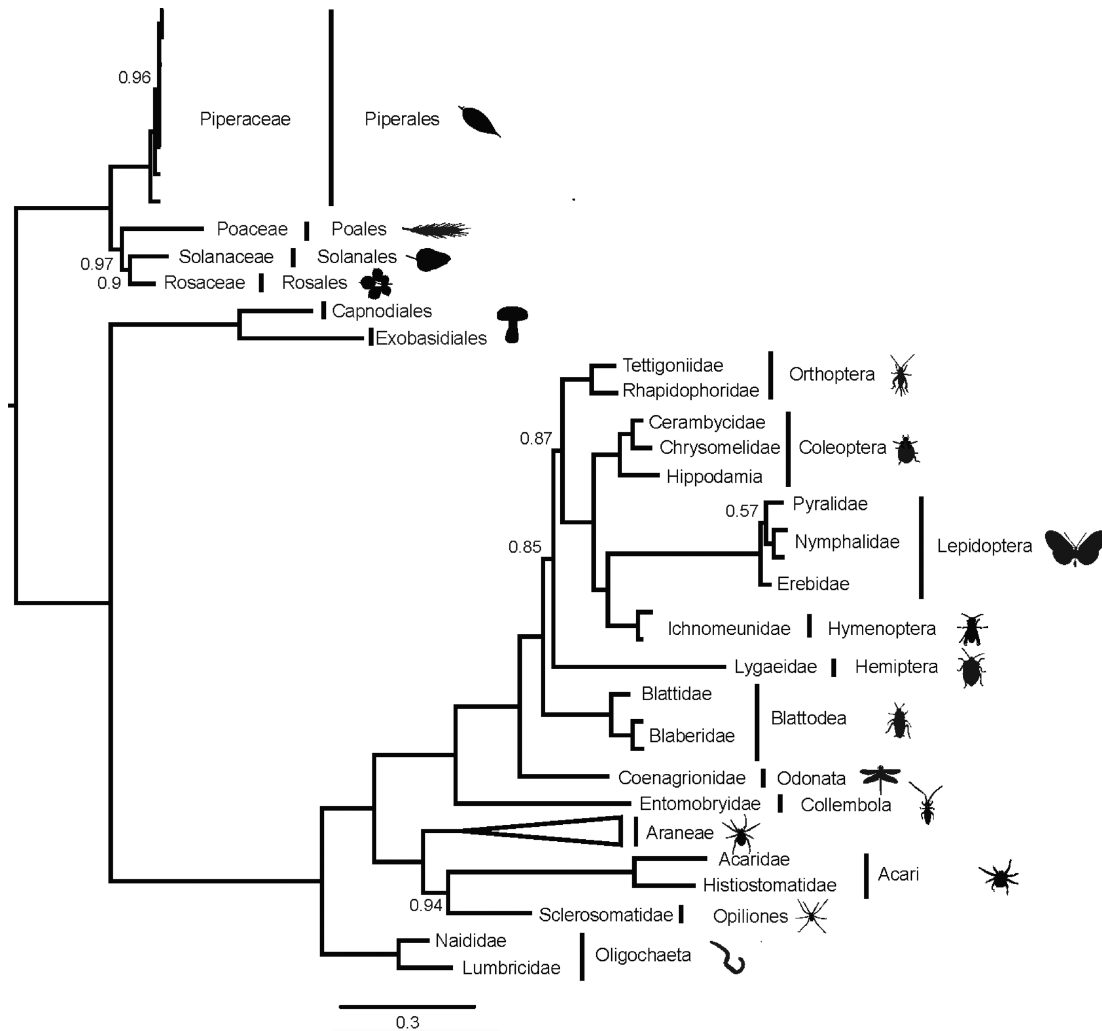
The length of the rDNA amplicon was quite variable between taxa. Arachnids, hexapods, and magnoliopsid plant specimens all showed significantly different amplicon lengths (pairwise

Wilcoxon test, FDR-corrected  $P < 0.05$ ). The length difference was found for the actual gene sequences (18S, 5.8S, 28S: plants:  $2781 \pm 4.96$ ; hexapods:  $3154 \pm 50.35$ ; arachnids:  $3047 \pm 10.77$ ; Supplementary Fig. S5A) as well as including the ITS sequences (plants:  $3243 \pm 11.78$ ; hexapods:  $4192 \pm 498.05$ ; arachnids:  $3644 \pm 129.07$ , Supplementary Fig. S5B). While most spiders showed very stable length distributions for the rDNA amplicon length (average length  $\pm$  standard deviation across all Araneae:  $3629 \text{ bp} \pm 81$ ), several hexapod orders had rDNA sequences of more variable length (Coleoptera:  $4488 \text{ bp} \pm 352$ ; Lepidoptera:  $4363 \text{ bp} \pm 603$ ).

In contrast to the variable length of the rDNA cluster, we found a very stable GC content across the whole taxonomic spectrum ( $46.75 \pm 2.67\%$  across all taxa). GC content of magnoliopsid plants, hexapods, and arachnids was highly similar (plants:  $46.01 \pm 1.66\%$ ; hexapods:  $46.67 \pm 3.73\%$ ; arachnids:  $46.93 \pm 2.47\%$ ) (Supplementary Fig. S5c).

### Phylogenetic reconstruction

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 the latest phylogenomic work [59]. 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 be-



**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.

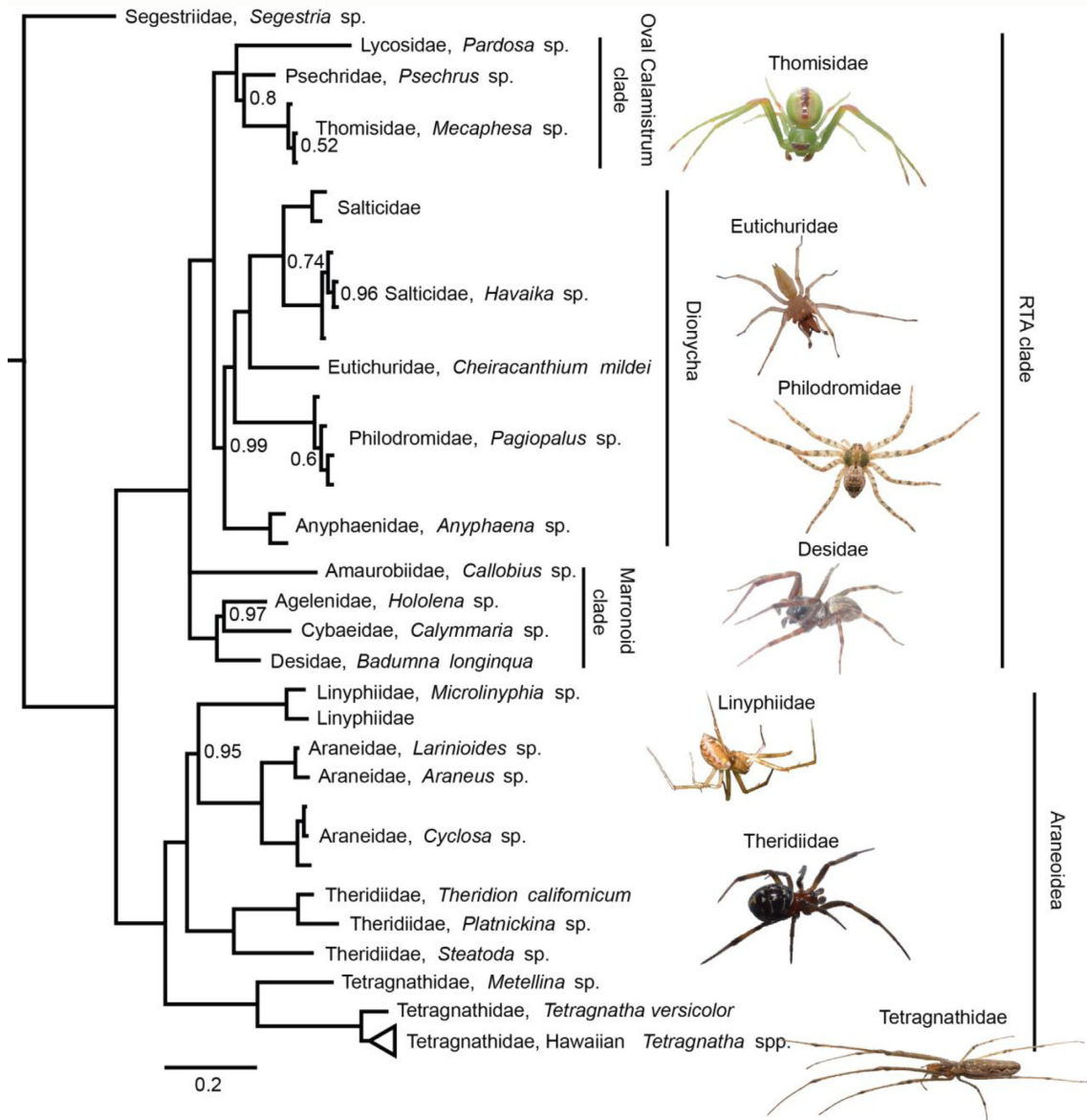
tween these three orders were comparatively low ( $\sim 0.85$ ). Finally, holometabolan insects (Hymenoptera, Coleoptera, and Lepidoptera) were recovered as monophyletic. The two Acari species, together with Opiliones, formed the sister clade to the monophyletic Araneae clade.

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 [60] 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, Psecridae, Salticidae, and Thomisidae) and the Araneoidea (Araneidae, Linyphiidae, Tetragnathidae, and 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 Psecridae and Thomisidae, and 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.

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 [61–64]: 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 [65]: “large brown” (*T. quasimodo*) with tree bark, “small brown” (*T. anuenuae*, *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 ecomorphs clustered phylogenetically, small brown species appeared in three separate clades on the tree. Within the orb weaving clade, *T. hawaiiensis*, a generalist species that occurs on all of the Hawaiian islands, fell basal. The characteristic web struc-





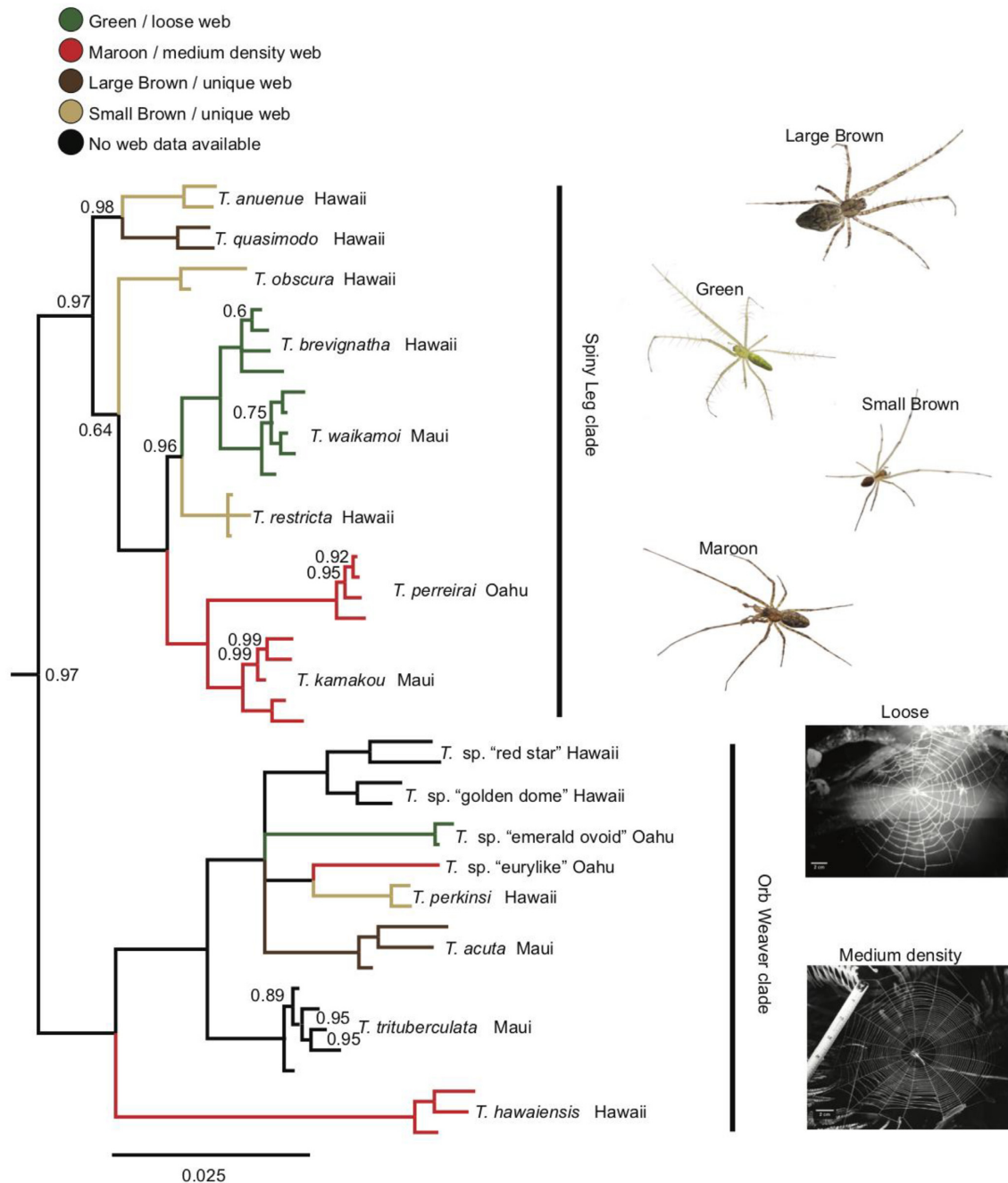
**Figure 4:** Bayesian consensus phylogeny of 83 spiders in 16 families, based on a 4,214 bp alignment of 18S, ITS1, 5.8S, ITS2, and 28S. The phylogeny is rooted using the basal haploglyne *Segestria* sp. The clade containing Hawaiian members of the genus *Tetragnatha* is collapsed (the uncollapsed clade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

tures of some of these species have been documented [66, 67]. We found a pattern of apparent convergence in web structure for some species. *Tetragnatha* sp. “emerald ovoid” spins a loose web with widely spaced rows of capture silk. *Tetragnatha 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.” *Tetragnatha 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 [67] and are thus classified as “unique”.

### Taxonomic resolution

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, 9,999 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, 6C). 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 S3). 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



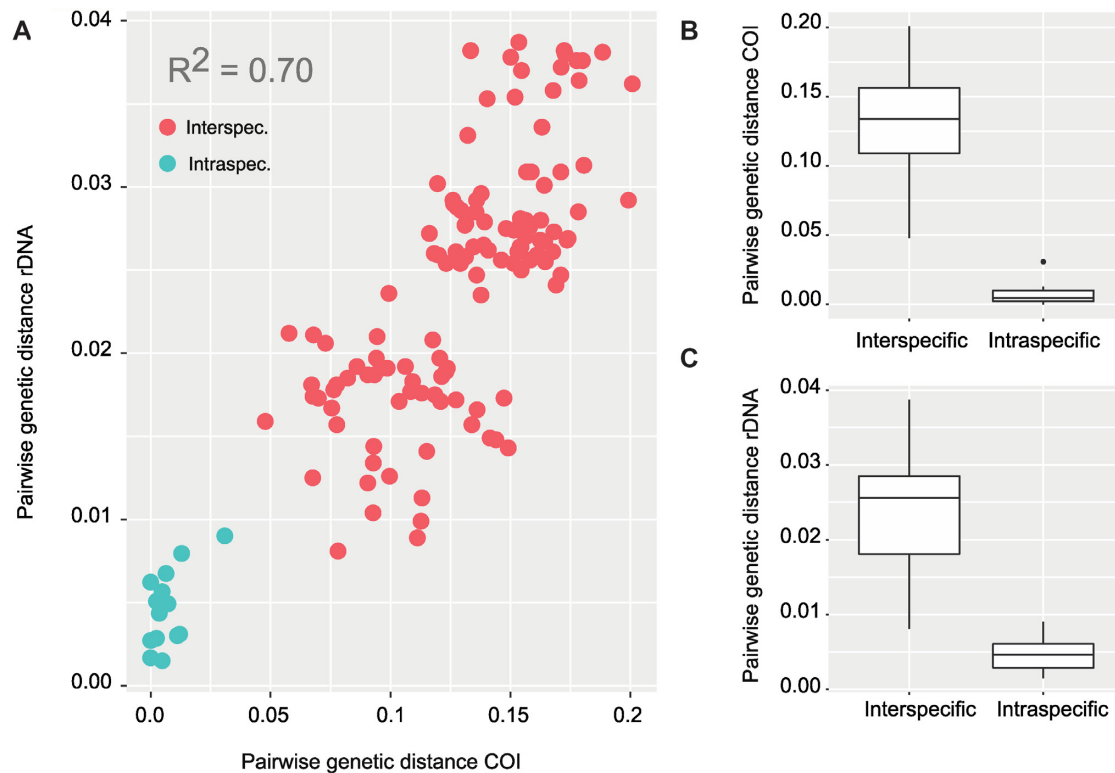
**Figure 5:** Section of the same phylogeny as Fig. 4, with expansion of the clade of 16 Hawaiian *Tetragnatha* species. Different “Spiny Leg” ecomorphs and web architectures are indicated by branch coloration. Only posterior probability values below 1 are displayed.

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. S6). While the 28S rDNA showed considerably lower distances than COI, the distances for ITS1 and ITS2 were more comparable to COI (Supplementary Fig. S6b-6d). Yet, interspecific and intraspecific distances for COI were significantly different from those

for any part of the rDNA cluster (pairwise Wilcoxon test, FDR-corrected  $P < 0.05$ ).

### Field trial in the Amazon rainforest

On 26 March 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. S7). This field site is a “Terra firme” forest in the sector of “Condonado,” approximately two and a half hours by boat up river from the native community of Inferno on the buffer zone of the Tambopata National Reserve. We col-



**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) Inter- and intraspecific genetic distances for the same spider species based on mitochondrial COI and (C) the whole rDNA amplicon.

lected plant and insect material, extracted DNA, amplified the rDNA cluster, and sequenced material on the MinION platform using the MinKNOW offline software (provided by ONT). The first run generated 17,149 reads and the second one generated 20,167 reads. We generated consensus sequences for five of the seven analyzed specimens. One plant sample and the grasshopper could not be assembled due to too low read coverage. Moreover, BLAST analysis of the reads assigned to the grasshopper suggested that we had sequenced a mite, instead of the grasshopper DNA. The unidentified insect eggs resulted in a butterfly consensus sequence, possibly a pierid species.

### Nanopore-based arthropod metabarcoding

On average, we recovered 2,645 reads for each Illumina sequenced mock community and 1,149 for each Nanopore mock community. The optimized Illumina amplicon sequencing based 18SrDNA protocol resulted in a very good taxon recovery. All nine taxa were recovered from all four mock communities (Fig. 7). Moreover, the Illumina-based protocol allowed relatively accurate predictions of taxon abundances. Even though no taxon's actual abundance was predicted by Illumina amplicon data, the average fold change between input DNA and recovered read count was closely distributed around zero (Supplementary Fig. S8). In contrast, the long-amplicon Nanopore protocol showed very biased qualitative and quantitative taxon recovery (Fig. 7). On average, only 83.33% of taxa were recovered per Nanopore sequenced mock community. Moreover, the fold changes of input DNA and recovered read count were highly biased between taxa. Some taxa were considerably over- or underrepresented among the read population. This led to a signif-

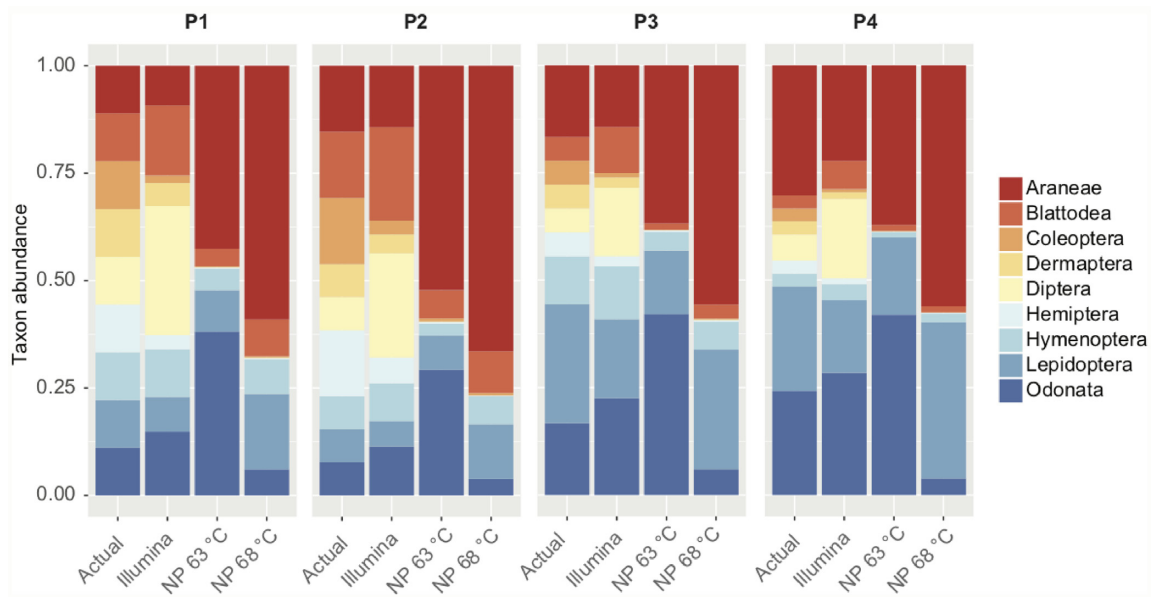
icantly higher variation of fold change between input DNA and read count compared to the Illumina amplicon-based protocol (Levene's test,  $P < 0.05$ ; Supplementary Fig. S8). A reduction of PCR annealing temperature did result in a considerable increase in Odonata sequences but overall did not have a strong effect on qualitative (77.78% of taxa recovered) or quantitative taxon recovery (Fig. 7). The variation of fold change between different PCR annealing temperatures was not significantly different (Levene's test,  $P > 0.05$ ). A reduction of PCR cycle number by 10 also did not yield any significant effect on qualitative (88.89% of taxa recovered) or quantitative taxon recovery (Supplementary Fig. S9).

### Discussion and potential implications

#### Phylogenetic and taxonomic utility of long rDNA amplicons

Developments in long-amplicon sequencing hold great promise for molecular taxonomy and phylogenetics across very broad taxonomic scales. We recovered phylogenetic relationships across the eukaryote tree of life, which were mostly consistent with the current state of research (e.g., [59]). Separate orders of arthropods all formed well-supported monophyletic groups. Our spider phylogeny was highly congruent with recent work based on whole transcriptomes [35] and multi-amplicon data [60]. Moreover, using the rDNA cluster allowed us to resolve young phylogenetic divergences; the relationships within the recent adaptive radiation of the genus *Tetragnatha* in Hawaii confirmed previous research [65, 67].

In addition to their high phylogenetic utility, long rDNA amplicons showed excellent support for taxonomic hypotheses. All morphologically identified species of Hawaiian *Tetragnatha*



**Figure 7:** Relative abundances for nine arthropod species in our four mock communities (actual) compared to an Illumina amplicon sequencing protocol and Nanopore protocols at 63 °C and 68 °C annealing temperature.

were recovered as monophyletic groups. The divergence patterns and taxonomic classifications of spiders based on rDNA were strongly correlated to those based on mitochondrial COI, the most commonly used animal barcode marker [4]. rDNA may thus be ideal to complement mitochondrial barcoding. A universal and variable nuclear marker as a supplement to COI barcoding will be particularly useful in cases of mito-nuclear discordance due to male biased gene flow [10, 68], hybridization [12], or infections with reproductive parasites [11].

Their high phylogenetic utility across very broad taxonomic categories also provides long rDNA amplicons with a distinct advantage over short read barcoding protocols, which are not well suited to support broad-scale phylogenetic hypotheses [69]. The inclusion of long amplicons would make it possible to scale up barcoding from simple taxon assignment to community-wide phylogenetic inferences [9]. It should be noted that the nuclear rDNA cluster is a single locus, and its divergence pattern does not necessarily reflect species divergence. Also, the multiple genomic rDNA copies do not necessarily all evolve in concert. rDNA genes may even be prone to pseudogenization.

Taxonomic and phylogenetic analyses based on rDNA may thus be affected by paralogues, and additional information from unlinked genomic regions would therefore be highly desirable to support taxonomic and phylogenetic hypotheses. The mitochondrial genome may be an ideal target for this purpose. Recently, the amplification of whole mitochondrial genomes was suggested for animal barcoding [70]. This would increase taxonomic and phylogenetic resolution and alleviate some disadvantages of short COI amplicons. However, it is challenging to develop truly universal primers to target mitochondrial genomes across a very wide range of taxonomic groups [71]. A straightforward way to achieve highly resolved phylogenies may be the combination of long rDNA amplicon sequencing with multiple PCR of short mitochondrial amplicons to amplify multiple mitochondrial DNA fragments [72]. Conserved stretches in mitochondrial rDNA may also allow the design of order- or even phylum-specific primers for long-range amplification [72]. A combination of long mitochondrial and nuclear rDNA ampli-

cons, possibly in a multiplex PCR, would be a desirable development for future DNA barcoding. With whole genome sequences of different taxa rapidly accumulating, it may also be possible to identify additional unlinked DNA barcoding markers.

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

Despite the high raw read error rate 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 potentially be amplified across diverse eukaryote taxa, here, largely demonstrated in arthropods and, in small scale, in fungi and plants. 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 a 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, amount to less than \$4 for each long barcode sequence generated.

#### **Pitfalls of Nanopore-based long-amplicon barcoding**

While our protocol was generally straightforward and reliable, we found several drawbacks that 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 [73]. 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 con-

centrations were loaded. The long-range polymerase may be more sensitive to PCR inhibitors present in some arthropod DNA extractions [74]. PCR conditions will have to be carefully optimized for reliable and consistent amplification. We also found that highly universal eukaryote primers may result in undesired amplification, e.g., plants from beetle and butterfly larval guts, phoretic mites, or fungal sequences. However, as long as the DNA of the target taxon is still dominating the resulting amplicon mixture, this undesired amplification will not affect consensus calling. It may be advisable to check the taxonomic composition of amplicon samples before assembly, e.g., by blasting against a reference library. To reduce non-target amplification, PCR primers could also be redesigned to exclude certain lineages from amplification.

It should also be noted that our approach results in only a single consensus sequence for each processed specimen. As a diploid marker, the rDNA cluster can contain heterozygous positions in some specimens, in particular within the ITS regions. This information is currently lost, and a different assembly approach may be necessary to recover heterozygosity as well. Furthermore, index length and edit distance are also important considerations. We used indexes of 15 bp and with a minimum distance of 10 bp to index both sides of our amplicons. Index edit distance of only 4 bp between samples already led to considerable cross-specimen index bleeding. It may thus be better to increase the length and edit distances of indexes. For example, indexes of 20 or 30 bp could be easily attached to the 5'-tails of PCR primers without strongly affecting PCR efficiency. We have used a relatively crude gel-based approach for pooling amplicon samples. This could have contributed to biased read abundance between some samples. Instead of gel electrophoresis, it may be advisable to use a more precise spectrophotometric quantification.

#### Nanopore-based arthropod metabarcoding

It is well known that Illumina amplicon sequencing of short 18S rDNA fragments can yield accurate taxon recovery in metabarcoding experiments [34], a finding that is confirmed by our results. Except for some outliers (e.g., *Diptera* were overrepresented), even the approximate relative abundance of all taxa was recovered. In contrast, little is known on the performance of long-amplicon Nanopore sequencing for community diversity assessments [32]. Our long barcode-based approach resulted in the dropout of several taxa and highly skewed relative taxon abundances. Skewed abundances were already found in microbial community analysis using Nanopore [32]. In the simplest case, primer mismatches may be responsible for biased amplification [32, 75]. However, the targeted priming sites in our study were extremely conserved. Also, a change of PCR cycle number and annealing temperature did not have a strong effect on taxon abundances, as would be expected in the case of PCR priming bias [76]. Another possibility is the preferential amplification of template molecules with a certain GC content by the DNA polymerase [33]. However, we found the GC content of the rDNA cluster to be very stable across taxa. Yet another potential explanation for the differential recovery of taxa in community samples is taxonomic bias in DNA degradation [77], but we do not expect DNA degradation to have played a role in our experiment because we used only high-quality DNA extractions (verified by gel electrophoresis) from fresh specimens. The most plausible explanation appears to be that variable rDNA lengths are found between different taxa. It is well known that shorter sequences are amplified preferentially in a PCR, especially after it reaches the plateau stage [78]. Such dominance of shorter amplicons could

explain the observed biases very well. In fact, the most abundant taxon in our pools was a spider, which also had the shortest amplicon length. The dominant amplification of shorter sequences may also explain the amplification of plant DNA from a butterfly and a flour beetle larva, as plants showed considerably shorter rDNA amplicons than insects. We found a very high variation of rDNA amplicon length within many taxonomic groups, which could be a considerable problem for long-read metabarcoding applications. This suggests that it may be worthwhile to focus on narrower taxonomic groups for long-amplicon metabarcoding. For example, all spiders in our study share rDNA amplicons of very similar size and would probably be less affected by amplification bias. However, with more closely related taxa in a community, the high error rate of raw reads may cause problems during read clustering and taxon assignments. It should also be noted that we used highly simplified mock community samples, not reflecting actual community composition in nature. Even with those simplified communities, we encountered considerable problems in taxon recovery. Metabarcoding with MinION sequencing may thus be much less trivial than single specimen sequencing. More research into the causes and possible mitigation of these biases will be required before long-amplicon sequencing can be routinely utilized for metabarcoding applications.

## Conclusion

Sequencing long dual indexed rDNA amplicons on ONT's MinION is a simple, cost-effective, accurate, and universal approach for eukaryote DNA barcoding. Long rDNA amplicons offer high phylogenetic and taxonomic resolution across broad taxonomic scales from kingdom down to species. They also prove to be an excellent complement to mitochondrial COI-based barcoding in arthropods. However, despite the long-amplicon advantages in the analysis of separate specimens, we found considerable biases associated with sequencing bulk community samples. The observed taxonomic bias is possibly a result of taxon-specific length variation of the rDNA cluster and preferential amplification of species with shorter rDNA. Further research into the sources of the observed bias is required before long rDNA amplicon sequencing can be utilized as a reliable resource for the analysis of bulk samples.

## Availability of source code and requirements

The program Minibar can be found at <https://github.com/calacademy-research/minibar>  
 Programming language: Python 2.7 (but can be run in Python 3)  
 Operating systems: MacOS, Linux and Windows  
 Other requirements: Edlib library module (<https://github.com/Martinosos/edlib>)  
 License: BSD 2-clause

## Availability of supporting data

The following data supporting the results of this article are available in the GigaScience repository [79]:

- Raw fastq read files from Nanopore sequencing runs and Illumina sequencing of arthropod mock communities for short 18S amplicons
- Fasta sequences of rDNA amplicon for all taxa, mitochondrial COI for Hawaiian *Tetragnatha* spp., as well as Illumina-derived consensus sequences for Hawaiian *Peperomia* spp.
- Newick tree files

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

## Additional file

SupplementaryTable1.SampleList.xlsx

## Abbreviations

BLAST: Basic Local Alignment Search Tool; COI: cytochrome oxidase subunit I; ITS: internal transcribed spacer; NGS: next-generation sequencing; ONT: Oxford Nanopore Technologies; PacBio: Pacific Bioscience; PCR: polymerase chain reaction; rDNA: ribosomal DNA; RTA: retrolateral tibial apophysis.

## Competing interests

The authors declare that they have no competing interests.

## Author contributions

H.K. and S.P. designed the study. H.K., A.P., S.R.K., J.Y.L., N.G., V.S., and J.D.S. collected the specimens. Laboratory work was carried out by H.K., A.P., and S.R.K., and the data were subsequently analyzed by H.K., A.P., J.B.H., S.R.K., and S.P. The article was written by H.K., A.P., J.B.H., S.R.K., J.Y.L., V.S., J.D.S., N.G., N.H.P., R.G.G., and S.P.

## Acknowledgements

We thank Taylor Liu for help during laboratory work, and Tara Gallant for help during specimen collection. Hitomi Asahara graciously provided access to a laboratory facility and the necessary software for our MinION sequencing run. We thank the State of Hawaii Department of Land and Natural Resources and the Servicio Nacional Forestal y de Fauna Silvestre, who provided collection permits and rainforest expeditions, and Gabriela Orihuela for providing assistance and support with fieldwork in Peru. We thank Anna Holmquist for providing the *Psechrus* sp. specimen. The specimen was collected with permits from the Indonesian Ministry of Research and Technology (KEMENRISTEK) and in collaboration with Pungki Lupiyaningdyah and Anang Achmadi from the Museum Zoologicum Bogoriense and funded by a National Science Foundation grant DEB 1457845.

## References

1. Sala OE, Chapin FS, Armesto JJ, et al. Global biodiversity scenarios for the year 2100. *Science* 2000;**287**:1770–4.
2. Pimm SL, Jenkins CN, Abell R, et al. The biodiversity of species and their rates of extinction, distribution, and protection. *Science* 2014;**344**:1246–752.
3. Rominger A, Goodman K, Lim J, et al. Community assembly on isolated islands: macroecology meets evolution. *Global Ecol Biogeogr* 2016;**25**:769–80.
4. Hebert PD, Ratnasingham S, de Waard JR. Barcoding animal life: cytochrome c oxidase subunit 1 divergences among closely related species. *Proc Biol Sci* 2003;**270**:S96–9.
5. Schoch CL, Seifert KA, Huhndorf S, et al. Nuclear ribosomal internal transcribed spacer (ITS) region as a universal DNA

- barcode marker for Fungi. *Proc Natl Acad Sci* 2012;**109**:6241–6.
6. Li D-Z, Gao L-MChina Plant BOL Group; et al.; China Plant BOL Group Comparative analysis of a large dataset indicates that internal transcribed spacer (ITS) should be incorporated into the core barcode for seed plants. *Proc Natl Acad Sci* 2011;**108**:19641–6.
7. Shokralla S, Porter TM, Gibson JF, et al. Massively parallel multiplex DNA sequencing for specimen identification using an Illumina MiSeq platform. *Sci Rep* 2015;**5**:9687.
8. Yu DW, Ji Y, Emerson BC, et al. Biodiversity soup: metabarcoding of arthropods for rapid biodiversity assessment and biomonitoring. *Methods Ecol Evol* 2012;**3**:613–23.
9. Graham CH, Fine PV. Phylogenetic beta diversity: linking ecological and evolutionary processes across space in time. *Ecol Lett* 2008;**11**:1265–77.
10. Krehenwinkel H, Graze M, Rödder D, et al. A phylogeographical survey of a highly dispersive spider reveals eastern Asia as a major glacial refugium for Palaearctic fauna. *J Biogeogr* 2016;**43**:1583–94.
11. Hurst GD, Jiggins FM. Problems with mitochondrial DNA as a marker in population, phylogeographic and phylogenetic studies: the effects of inherited symbionts. *Proc Biol Sci* 2005;**272**:1525–34.
12. Bernatchez L, Glémet H, Wilson CC, et al. Introgression and fixation of Arctic char (*Salvelinus alpinus*) mitochondrial genome in an allopatric population of brook trout (*Salvelinus fontinalis*). *Can J Fish Aquat Sci* 1995;**52**:179–85.
13. Melo-Ferreira J, Boursot P, Suchentrunk F, et al. Invasion from the cold past: extensive introgression of mountain hare (*Lepus timidus*) mitochondrial DNA into three other hare species in northern Iberia. *Mol Ecol* 2005;**14**:2459–64.
14. Soltis PS, Soltis DE. Molecular evolution of 18S rDNA in angiosperms: implications for character weighting in phylogenetic analysis. In: *Molecular systematics of plants II*. Springer; 1998. pp. 188–210.
15. Hillis DM, Dixon MT. Ribosomal DNA: molecular evolution and phylogenetic inference. *Q Rev Biol* 1991;**66**:411–53.
16. Black WC, IV, Klompen J, Keirans JE. Phylogenetic relationships among tick subfamilies (Ixodida: Ixodidae: Argasidae) based on the 18S nuclear rDNA gene. *Mol Phylogenet Evol* 1997;**7**:129–44.
17. Powers TO, Todd T, Burnell A, et al. The rDNA internal transcribed spacer region as a taxonomic marker for nematodes. *J Nematol* 1997;**29**:441.
18. Sonnenberg R, Nolte AW, Tautz D. An evaluation of LSU rDNA D1-D2 sequences for their use in species identification. *Front Zool* 2007;**4**:6.
19. Tang CQ, Leasi F, Obertegger U, et al. The widely used small subunit 18S rDNA molecule greatly underestimates true diversity in biodiversity surveys of the meiofauna. *Proc Natl Acad Sci* 2012;**109**:16208–12.
20. von der Schulenburg JHG, Hancock JM, Pagnamenta A, et al. Extreme length and length variation in the first ribosomal internal transcribed spacer of ladybird beetles (Coleoptera: Coccinellidae). *Mol Biol Evol* 2001;**18**:648–60.
21. Jain M, Koren S, Miga KH, et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat Biotechnol* 2018;**36**:338.
22. Heeger F, Bourne EC, Baschien C, et al. Long-amplicon DNA metabarcoding of ribosomal rRNA in the analysis of fungi from aquatic environments. *Mol Ecol Resour* 2018;**18**:1500–14.
23. Tedersoo L, Tooming-Klunderud A, Anslan S. PacBio

- metabarcoding of Fungi and other eukaryotes: errors, biases and perspectives. *New Phytol* 2018;**217**:1370–85.
24. Giordano F, Aigrain L, Quail MA, et al. De novo yeast genome assemblies from MinION, PacBio and MiSeq platforms. *Sci Rep* 2017;**7**:3935.
  25. Pomerantz A, Peñafiel N, Arteaga A, et al. Real-time DNA barcoding in a rainforest using Nanopore sequencing: opportunities for rapid biodiversity assessments and local capacity building. *GigaScience* 2018;**7**:1–14.
  26. Wurzbacher C, Larsson E, Bengtsson-Palme J, et al. Introducing ribosomal tandem repeat barcoding for fungi. *Mol Ecol Resour* 2019;**19**:118–27.
  27. Srivathsan A, Baloğlu B, Wang W, et al. A MinION™-based pipeline for fast and cost-effective DNA barcoding. *Mol Ecol Resour* 2018;**18**:1035–49.
  28. Quick J, Loman NJ, Duraffour S, et al. Real-time, portable genome sequencing for Ebola surveillance. *Nature* 2016;**530**:228.
  29. Edwards A, Debonnaire AR, Sattler B, et al. Extreme metagenomics using Nanopore DNA sequencing: a field report from Svalbard, 78 N. *bioRxiv* 2016;073965.
  30. Giribet G, Edgecombe GD. Reevaluating the arthropod tree of life. *Annu Rev Entomol* 2012;**57**:167–86.
  31. Hochkirch A. The insect crisis we can't ignore. *Nature News* 2016;**539**:141.
  32. Benítez-Páez A, Portune KJ, Sanz Y. Species-level resolution of 16S rRNA gene amplicons sequenced through the MinION™ portable Nanopore sequencer. *GigaScience* 2016;**5**:4.
  33. Nichols RV, Vollmers C, Newsom LA, et al. Minimizing polymerase biases in metabarcoding. *Mol Ecol Resour* 2018;**18**:927–39.
  34. Krehenwinkel H, Wolf M, Lim JY, et al. Estimating and mitigating amplification bias in qualitative and quantitative arthropod metabarcoding. *Sci Rep* 2017;**7**:17668.
  35. Fernández R, Kallal RJ, Dimitrov D, et al. Phylogenomics, diversification dynamics, and comparative transcriptomics across the spider tree of life. *Curr Biol* 2018;**28**:1489–97. e1485.
  36. NEB Tm Calculator. <https://tcalculator.neb.com/#/main>, Accessed 24 April 2018.
  37. Barcode Generator - Comaiwiki. [http://comailab.genomecenter.ucdavis.edu/index.php/Barcode\\_generator](http://comailab.genomecenter.ucdavis.edu/index.php/Barcode_generator), Accessed 07 February 2018.
  38. De Coster W, D'Hert S, Schultz DT, et al. NanoPack: visualizing and processing long-read sequencing data. *Bioinformatics* 2018;**34**:2666–9.
  39. Filtering and trimming of long read sequencing data. <https://github.com/wdecoster/nanofilt>, Accessed 24 April 2018.
  40. For wrangling HLA alleles. <https://github.com/transplantaton-immunology/allele-wrangler/>, Accessed 27 April 2018.
  41. Vaser R, Sović I, Nagarajan N, et al. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res* 2017;**27**:737–46.
  42. Ultrafast consensus module for raw de novo genome assembly of long uncorrected reads. <https://github.com/isovic/racoon>, Accessed 01 May 2018.
  43. A versatile pairwise aligner for genomic and spliced nucleotide sequences. <https://github.com/lh3/minimap2>, Accessed 25 April 2018.
  44. Tamura K, Stecher G, Peterson D, et al. MEGA6: molecular evolutionary genetics analysis version 6.0. *Mol Biol Evol* 2013;**30**:2725–9.
  45. Straub SC, Parks M, Weitemier K, et al. Navigating the tip of the genomic iceberg: next-generation sequencing for plant systematics. *Am J Bot* 2012;**99**:349–64.
  46. Bolger A, Giorgi F. Trimmomatic: A Flexible Read Trimming Tool for Illumina NGS Data. URL <http://www.usadellab.org/cms/index.php>, Accessed 23 May 2018.
  47. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* 2009;**25**:1754–60.
  48. Li H, Handsaker B, Wysoker A, et al. The sequence alignment/map format and SAMtools. *Bioinformatics* 2009;**25**:2078–9.
  49. Lanfear R, Calcott B, Ho SY, et al. PartitionFinder: combined selection of partitioning schemes and substitution models for phylogenetic analyses. *Mol Biol Evol* 2012;**29**:1695–701.
  50. Huelsenbeck JP, Ronquist F. MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics* 2001;**17**:754–5.
  51. Dray S, Dufour A-B. The ade4 package: implementing the duality diagram for ecologists. *J Stat Softw* 2007;**22**:1–20.
  52. Machida RJ, Knowlton N. PCR primers for metazoan nuclear 18S and 28S ribosomal DNA sequences. *PLoS One* 2012;**7**:e46180.
  53. Krehenwinkel H, Kennedy S, Pekár S, et al. A cost-efficient and simple protocol to enrich prey DNA from extractions of predatory arthropods for large-scale gut content analysis by Illumina sequencing. *Methods Ecol Evol* 2017;**8**:126–34.
  54. Altschul SF, Gish W, Miller W, et al. Basic local alignment search tool. *J Mol Biol* 1990;**215**:403–10.
  55. Levenshtein VI. Binary codes capable of correcting deletions, insertions, and reversals. In: *Soviet physics doklady*, Volume 10; 1966. pp. 707–10.
  56. Šošić M, Šikić M. Edlib: a C/C++ library for fast, exact sequence alignment using edit distance. *Bioinformatics* 2017;**33**:1394–5.
  57. Dual barcode and primer demultiplexing for MinION sequenced reads. <https://github.com/calacademy-research/minibar>, Accessed 01 June 2018.
  58. Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using only Nanopore sequencing data. *Nature Methods* 2015;**12**:733–5.
  59. Misof B, Liu S, Meusemann K, et al. Phylogenomics resolves the timing and pattern of insect evolution. *Science* 2014;**346**:763–7.
  60. Wheeler WC, Coddington JA, Crowley LM, et al. The spider tree of life: phylogeny of Araneae based on target-gene analyses from an extensive taxon sampling. *Cladistics* 2017;**33**:574–616.
  61. Gillespie RG. Hawaiian spiders of the genus Tetragnatha: I. Spiny leg clade. *J Arachnol* 1991;**19**:174–209.
  62. Gillespie RG. Comparison of rates of speciation in web-building and non-web-building groups within a Hawaiian spider radiation. *J Arachnol* 1999;**27**:79–85.
  63. Gillespie RG. Island time and the interplay between ecology and evolution in species diversification. *Evol Appl* 2016;**9**:53–73.
  64. Gillespie RG, Croom HB, Hasty GL. Phylogenetic relationships and adaptive shifts among major clades of Tetragnatha spiders (Araneae: Tetragnathidae) in Hawai'i. *Pac Sci* 1997;**51**:380–94.
  65. Gillespie R. Community assembly through adaptive radiation in Hawaiian spiders. *Science* 2004;**303**:356–9.
  66. Blackledge TA, Binford GJ, Gillespie RG. Resource use within a community of Hawaiian spiders (Araneae: Tetragnathidae). In: *Annales Zoologici Fennici*. JSTOR; 2003. pp. 293–303.
  67. Blackledge TA, Gillespie RG. Convergent evolution of behav-

- ior in an adaptive radiation of Hawaiian web-building spiders. *Proc Natl Acad Sci USA* 2004;**101**:16228–33.
68. Wilmer JW, Hall L, Barratt E, et al. Genetic structure and male-mediated gene flow in the ghost bat (*Macroderma gigas*). *Evolution* 1999;**53**:1582–91.
69. Kjer KM, Zhou X, Frandsen PB, et al. Moving toward species-level phylogeny using ribosomal DNA and COI barcodes: an example from the diverse caddisfly genus *Chimarra* (Trichoptera: Philopotamidae). *Arthropod Syst Phylogeny* 2014;**72**:345–54.
70. Deiner K, Renshaw MA, Li Y, et al. Long-range PCR allows sequencing of mitochondrial genomes from environmental DNA. *Methods Ecol Evol* 2017;**8**:1888–98.
71. Briscoe AG, Goodacre S, Masta SE, et al. Can long-range PCR be used to amplify genetically divergent mitochondrial genomes for comparative phylogenetics? A case study within spiders (Arthropoda: Araneae). *PLoS One* 2013;**8**:e62404.
72. Krehenwinkel H, Kennedy SR, Rueda A, et al. Scaling up DNA barcoding—Primer sets for simple and cost-efficient arthropod systematics by multiplex PCR and Illumina amplicon sequencing. *Methods Ecol Evol* 2018;**9**(11):2181–93.
73. Krehenwinkel H, Pekar S. 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* 2015;**10**:e0136337.
74. Margam VM, Gachomo EW, Shukle JH, et al. A simplified arthropod genomic-DNA extraction protocol for polymerase chain reaction (PCR)-based specimen identification through barcoding. *Mol Biol Rep* 2010;**37**:3631–5.
75. Sipos R, Székely AJ, Palatinszky M, et al. Effect of primer mismatch, annealing temperature and PCR cycle number on 16S rRNA gene-targeting bacterial community analysis. *FEMS Microbiol Ecol* 2007;**60**:341–50.
76. Suzuki MT, Giovannoni SJ. Bias caused by template annealing in the amplification of mixtures of 16S rRNA genes by PCR. *Appl Environ Microbiol* 1996;**62**:625–30.
77. Krehenwinkel H, Fong M, Kennedy S, et al. The effect of DNA degradation bias in passive sampling devices on metabarcoding studies of arthropod communities and their associated microbiota. *PLoS One* 2018;**13**:e0189188.
78. Wattier R, Engel C, Saumitou Laprade P, et al. Short allele dominance as a source of heterozygote deficiency at microsatellite loci: experimental evidence at the dinucleotide locus Gv1CT in *Gracilaria gracilis* (Rhodophyta). *Mol Ecol* 1998;**7**:1569–73.
79. Krehenwinke LH, Pomerantz A, Henderson JB, et al. 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* 2019; <http://dx.doi.org/10.5524/100552>.