

1 **Title: Cost-efficient high throughput capture of museum arthropod specimen DNA using**

2 **PCR-generated baits**

3

4 **Running title: Capture of museum specimens using PCR baits**

5

6 Alexander Knyshev, University of California Riverside, Entomology, Riverside, CA, USA,

7 corresponding author email and ORCID: aknys001@ucr.edu, orcid.org/0000-0002-2141-9447

8

9 Eric R.L. Gordon¹, University of California Riverside, Entomology, Riverside, CA, USA,

10

11 Christiane Weirauch, University of California Riverside, Entomology, Riverside, CA, USA

12

¹ Current affiliation: University of Connecticut, Ecology and Evolutionary Biology, Storrs, CT, USA

13 **Abstract:**

- 14 1. Gathering genetic data for rare species is one of the biggest remaining obstacles in
15 modern phylogenetics, particularly for megadiverse groups such as arthropods. Next
16 generation sequencing techniques allow for sequencing of short DNA fragments
17 contained in preserved specimens >20 years old, but approaches such as whole genome
18 sequencing are often too expensive for projects including many taxa. Several methods of
19 reduced representation sequencing have been proposed that lower the cost of sequencing
20 per specimen, but many remain costly because they involve synthesizing nucleotide
21 probes and target hundreds of loci. These datasets are also frequently unique for each
22 project and thus generally incompatible with other similar datasets.
- 23 2. Here, we explore utilization of in-house generated DNA baits to capture commonly
24 utilized mitochondrial and ribosomal DNA loci from insect museum specimens of various
25 age and preservation types without the a priori need to know the sequence of the target
26 loci. Both within species and cross-species capture are explored, on preserved specimens
27 ranging in age from one to 54 years old.
- 28 3. We found most samples produced sufficient amounts of data to assemble the nuclear
29 ribosomal rRNA genes and near complete mitochondrial genomes and produce well-
30 resolved phylogenies in line with expected results. The dataset obtained can be
31 straightforwardly combined with the large cache of existing Sanger-sequencing-generated
32 data built up over the past 30 years and targeted loci can be easily modified to those
33 commonly used in different taxa. Furthermore, the protocol we describe allows for
34 inexpensive data generation (as low as ~\$35/sample), of at least 20 kilobases per

35 specimen, for specimens at least as old as ~1965, and can be easily conducted in most
36 laboratories.

37 4. If widely applied, this technique will accelerate the accurate resolution of the Tree of Life
38 especially on non-model organisms with limited existing genomic resources.

39 **Keywords: [4-6]**

40 Insects, phylogeny, host plant, Miridae

41 **Introduction:**

42 Natural history museums host troves of biological material and sometimes the only known
43 representatives of extinct or rare species (Coddington, Agnarsson, Miller, Kuntner, & Hormiga,
44 2009; Lim, Balke, & Meier, 2011). In these cases, museum specimens represent the only
45 accessible sources of genetic data for a given species and gathering data from such specimens in
46 a cost-effective way is one of the primary obstacles yet to be overcome in modern phylogenetics.
47 Specimens in museums may also allow for the inclusion of a temporal variable into analyses by
48 comparing DNA sequence of individuals across different sampling dates and can even be used
49 for the analysis of short-term evolutionary trends (Hartley et al., 2006; DiEuliis, Johnson, Morse,
50 & Schindel, 2016).

51 Preservation conditions of museum material can dramatically impact the viability of obtaining
52 DNA sequence data. Traditional approaches used amplification of target regions of DNA
53 followed by Sanger sequencing. This method is highly dependent on residual DNA fragment size
54 and the proportion of endogenous DNA in the extract. While targeting shorter gene regions can
55 mitigate the issue of DNA fragmentation, a low endogenous content is harder to overcome

56 (Burrell, Disotell, & Bergey, 2015), and even innovative new PCR techniques are only capable of
57 somewhat reliably amplifying fragments of less than 600 bp (Mitchell, 2015). The development
58 of next generation sequencing (NGS) has expanded the array of methods for DNA sequencing
59 from museum specimens. For whole genome sequencing, an NGS library is prepared from the
60 original DNA extract and this library is then combined with other samples for multiplex
61 sequencing and allocated a certain proportion of reads on a sequencer, depending on desired
62 sequencing depth and the total budget (Cridland, Ramirez, Dean, Sciligo, & Tsutsui, 2018;
63 Kanda, Pflug, Sproul, Dasenko, & Maddison, 2015; Maddison & Cooper, 2014). However, even
64 low-coverage whole genome sequencing is currently still prohibitively expensive for all but very
65 well-funded projects or studies focusing on relatively few samples.

66 As a way to decrease the cost per sample while still generating sufficient amounts of data for
67 accurate phylogenetic placement, several methods of reduced representation sequencing have
68 been proposed. Typically, these methods include selective hybrid capture of target loci, where the
69 type and number of loci being captured depends on the scope and context of the study. During the
70 past few years, utilization of commercially synthesized probes or microarray kits for the capture
71 of conserved DNA regions has become popular for phylogenetic studies and can be applied to
72 historical museum specimens (Bi et al., 2013; Blaimer, Lloyd, Guillory, & Brady, 2016;
73 McCormack, Tsai, & Faircloth, 2015). These kits are designed based on existing reference
74 genomes or transcriptomes and typically enrich many loci (~500-5000), thus a large amount of
75 data is generated for each sample, but the cost per sample is relatively high. Other kits are
76 designed to enrich mitochondrial genomes, including a kit specifically designed for
77 mitochondrial DNA across insects (Liu et al., 2016). However, all methods relying on

78 commercially synthesized kits are relatively expensive and might not be feasible for low-budget
79 projects. These kits are also limited by the original design and probe composition cannot be
80 adjusted after synthesis.

81 These limitations led us to explore an approach that uses in-house generated DNA baits for
82 hybrid enrichment (Maricic, Whitten, & Pääbo, 2010). These baits can be produced from
83 amplicons generated by PCR of short gene regions (Peñalba et al., 2014), or by long-range PCR
84 of complete mitochondrial (Li et al., 2015; Maricic et al., 2010) or chloroplast genomes (Mariac
85 et al., 2014), or even from ddRAD library fragments (Suchan et al., 2016). PCR-generated baits
86 have so far only been applied to vertebrates and plants, and only in a few cases tested on archival
87 specimens (Li et al., 2015). Appealing features of this approach include affordable synthesis of
88 baits, independence from the need of a good quality reference, and flexibility of the synthesis
89 workflow for low-cost modifications of the bait set (e.g., pooling different combinations of bait
90 amplicons, using same primers to obtain bait amplicons from different taxa, or generating
91 additional baits with new sets of primers).

92 The diversity of arthropods is staggering, with estimates of about 80% of species still
93 undescribed (Stork, 2018) and an increasing number of species going extinct every day
94 (Hallmann et al., 2017). While modern phylogenetic studies of vertebrates sometimes approach
95 complete sampling of extant diversity, complete extant sampling of any large clade of arthropods
96 is almost impossible due to the abundance of rare species, limited material, and the huge diversity
97 of arthropods (Coddington et al., 2009; Lim et al., 2011). However, near-complete sampling is
98 useful for many downstream analyses, including unbiased estimation of lineage diversification
99 rates (Cusimano, & Renner, 2010; Cusimano, Stadler, & Renner, 2012; Höhna, Stadler, Ronquist,

100 & Britton, 2011). Scientists have just started to utilize the enormous resources of arthropod
101 specimens deposited in natural history collections for gathering large DNA datasets (Stork,
102 McBroom, Gely, & Hamilton, 2015; Stork, 2018). We argue that insects in particular are an apt
103 test case for the application of new NGS approaches to illuminating the dark areas in the Tree of
104 Life, because most material in entomological collections is stored as dried and pinned or point-
105 mounted specimens, which are often suitable for the retrieval of fragmented DNA. Previous
106 applications of this approach on vertebrate and plant samples employed destructive extraction
107 protocols to generate adequate amounts of DNA for capture. But DNA extraction can be
108 performed without destroying external or genitalic morphological features and from individual
109 and small specimens as in many insects. For a complete taxonomic sampling of large clades,
110 already existing data should be compatible with character-rich new datasets generated at low
111 costs.

112
113 Here, we test the efficiency of PCR-generated DNA baits (targeting the mitochondrial genome,
114 nuclear ribosomal operon, and one nuclear protein-coding gene) to capture DNA sequences from
115 museum-deposited insect specimens with different collection dates, preservation methods, and
116 evolutionary relatedness, using phyline plant bugs (Insecta: Hemiptera: Miridae: Phylinae) as our
117 test case. These loci were selected for optimal integration with existing, Sanger-based sequence
118 data and to allow adequate coverage when multiplexing hundreds of libraries. Plant bugs are a
119 group of > 11,000 described species that include serious plant pests and beneficial insects (Cassis
120 & Schuh, 2012). Phylogenetic hypotheses for the entire group are in their infancy (Jung & Lee,
121 2011), but studies targeting selected subfamilies including the Phylinae now provide testable

122 hypotheses (Konstantinov & Knyshov, 2015; Menard, Schuh, & Woolley, 2014; Namyatova,
123 Konstantinov, & Cassis, 2015; Tataric & Cassis, 2012). The taxonomic diversity of plant bugs
124 in the Western U.S. is fairly well understood (Cassis & Schuh, 2012; Weirauch et al., 2016), but
125 few species have been incorporated into phylogenetic analyses, and some are only known from
126 the type specimen(s). As the first test case, we selected a putatively monophyletic group of native
127 oak-associated plant bugs, the so called “Orange Oak Bugs” (OOB) (Weirauch, 2006a, 2006b),
128 where some species may be monophagous on specific species of oaks, while at least two
129 widespread and polymorphic species (*Phallospinophylus setosus* Weirauch and *Pygovepres*
130 *vaccinicola* (Knight)) feed on a variety of host plants (including Fagaceae, Rhamnaceae, and
131 Rosaceae). We sampled specimens of these two species from a range of localities and host plants,
132 together with several additional species of OOB that had not yet been included in phylogenetic
133 analyses (Menard et al., 2014) to test efficacy of capture across closely related samples and to
134 investigate potential cryptic host plant races. As second test case, we selected the genus *Tuxedo*
135 Schuh with seven described species associated with host plants in several families (Schuh, 2004);
136 phylogenetic relationships within this genus are unknown. We aimed to sample several
137 individuals from each of the seven species, including paratype specimens, to investigate capture
138 efficiency at deeper phylogenetic levels and to explore host plant shifts within the genus. Both
139 datasets were analyzed together with a Sanger-derived phylogenetic dataset of Phylinae (Menard
140 et al., 2014), demonstrating the feasibility of combining existing and newly generated NGS data.

141 **Material and methods:**

142 *Taxon Sampling and Vouchering*

143 Specimens for this study were loaned from the American Museum of Natural History (AMNH),
144 the Entomology Research Museum (UCRC), and the Zoological Institute, Russian Academy of
145 Sciences (ZISP). Tentative voucher identification was done based on habitus and host association
146 data using Weirauch (2006a, 2006b) and Schuh (2004). Age of specimens at the moment of DNA
147 extraction varied from one to 54 years. Specimens of *Tuxedo*, *Leucophoroptera* Poppius,
148 *Ausejanus* Menard and Schuh, and *Pseudophylus* Yasunaga were imaged using a Leica DFC 450
149 C imaging system. Image vouchers and specimen information are available through the
150 Heteroptera Species Pages (<http://research.amnh.org/pbi/heteropteraspeciespage/>). After clearing
151 soft abdominal tissues during the DNA extraction process, we examined male genitalic characters
152 to confirm our tentative identifications. In cases where different diagnostic characters were in
153 conflict (e.g., in some *Tuxedo* spp., see results and discussion), we based our identification on
154 genitalic characters.

155 *DNA Extraction*

156 In most cases, only the abdomen (1-1.5 mm in length) was used for non-destructive DNA
157 extraction, which was performed using a Qiagen DNeasy® Blood & Tissue kit (for relatively
158 fresh ethanol specimens) or a combination of the previous kit with a Qiagen QIAquick® PCR
159 purification kit (for dry specimens, see supplemental text S1), since the latter is commonly used
160 for DNA extraction from degraded samples (Lee et al., 2010; Yang, Eng, Waye, Dudar, &
161 Saunders, 1998). Abdomens were soaked in the extraction buffer, such that cuticular structures
162 remain undamaged, and mirrored standard dissection procedures for plant bug specimens. This
163 approach allows for subsequent remounting of the abdominal cuticle and genitalia with the rest of
164 the specimen, or in a genitalic vial.

165 *Bait Synthesis*

166 Freshly collected specimens of *Phallospinophylus setosus* and *Tuxedo drakei* Schuh were
167 selected as bait donors for the OOB and *Tuxedo* subprojects, respectively. Primers for obtaining
168 long range PCR products are listed in Table S1. Details on the primer design are available in
169 supplemental text S1. Target regions included mitochondrion, nuclear ribosomal operon, and a
170 fragment of the cytoplasmic dynein heavy chain gene.

171 To prepare baits, six long-range (LR) PCRs per specimen were performed. For this and all
172 subsequent PCR described in this paper, we used Takara PrimeSTAR® GXL polymerase, a hot-
173 start high-fidelity enzyme that is able to amplify long products. The PCR mix contained 10 µl
174 PrimeSTAR® GXL buffer, 4 µl 2.5M dNTPs, 1 µl PrimeSTAR® GXL polymerase, 32 µl water,
175 1.5 µl of each primer (10 µM), and 1 µl of DNA template. The thermocycler program included
176 initial denaturation at 98° for 3 min, 35 cycles of denaturation for 10 sec at 98°, followed by
177 annealing at variable temperatures for 15 sec, followed by elongation at 68° for a variable amount
178 time, and with the final incubation at 68° for 15 min. Additional details on long-range PCR
179 conditions are available in Table S1.

180 After clean up with custom Solid Phase Reversible Immobilization (SPRI) beads (Glenn et al.,
181 2016; Rohland, & Reich, 2012), mitochondrial, nuclear ribosomal, and nuclear protein-coding
182 products were mixed in molar ratios of 1:1:5, following recommendations of Peñalba et al.
183 (2014) regarding capture of low copy nuclear genes. Mixtures were diluted to the volume of 100
184 µl and sonicated on a Diagenode Bioruptor® UCD-200 with 30/30 cycles for 6 runs of 5 minutes.
185 Sheared PCR products were subjected to a bait library preparation generally following the
186 protocol of Li et al. (2015) with the exception that regular dNTPs instead of a dUTP-containing

187 mixture were used, since NaOH melting was used to subsequently elute captured libraries instead
188 of off-bead amplification. Three pools of ready-to-use bait were produced by amplifying M13-
189 adaptor-ligated bait libraries with 5' biotinylated primers using PCR conditions outlined above
190 with the following modifications: 6 µl of template was used, and annealing temperature set to
191 55°.

192 *Preparation of Illumina-compatible Libraries*

193 Since DNA sequence of bait donors was also of interest in this project, we also sequenced
194 amplicons used for bait production. These LR PCR products were mixed in equimolar ratios and
195 sonicated as described above. Following sonication, Illumina®-compatible libraries were
196 prepared using the protocol from Li, Hofreiter, Straube, Corrigan, and Naylor (2013), with the
197 following modifications: end prep mix contained 50% 2X Takara EmeraldAmp® GT PCR mix
198 and after incubation at 25° for 15 min and 12° for 5 min was incubated at 72° for 20 min in order
199 to obtain a-tailed fragments. We utilized with-bead SPRI method as originally described in Fisher
200 et al. (2011), carrying same SPRI beads through the library preparation steps. T-tailed loop
201 adaptors from NEBNext® Multiplex Oligos for Illumina® kit E7600s were ligated to the DNA
202 and a PCR with indexing primers from the same kit was conducted using PCR conditions
203 outlined above with the following modifications: 6 µl of template was used, annealing
204 temperature set to 60°, number of cycles set to 16.

205 To prepare target libraries, DNA extracts were run on a gel with Biotium GelRed® premixed
206 loading buffer in ratios 1:2 to check average fragment size and determine if sonication was
207 needed (i.e., for younger samples). These DNA extracts were quantified using Qubit™
208 fluorometer, and for more consistent sonication results approximately 70 ng of DNA (where

209 possible, also see Table 1) were used for sonication. Library preparation followed the protocol
210 outlined above with the exception that after adaptor ligation, libraries were amplified with short
211 IS7/IS8 primers following Li et al. (2013). The same PCR conditions as above were used,
212 however number of cycles were varied from 16 to 21 depending on the amount of starting
213 material.

214 *First Sequencing Run – Target Capture, Pooling and Sequencing*

215 In our first sequencing run, target captures generally followed the protocol of Li et al. (2015).
216 Every sample was captured individually as in Li et al. (2015), 10 µl of Invitrogen Dynabeads®
217 M-270 and 10 ng of bait library was used for most samples, whereas all remaining bait library
218 was used for the last few captured samples (for details on bait amount used, see Table 1). DNA
219 concentration of input target library was not quantified, and we used 6 µl of target library in each
220 capture reaction. Elution was conducted with NaOH melting as in Maricic et al. (2010), and
221 double capture was performed following suggestions of Peñalba et al. (2014). After the second
222 round of capture, the supernatant was cleaned, and eluted in 50 µl of 10 mM Tris-HCl. Post-
223 capture PCR followed the same PCR procedure as outlined above, however indexing primers and
224 20 µl of template were used, and variable number of cycles was performed (16-24).

225 After indexing PCR, products were cleaned and normalized with Just-a-Plate™ 96 PCR
226 Purification and Normalization Kit. Since using Bioanalyzer on all 60 samples was prohibitively
227 expensive, libraries were first run on a gel with GelRed® to check average fragment size, pooled
228 together into nine groups according to their size, which were then analyzed on a single
229 Bioanalyzer chip to obtain more accurate fragment size distribution. Then libraries were pooled
230 equimolarly with the exception of sheared amplicon libraries (samples ph32 and ph47), which

231 were pooled at twice higher concentrations. The pool of 60 indexed libraries then was mixed in
232 molar proportion of 50:50 with unrelated samples from other projects and sequenced on a single
233 run of Illumina® MiSeq® V3 2x300bp at the UCR IIGB Core Facility.

234 *Second Sequencing Run – Library Preparation, Target Capture and Sequencing*

235 In the second sequencing run, we followed the protocol of Maricic et al. (2010) with
236 modifications. DNA extracts from the same specimen of *Tuxedo drakei* as above was used as a
237 source for bait preparation. The procedure differed from described above in that only nuclear
238 rRNA operon and mitochondrial PCR products were used. We extracted one more specimen of
239 *Pseudophylus* and prepared a library as outlined above. Five libraries (samples ph45, ph54, ph57,
240 ph59, and a new *Pseudophylus* library) were carried through indexing PCR, quantified using
241 Qubit™, checked on an agarose gel, and pooled equimolarly to obtain about 450 ng of DNA.
242 Because indexed libraries were used, we added additional blocking oligos as in Maricic et al.
243 (2010) to block longer adaptor fragments. Approximately 500 ng of bait and 5 µl of Dynabeads®
244 as in Maricic et al. (2010) were used for each round of capture (two rounds total as in the first
245 sequencing run). Post-capture amplification was done using IS5 and IS6 primers and was carried
246 over in two aliquots. After PCR, the products were combined and purified, they were then
247 sequenced on 5% of another Illumina® MiSeq® V3 2x300bp run at the UCR IIGB Core Facility.

248 *Post-Sequencing Data Processing*

249 Raw sequences were demultiplexed and adaptors were removed using bcl2fastq software
250 (Illumina®) at the UCR IIGB Core Facility. Trimmomatic v0.36 (Bolger, Lohse, & Usadel,
251 2014) was used to trim off low quality ends of the sequences as well as perform more thorough

252 adaptor trimming. Reads were assembled into contigs with SPAdes (Bankevich et al., 2012). In
253 cases where assembly did not yield complete target regions, we obtained them by mapping
254 shorter contigs onto full length assemblies of other related samples. Assembled contigs were
255 checked for misassembled regions and manually curated in Geneious v.10
256 (<https://www.geneious.com>, Kearse et al., 2012). We mapped reads on these contigs using BWA
257 (Li & Durbin, 2009) to assess the coverage depth (see Table 1), prior to average coverage
258 calculations, reads were deduplicated using PRINSEQ (Schmieder & Edwards, 2011).

259 We aligned all resulting 18S, 28S, and mitochondrial contigs using MAFFT v.7 (Kato &
260 Standley, 2013). Manual inspection of alignments and trimming was performed. Since accurate
261 assembly of the mitochondrial control region with short reads without a close reference was
262 problematic due to presence of repeats, we excluded it from the analysis. The remainder of the
263 mitochondrion was annotated by aligning it with mitochondrial genome of another plant bug
264 available on GenBank (NC_024641.1).

265 *Phylogenetic Analysis*

266 For phylogenetic analysis, the dataset was concatenated and divided into 18 partitions with
267 protein coding genes split further into codon positions. Substitution models and partitioning
268 scheme were optimized using PartitionFinder 2.1.1 (Lanfear, Frandsen, Wright, Senfeld, &
269 Calcott, 2016) or ModelFinder (Kalyaanamoorthy, Minh, Wong, von Haeseler, & Jermiin, 2017),
270 which is an IQ-TREE built-in model and partition test. Phylogeny estimation was performed in
271 RAxML v8.2.11 (Stamatakis, 2014) and IQ-TREE v1.5.4 (Nguyen, Schmidt, von Haeseler, &
272 Minh, 2014). Branch support was calculated using Rapid Bootstrap (Stamatakis, Hoover, &
273 Rougemont, 2008) which is shown on Fig. 2, Ultrafast Bootstrap (Minh, Nguyen, & von

274 Haeseler, 2013), and SH-aLRT (Anisimova & Gascuel, 2006), which are shown on Figs S2 and
275 S3.

276 To test how well our data can be combined with previously generated data, we combined our data
277 with the dataset of Menard et al. (2014) which is the most comprehensive set of genetic data for
278 related species. After downloading the sequences from GenBank, we extracted only 18S, 28S,
279 16S and COX1 sequences from our data, performed alignment and manual trimming. Alignments
280 were then concatenated, optimized for model and partitioning scheme, and phylogenetically
281 analyzed as above.

282 Illustrations for Figs 2, 3, and S1 were drafted using R v3.4.3 and packages APE (Paradis,
283 Claude, & Strimmer, 2004), phytools (Revell, 2012), ggtree (Yu, Smith, Zhu, Guan, & Lam,
284 2016) and ggbio (Yin, Cook, & Lawrence, 2012). Relief image for Fig. 3 was taken from the
285 SimpleMappr website (<http://www.simplemappr.net/>).

286 **Results and Discussion**

287 *Expenses*

288 Total expenses after bait and target library preparations, target capture, and sequencing and
289 including all reagents and supplies came to about \$54 per specimen or about \$2.8 per 1 Kb of
290 data in the first sequencing run, and about \$39 per specimen or about \$2.1 per 1 Kb of data in the
291 second sequencing run (Table S2). Our estimates suggest that pooled capture together with using
292 a higher throughput sequencer (e.g., a HiSeq® lane or a NextSeq® run) can generate the same
293 amount of data for about half the price (up to \$25 per specimen), however a greater number of
294 samples (at least 360) need to be pooled together to efficiently utilize the sequencer.

295 *DNA extraction*

296 The amount of DNA extracted greatly varied across samples (see Table 1). The minimum amount
297 of DNA that was used for library preparation was 2.75 ng (sample 42). The average fragment size
298 for ethanol preserved material was large: we always detected a bright band larger than 10 Kb in
299 size, with many extracts also with a smear of fragments spanning down to 300 bp. For dry point-
300 mounted material we observed two types of fragmentation: extracts that had fragments of 500-
301 700 bp on average in addition to long (~8Kb) fragments (dry specimens collected within past ten
302 years), and extracts with only fragments shorter than 1000 bp (dry specimens collected more than
303 ten years ago).

304 *Sequencing and assembly*

305 A total of 45% of an Illumina® MiSeq® V3 lane was used for the samples in the first sequencing
306 run. The amount of reads obtained per sample is listed in Table 1 (average of 152670, $\sigma =$
307 39070). For bait samples, we obtained full bait contigs for the nuclear ribosomal operon, the
308 dynein fragment, and entire mitochondrial genome, although unambiguous assembly of the
309 control region was problematic due to lack of a close reference and long read data. For other
310 samples, we obtained full or partial mitochondrial contigs and nuclear ribosomal gene contigs for
311 the majority of samples (see Table 1). Mitochondrial completeness is indicated on Fig. 2B and
312 excludes the control region, and mitochondrial average coverage depth is indicated in Fig. 2C.
313 We were able to obtain reliable ribosomal data for 48 taxa, but some sequences exhibited cross
314 contamination of about 1% of reads by the bait taxon as detailed in supplemental text S1. In the
315 second sequencing run, we observed a higher percent of ribosomal operon reads on target (on
316 average 7.33% in the second run compared to 2.1% in the first run for the same samples) and

317 both for recaptured libraries, as well as for the library prepared after the first sequencing run was
318 complete (ph61), we have not detected contaminating reads.

319 *Capture efficiency*

320 Percent of reads on target varied from 0.61% to 33.95% and was on average 8.19% in the OOB
321 subproject and 4.02% in *Tuxedo* subproject. The percent of reads on target was slightly larger for
322 samples that are close to bait specimens (Figs 2A, 2E, Fig. S4). We also observed a significant
323 variation of percentage on-target across samples of close phylogenetic relatedness, which may be
324 attributed to variation in total amount of target DNA submitted to capture reactions (equal
325 volumes of target libraries were used in all reactions). Capture in the *Tuxedo* subproject
326 performed worse, which could be attributed to the higher sequence divergence from the bait (Fig.
327 2D).

328 Baits for a nuclear protein-coding gene (dynein) performed unsatisfactorily, even though they
329 were five times more concentrated. Although we do not have a clear explanation as to why this
330 bait performed suboptimally (see supplemental text S1), the large middle intron may have been
331 detrimental for bait efficacy.

332 On the contrary, mitochondrial baits were only 14.3% of total bait pool, yet were able to
333 considerably enrich for mitochondrial DNA. Typical sequencing of non-enriched DNA libraries
334 from insect museum specimens yields from 0.002% to 0.08% of total reads mapping to the
335 mitochondrial genome (Staats et al., 2013), however, we recovered on average of 2.13% ($\sigma =$
336 2.59%, range 0.24%-15.64%) representing an enrichment of at least 25x on average for our first
337 sequencing run. Given the amount of reads we allocated for our samples, an unenriched library

338 would produce only about 120 mitochondrial reads, where we achieved on average ~3,500 reads
339 (an enrichment of 29x), sufficient for assembling the whole mitochondrial genome.

340 Suboptimal capture performance in our first sequencing run could be also attributed to the
341 amount of bait. Overall, we observed an increase in the amount of reads on target in capture
342 reactions where more bait was used (Table 1, samples ph29-ph31, ph33-ph35, and ph55-ph60).
343 Thus, we repeated sequencing of five selected samples captured with a modified protocol (see
344 Materials & Methods) where more bait was used. We also explored a pooled capture approach,
345 which is significantly cheaper than the individual sample captures. In the result of the second
346 sequencing run, we observed on average 8.65% on target reads as opposed to 3.42% for the same
347 samples in the first sequencing run (see Table 1). We also noticed a larger variation of total
348 amount of reads received for a given sample in the pool. This might be due to unequal divergence
349 of samples in the pool with the respect to the bait or difference in library quality due to the age of
350 the specimens. Because of this, we recommend balancing sample pools prior to capture and
351 performing individual captures for sensitive samples.

352 Our results show no difference in capture efficiency as related to the age of the specimen (Fig. 2,
353 specimens older than 20 years denoted with red asterisks). We thus expect that even older
354 specimens can be used (Blaimer et al., 2016), but for this pilot study the youngest available
355 specimens were chosen. Further adjustments of hybridization temperature and duration may
356 further improve capture success, however need to be modified on an individual basis.

357 *Phylogenetic analyses*

358 Using the obtained data, we reconstructed a well resolved phylogeny, contributing new insights.
359 Our phylogenetic analysis supports the monophyly of *Tuxedo* + *Pseudophylus*, the OOB clade,
360 *Phallospinophylus setosus* and *Pygovepres vaccinicola* with the highest branch support (Fig. 2A,
361 Fig. S2). As part of the *Tuxedo* subproject, we sampled two specimens of *Pseudophylus stundjuki*
362 (Kulik) since this species from Far East Asia rendered the Western Nearctic *Tuxedo* paraphyletic
363 in a previous analysis (Menard et al., 2014). “*Tuxedo*” is here confirmed to be paraphyletic with
364 respect to *Pseudophylus*, after thorough examination of our sequence data and comparison with
365 data from Menard et al. (2014) and Jung and Lee (2011). All primarily Fagaceae-feeding species
366 of “*Tuxedo*” form a well-supported monophyletic group. Species other than *Tuxedo flavicollis*
367 (Knight) and *Tuxedo susansolomoniae* Schuh were recovered as monophyletic and conform with
368 genitalic-based identifications. Phylogenetic analysis recovered two highly supported
369 monophyletic groups within the *T. flavicollis/susansolomoniae* species group, however
370 composition of each group is not congruent with either genitalic structure or coloration. One
371 specimen (ph57) initially identified as *T. susansolomoniae* is distantly related from other members
372 of *T. flavicollis/susansolomoniae* clade and is recovered as sister taxon to *T. nicholi* (Knight), and
373 likely represents an undescribed species. Species within the OOB clade represented by multiple
374 specimens are monophyletic with high support. Our analysis did not find support for our
375 hypothesis on the presence of host plant races within each of the widespread and polyphagous
376 OOB species (Fig. 3). In contrast, the phylogenetic structure in *Phallospinophylus setosus* is
377 more likely explained by geographic proximity between sampled localities.

378 Combined with existing data of Menard et al. (2014), phylogenetic hypotheses inferred from our
379 dataset are congruent with those presented in prior studies (Fig. 4, Fig. S3). Deep level

380 relationships within Oncotylina as well as the monophyly of the subtribe itself remain poorly
381 supported based on this data set. As in Menard et al. (2014), “*Tuxedo*” + *Pseudophylus* are
382 recovered as the sister group to Leucophoropterini, although with low support. Sampled species
383 of Leucophoropterini were recovered in expected phylogenetic positions.

384 *Conclusions*

385 In conclusion, we were able to cost-efficiently (\$2.8/sample/Kb) sequence long-range PCR
386 products as well as perform hybrid enrichment using in-house generated baits and obtain DNA
387 sequences (~20 Kb) from archival specimens (up to 54 years old) using a minimal amount of
388 DNA. This approach offers a much lower cost of bait production than other approaches, however,
389 especially if LR PCR is chosen for amplicon generation, a high-quality sample of a related
390 species is needed. While it is hard to scale up this method to produce baits for 500 targets, it is
391 well suited to generate commonly used high-copy gene sequences for both archival and recently
392 collected samples. It fits within a narrow ‘Goldilocks’ zone in terms of adequate data for
393 accurately reconstructing phylogenies and relative cost effectiveness with the ability to multiplex
394 at least ~120 individuals per MiSeq® run given the number of loci captured. While the amount of
395 reads on target in our project was not high, we were able to assemble genes of interest for most
396 captured samples.

397 Data obtained showed no evidence for host plant races in OOB. For *Pygovepres*, we could not
398 detect any phylogenetic structure within the species, whereas the structure within
399 *Phallospinophylus* could be explained by distribution. We also reconstructed the phylogeny of
400 the genus *Tuxedo* and sampled all described species, some of which were rarely collected species
401 that are based on specimens from type series.

402 Finally, it is straightforward to combine such data with previously generated data using
403 conventional Sanger sequencing. Commonly used primers for different genes for use in
404 phylogenetic analysis of other groups are easy to add to our protocol. When applied to museum
405 specimens, this approach is optimal for generating complete phylogenetic sampling for clades of
406 interest and relatively cheaply contributing confidently resolved twigs to the Tree of Life.

407 **Acknowledgements**

408 The UCR seed grant “Unlocking the Vault of SoCal Biota” awarded to CW, Amy Litt, and John
409 Heraty and a Dr. Mir S. Mulla and Lelia Mulla Endowed Scholarship awarded to AK are
410 acknowledged for supporting this project. We would like to thank Randall Schuh (AMNH),
411 Serguei Triapitsyn (UCRC), and Fedor Konstantinov (ZISP) for loaning material and permitting
412 dissections and DNA extractions, Randall Schuh (AMNH) and members of the Weirauch lab for
413 reviewing the manuscript.

414 **References**

- 415 Anisimova, M., & Gascuel, O. (2006). Approximate likelihood-ratio test for branches: a fast,
416 accurate, and powerful alternative. *Systematic biology*, 55(4), 539-552. doi:
417 10.1080/10635150600755453
- 418 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., ... &
419 Pevzner P. A. (2012). SPAdes: a new genome assembly algorithm and its applications to
420 single-cell sequencing. *Journal of computational biology*, 19(5), 455-477. doi:
421 10.1089/cmb.2012.0021
- 422 Bi, K., Linderoth, T., Vanderpool, D., Good, J. M., Nielsen, R., & Moritz, C. (2013). Unlocking
423 the vault: next-generation museum population genomics. *Molecular ecology*, 22(24),
424 6018-6032. doi: 10.1111/mec.12516
- 425 Blaimer, B. B., Lloyd, M. W., Guillory, W. X., & Brady, S. G. (2016). Sequence capture and
426 phylogenetic utility of genomic ultraconserved elements obtained from pinned insect
427 specimens. *PloS one*, 11(8), e0161531. doi: 10.1371/journal.pone.0161531
- 428 Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
429 sequence data. *Bioinformatics*, 30(15), 2114-2120. doi: 10.1093/bioinformatics/btu170

- 430 Burrell, A. S., Disotell, T. R., & Bergey, C. M. (2015). The use of museum specimens with high-
431 throughput DNA sequencers. *Journal of human evolution*, 79, 35-44. doi:
432 10.1016/j.jhevol.2014.10.015
- 433 Cassis, G., & Schuh, R. T. (2012). Systematics, biodiversity, biogeography, and host associations
434 of the Miridae (Insecta: Hemiptera: Heteroptera: Cimicomorpha). *Annual review of*
435 *entomology*, 57, 377-404. doi: 10.1146/annurev-ento-121510-133533
- 436 Coddington, J. A., Agnarsson, I., Miller, J. A., Kuntner, M., & Hormiga, G. (2009).
437 Undersampling bias: the null hypothesis for singleton species in tropical arthropod
438 surveys. *Journal of animal ecology*, 78(3), 573-584. doi: 10.1111/j.1365-
439 2656.2009.01525.x
- 440 Cridland, J. M., Ramirez, S. R., Dean, C. A., Sciligo, A., & Tsutsui, N. D. (2018). Genome
441 sequencing of museum specimens reveals rapid changes in the genetic composition of
442 honey bees in California. *Genome biology and evolution*, 10(2), 458-472. doi:
443 10.1093/gbe/evy007
- 444 Cusimano, N., & Renner, S. S. (2010). Slowdowns in diversification rates from real phylogenies
445 may not be real. *Systematic biology*, 59(4), 458-464. doi: 10.1093/sysbio/syq032
- 446 Cusimano, N., Stadler, T., & Renner, S. S. (2012). A new method for handling missing species in
447 diversification analysis applicable to randomly or nonrandomly sampled phylogenies.
448 *Systematic biology*, 61(5), 785-792. doi: 10.1093/sysbio/sys031
- 449 DiEuliis, D., Johnson, K. R., Morse, S. S., & Schindel, D. E. (2016). Opinion: Specimen
450 collections should have a much bigger role in infectious disease research and response.
451 *Proceedings of the National Academy of Sciences*, 113(1), 4-7. doi:
452 10.1073/pnas.1522680112
- 453 Fisher, S., Barry, A., Abreu, J., Minie, B., Nolan, J., Delorey, T. M., ... & Berlin, A. M. (2011). A
454 scalable, fully automated process for construction of sequence-ready human exome
455 targeted capture libraries. *Genome biology*, 12(1), R1. doi: 10.1186/gb-2011-12-1-r1
- 456 Glenn, T. C., Nilsen, R., Kieran, T. J., Finger, J. W., Pierson, T. W., Bentley, K. E., ... &
457 Faircloth, B. C. (2016). Adapterama I: universal stubs and primers for thousands of dual-
458 indexed Illumina libraries (iTru & iNext). *BioRxiv*, 049114. doi: 10.1101/049114
- 459 Hallmann, C. A., Sorg, M., Jongejans, E., Siepel, H., Hofland, N., Schwan, H., ... & de Kroon, H.
460 (2017). More than 75 percent decline over 27 years in total flying insect biomass in
461 protected areas. *PloS one*, 12(10), e0185809. doi: 10.1371/journal.pone.0185809
- 462 Hartley, C. J., Newcomb, R. D., Russell, R. J., Yong, C. G., Stevens, J. R., Yeates, D. K., ... &
463 Oakeshott, J. G. (2006). Amplification of DNA from preserved specimens shows
464 blowflies were preadapted for the rapid evolution of insecticide resistance. *Proceedings of*
465 *the National Academy of Sciences*, 103(23), 8757-8762. doi: 10.1073/pnas.0509590103
- 466 Höhna, S., Stadler, T., Ronquist, F., & Britton, T. (2011). Inferring speciation and extinction rates
467 under different sampling schemes. *Molecular biology and evolution*, 28(9), 2577-2589.
468 doi: 10.1093/molbev/msr095
- 469 Jung, S., & Lee, S. (2011). Molecular phylogeny of the plant bugs (Heteroptera: Miridae) and the
470 evolution of feeding habits. *Cladistics*, 28(1), 50-79. doi: 10.1111/j.1096-
471 0031.2011.00365.x
- 472 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K., von Haeseler, A., & Jermini, L. S. (2017).
473 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature methods*,
474 14(6), 587. doi: 10.1038/nmeth.4285

- 475 Kanda, K., Pflug, J. M., Sproul, J. S., Dasenko, M. A., & Maddison, D. R. (2015). Successful
476 recovery of nuclear protein-coding genes from small insects in museums using Illumina
477 sequencing. *PLoS One*, *10*(12), e0143929. doi: 10.1371/journal.pone.0143929
- 478 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:
479 improvements in performance and usability. *Molecular biology and evolution*, *30*(4), 772-
480 780. doi: 10.1093/molbev/mst010
- 481 Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., ... & Drummond
482 A. (2012). Geneious Basic: an integrated and extendable desktop software platform for
483 the organization and analysis of sequence data. *Bioinformatics*, *28*(12), 1647-1649.
484 10.1093/bioinformatics/bts199
- 485 Konstantinov, F. V., & Knyshov, A. A. (2015). The tribe Bryocorini (Insecta: Heteroptera:
486 Miridae: Bryocorinae): phylogeny, description of a new genus, and adaptive radiation on
487 ferns. *Zoological Journal of the Linnean Society*, *175*(3), 441-472. doi: 10.1111/zoj.12283
- 488 Lanfear, R., Frandsen, P. B., Wright, A. M., Senfeld, T., & Calcott, B. (2016). PartitionFinder 2:
489 new methods for selecting partitioned models of evolution for molecular and
490 morphological phylogenetic analyses. *Molecular Biology and Evolution*, *34*(3), 772-773.
491 doi: 10.1093/molbev/msw260
- 492 Lee, H. Y., Park, M. J., Kim, N. Y., Sim, J. E., Yang, W. I., & Shin, K. J. (2010). Simple and
493 highly effective DNA extraction methods from old skeletal remains using silica columns.
494 *Forensic Science International: Genetics*, *4*(5), 275-280. doi:
495 10.1016/j.fsigen.2009.10.014
- 496 Li, C., Corrigan, S., Yang, L., Straube, N., Harris, M., Hofreiter, M., ... & Naylor, G. J. (2015).
497 DNA capture reveals transoceanic gene flow in endangered river sharks. *Proceedings of*
498 *the National Academy of Sciences*, *112*(43), 13302-13307. doi: 10.1073/pnas.1508735112
- 499 Li, C., Hofreiter, M., Straube, N., Corrigan, S., & Naylor, G. J. (2013). Capturing protein-coding
500 genes across highly divergent species. *Biotechniques*, *54*(6), 321-326. doi:
501 10.2144/000114039
- 502 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler
503 transform. *Bioinformatics*, *25*(14), 1754-1760. doi: 10.1093/bioinformatics/btp324
- 504 Lim, G. S., Balke, M., & Meier, R. (2011). Determining species boundaries in a world full of
505 rarity: singletons, species delimitation methods. *Systematic biology*, *61*(1), 165-169. doi:
506 10.1093/sysbio/syr030
- 507 Liu, S., Wang, X., Xie, L., Tan, M., Li, Z., Su, X., ... & Zhou, X. (2016). Mitochondrial capture
508 enriches mito-DNA 100 fold, enabling PCR-free mitogenomics biodiversity analysis.
509 *Molecular ecology resources*, *16*(2), 470-479. doi: 10.1111/1755-0998.12472
- 510 Maddison, D. R., & Cooper, K. W. (2014). Species delimitation in the ground beetle subgenus
511 *Liocosmius* (Coleoptera: Carabidae: Bembidion), including standard and next-generation
512 sequencing of museum specimens. *Zoological Journal of the Linnean Society*, *172*(4),
513 741-770. doi: 10.1111/zoj.12188
- 514 Mariac, C., Scarcelli, N., Pouzadou, J., Barnaud, A., Billot, C., Faye, A., ... & Couvreur, T. L. P.
515 (2014). Cost-effective enrichment hybridization capture of chloroplast genomes at deep
516 multiplexing levels for population genetics and phylogeography studies. *Molecular*
517 *Ecology Resources*, *14*(6), 1103-1113. doi: 10.1111/1755-0998.12258

- 518 Maricic, T., Whitten, M., & Pääbo, S. (2010). Multiplexed DNA sequence capture of
519 mitochondrial genomes using PCR products. *PloS one*, 5(11), e14004. doi:
520 10.1371/journal.pone.0014004
- 521 McCormack, J. E., Tsai, W. L., & Faircloth, B. C. (2015). Sequence capture of ultraconserved
522 elements from bird museum specimens. *Molecular ecology resources*, 16(5), 1189-1203.
523 doi: 10.1111/1755-0998.12466
- 524 Menard, K. L., Schuh, R. T., & Woolley, J. B. (2014). Total-evidence phylogenetic analysis and
525 reclassification of the Phylinae (Insecta: Heteroptera: Miridae), with the recognition of
526 new tribes and subtribes and a redefinition of Phylini. *Cladistics*, 30(4), 391-427. doi:
527 10.1111/cla.12052
- 528 Minh, B. Q., Nguyen, M. A. T., & von Haeseler, A. (2013). Ultrafast approximation for
529 phylogenetic bootstrap. *Molecular biology and evolution*, 30(5), 1188-1195. doi:
530 10.1093/molbev/mst024
- 531 Mitchell, A. (2015). Collecting in collections: a PCR strategy and primer set for DNA barcoding
532 of decades-old dried museum specimens. *Molecular ecology resources*, 15(5), 1102-1111.
533 doi: 10.1111/1755-0998.12380
- 534 Namyatova, A. A., Konstantinov, F. V., & Cassis, G. (2015). Phylogeny and systematics of the
535 subfamily Bryocorinae with the emphasis on the tribe Dicyphini sensu Schuh, 1976
536 derived from morphological characters. *Systematic Entomology*, 41, 3-40. doi:
537 10.1111/syen.12140
- 538 Nguyen, L. T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2014). IQ-TREE: a fast and
539 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular
540 biology and evolution*, 32(1), 268-274. doi: 10.1093/molbev/msu300
- 541 Paradis, E., Claude, J., & Strimmer, K. (2004). APE: analyses of phylogenetics and evolution in
542 R language. *Bioinformatics*, 20(2), 289-290. doi: 10.1093/bioinformatics/btg412
- 543 Peñalba, J. V., Smith, L. L., Tonione, M. A., Sass, C., Hykin, S. M., Skipwith, P. L., ... & Moritz,
544 C. (2014). Sequence capture using PCR-generated probes: a cost-effective method of
545 targeted high-throughput sequencing for nonmodel organisms. *Molecular Ecology
546 Resources*, 14(5), 1000-1010. doi: 10.1111/1755-0998.12249
- 547 Revell, L. J. (2012). phytools: an R package for phylogenetic comparative biology (and other
548 things). *Methods in Ecology and Evolution*, 3(2), 217-223. doi: 10.1111/j.2041-
549 210X.2011.00169.x
- 550 Rohland, N., & Reich, D. (2012). Cost-effective, high-throughput DNA sequencing libraries for
551 multiplexed target capture. *Genome research*, 22(5), 939-946. doi:
552 10.1101/gr.128124.111
- 553 Schmieder, R., & Edwards, R. (2011). Quality control and preprocessing of metagenomic
554 datasets. *Bioinformatics*, 27(6), 863-864. doi: 10.1093/bioinformatics/btr026
- 555 Schuh, R. T. (2004). Revision of Tuxedo Schuh (Hemiptera: Miridae: Phylinae). *American
556 Museum Novitates*, 3435, 1-26. doi: 10.1206/0003-
557 0082(2004)435<0001:ROTSHM>2.0.CO;2
- 558 Staats, M., Erkens, R. H., van de Vossen, B., Wieringa, J. J., Kraaijeveld, K., Stielow, B., ...
559 & Bakker, F. T. (2013). Genomic treasure troves: complete genome sequencing of
560 herbarium and insect museum specimens. *PLoS One*, 8(7), e69189.
- 561 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of
562 large phylogenies. *Bioinformatics*, 30(9), 1312-1313. doi: 10.1093/bioinformatics/btu033

- 563 Stamatakis, A., Hoover, P., & Rougemont, J. (2008). A rapid bootstrap algorithm for the RAxML
564 web servers. *Systematic biology*, 57(5), 758-771. doi: 10.1080/10635150802429642
- 565 Stork, N. E. (2018). How Many Species of Insects and Other Terrestrial Arthropods Are There on
566 Earth?. *Annual review of entomology*, 63. doi: 10.1146/annurev-ento-020117-043348
- 567 Stork, N. E., McBroom, J., Gely, C., & Hamilton, A. J. (2015). New approaches narrow global
568 species estimates for beetles, insects, and terrestrial arthropods. *Proceedings of the*
569 *National Academy of Sciences*, 112(24), 7519-7523. doi: 10.1073/pnas.1502408112
- 570 Suchan, T., Pitteloud, C., Gerasimova, N. S., Kostikova, A., Schmid, S., Arrigo, N., ... &
571 Alvarez, N. (2016). Hybridization capture using RAD probes (hyRAD), a new tool for
572 performing genomic analyses on collection specimens. *PloS one*, 11(3), e0151651. doi:
573 10.1371/journal.pone.0151651
- 574 Tatarinic, N. J., & Cassis, G. (2012). The Halticini of the world (Insecta: Heteroptera: Miridae:
575 Orthotylinae): generic reclassification, phylogeny, and host plant associations. *Zoological*
576 *Journal of the Linnean Society*, 164(3), 558-658. doi: 10.1111/j.1096-3642.2011.00770.x
- 577 Weirauch, C. (2006a). New genera, new species, and new combinations in western Nearctic
578 Phylini (Heteroptera: Miridae: Phylinae). *American Museum Novitates*, 3521, 1-41. doi:
579 10.1206/0003-0082(2006)3521[1:NGNSAN]2.0.CO;2
- 580 Weirauch, C. (2006b). New genera and species of oak-associated Phylini (Heteroptera: Miridae:
581 Phylinae) from western North America. *American Museum Novitates*, 3522, 1-54. doi:
582 10.1206/0003-0082(2006)3522[1:NGASOO]2.0.CO;2
- 583 Weirauch, C., Seltmann, K. C., Schuh, R. T., Schwartz, M. D., Johnson, C., Feist, M. A., &
584 Soltis, P. S. (2016). Areas of endemism in the Nearctic: a case study of 1339 species of
585 Miridae (Insecta: Hemiptera) and their plant hosts. *Cladistics*, 33(3), 279-294. doi:
586 10.1111/cla.12169
- 587 Yang, D. Y., Eng, B., Wayne, J. S., Dudar, J. C., & Saunders, S. R. (1998). Improved DNA
588 extraction from ancient bones using silica-based spin columns. *American journal of*
589 *physical anthropology*, 105(4), 539-543. doi: 10.1002/(SICI)1096-
590 8644(199804)105:4<539::AID-AJPA10>3.0.CO;2-1
- 591 Yin, T., Cook, D., & Lawrence, M. (2012). ggbio: an R package for extending the grammar of
592 graphics for genomic data. *Genome biology*, 13(8), R77. doi: 10.1186/gb-2012-13-8-r77
- 593 Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T. Y. (2016). ggtree: an R package for
594 visualization and annotation of phylogenetic trees with their covariates and other
595 associated data. *Methods in Ecology and Evolution*, 8(1), 28-36. doi: 10.1111/2041-
596 210X.12628
- 597

598 **Data Accessibility:**

599 - DNA sequences: GenBank accessions [annotated mitochondrial genomes, ribosomal genes, and
600 dynein fragments for baits will be uploaded to GenBank, and accession numbers will be indicated
601 in Table S3]; NCBI SRA: SRP136090, accession numbers for individual samples are indicated in
602 the Table S3.

603 - Final DNA sequence alignments and partitioning schemes: will be uploaded to Dryad
604 repository.
605 - Voucher specimen information including photographs: available through the Plant Bug
606 Planetary Biodiversity Inventory Project website
607 (<http://research.amnh.org/pbi/heteropteraspeciespage/>), linked to the unique specimen identifier
608 (See Table 1, the USI column) [photographs are in the process of being uploaded].

609
610 **Author contributions**
611 AK, ERLG & CW designed the research, AK and CW performed the research, AK analyzed the
612 data and AK, ERLG and CW wrote the paper.

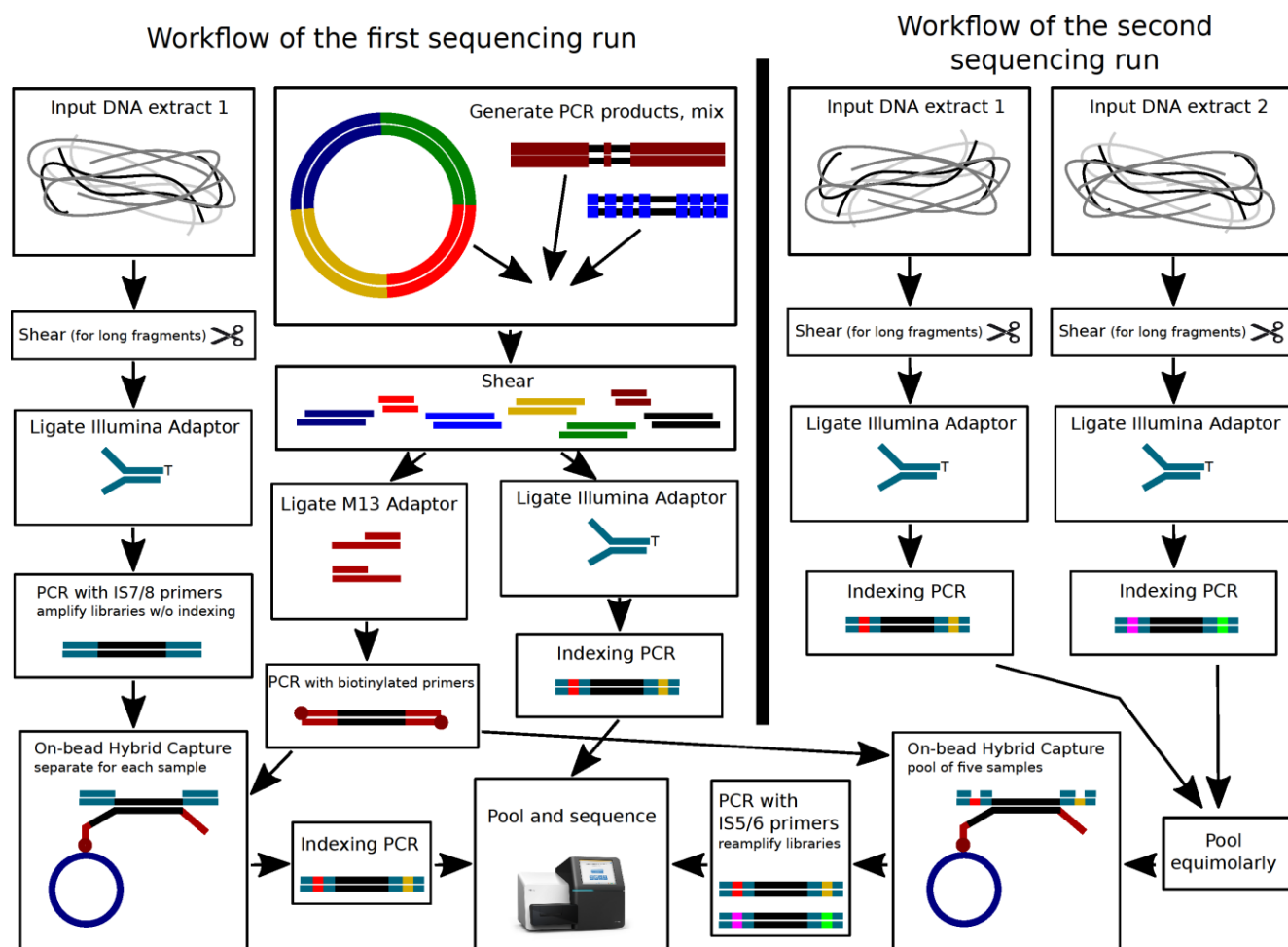
613

614 Tables and Figures

615 Table 1. List of samples used in the project, voucher specimen information, and sequencing
616 information.

PIH	USI	Species	Type status	Geographic region	Host plant	Preservation method	Collecting year	Collecting year	(d/d/M), ng/dl	DNA mass used for library prep (ng)	Amount of bait used, ng	PF Clusters	Filtered and merged reads	Filtered and merged reads (duplet)	Duplication rate	Mito: mapped reads	Mito: percent mapped	Mito: coverage depth	18S: mapped reads	18S: percent mapped	18S: coverage depth	28S: mapped reads	28S: percent mapped	28S: coverage depth	Dyname: mapped reads	Dyname: percent mapped	Total reads mapped	CGM distance between a bait and a sample
617	LCKr EN1 00127382	<i>Rufocornutus elegans</i>	juv	Cleveland NF	Quercus	HOH	2009	0.78	11.04	10	130678	130677	130680	0.90%	1041	0.81%	118616	85	0.02%	309444	41	0.01%	844417	7	0.00%	1116	27.09%	
618	LCKr EN1 00127383	<i>Chrysophyllum ruminatum</i>	juv	Tahiti	Quercus	dry	2009	0.52	18.24	10	131628	131630	131743	0.93%	1394	0.82%	110312	847	0.04%	307441	847	0.04%	847	219	0.01%	1626	11.51%	
619	LCKr EN1 00127384	<i>Rufocornutus elegans</i>	juv	San Antonio TX	Quercus	dry	2009	0.73	14.4	10	131774	131783	131930	0.92%	105	0.04%	127733	845	0.04%	307	258	0.01%	21727	1439	0.01%	1539	15.98%	
620	LCKr EN1 00127379	<i>Rufocornutus elegans</i>	juv	San Antonio TX	Quercus	dry	2009	0.616	8.24	10	130228	130236	130274	1.70%	1735	1.37%	112083	847	0.04%	307	309	0.01%	25286	1439	0.01%	4495	10.31%	
621	LCKr EN1 00127375	<i>Chrysophyllum ruminatum</i>	juv	Cleveland NF	Quercus	dry	2009	0.68	19.36	10	130484	130484	130493	0.92%	1122	0.82%	112063	845	0.04%	307	847	0.04%	847	219	0.01%	1626	11.51%	
622	LCKr EN1 00127372	<i>Rufocornutus elegans</i>	juv	Raja	Quercus	dry	2009	0.308	8.24	10	131679	131680	131684	0.93%	1137	0.83%	110284	846	0.04%	307	48137	378	0.01%	251485	1209	0.01%	1810	14.82%
623	LCKr EN1 00127380	<i>Chrysophyllum ruminatum</i>	juv	Lak Forest	Quercus	dry	2009	0.57	7.6	10	132670	132670	132680	0.93%	1020	1.14%	20288	847	0.04%	307	318	0.01%	16137	1439	0.01%	1626	11.51%	
624	LCKr EN1 00127381	<i>Chrysophyllum ruminatum</i>	juv	Cleveland NF	Quercus	dry	2009	1.26	16.36	10	131511	131511	131511	1.01%	13172	1.01%	112167	845	0.04%	307	847	0.04%	847	219	0.01%	1626	11.51%	
625	LCKr EN1 00127378	<i>Rufocornutus elegans</i>	juv	Cleveland NF	Quercus	dry	2009	0.748	19.8	10	132130	132132	132135	2.61%	1839	1.43%	21774	847	0.04%	307	74787	1439	0.01%	118	4329	0.04%	1131	9.43%
626	LCKr EN1 00127376	<i>Chrysophyllum ruminatum</i>	juv	Cleveland NF	Quercus	dry	2009	0.606	19.36	10	132134	132132	132135	2.61%	1839	1.43%	21774	847	0.04%	307	74787	1439	0.01%	118	4329	0.04%	1131	9.43%
627	MANNr PB 00082132	<i>Rufocornutus elegans</i>	juv	Maharaj	Quercus	dry	1978	2.2	18	10	134450	134450	134458	1.49%	398	0.44%	10389	1327	0.04%	307	847	0.04%	847	219	0.01%	1626	11.51%	
628	LCKr EN1 00127385	<i>Chrysophyllum ruminatum</i>	juv	San Bernardino CA	Quercus	dry	2009	0.3	47.66	10	130833	130836	130838	0.93%	1050	1.16%	24386	845	0.04%	307	847	0.04%	847	219	0.01%	1626	11.51%	
629	LCKr EN1 00127386	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	0.375	18	10	131511	131511	131511	1.01%	13172	1.01%	112167	845	0.04%	307	847	0.04%	847	219	0.01%	1626	11.51%	
630	LCKr EN1 00127389	<i>Rufocornutus elegans</i>	juv	San Antonio TX	Quercus	dry	2009	2.44	19.24	10	130278	130279	130280	0.93%	1043	1.04%	21243	110	0.02%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
631	LCKr EN1 00127390	<i>Rufocornutus elegans</i>	juv	San Antonio TX	Quercus	dry	2009	1.84	19.48	10	130317	130318	130320	0.92%	1029	0.82%	18263	845	0.04%	307	84446	117	0.01%	11029	1439	0.01%	1626	11.51%
632	LCKr EN1 00127374	<i>Rufocornutus elegans</i>	juv	Raja	Quercus	dry	2009	0.276	9.04	10	130634	130631	130634	10.10%	10941	1.01%	98302	110	0.00%	11355	306	0.01%	11544	1439	0.01%	1626	11.51%	
633	LCKr EN1 00127387	<i>Chrysophyllum ruminatum</i>	juv	Maharaj	Quercus	dry	2009	1.18	14.74	10	130444	130444	130444	0.92%	1115	0.82%	110366	274	0.01%	307	847	0.04%	847	219	0.01%	1626	11.51%	
634	LCKr EN1 00127384	<i>Chrysophyllum ruminatum</i>	juv	Maharaj	Quercus	dry	2009	1.18	14.74	10	130444	130444	130444	0.92%	1115	0.82%	110366	274	0.01%	307	847	0.04%	847	219	0.01%	1626	11.51%	
635	MANNr PB 00082138	<i>Rufocornutus elegans</i>	juv	Sharda CA	Quercus	dry	1984	7.72	77.2	10	130644	130649	130648	0.73%	1888	0.88%	110113	1151	0.02%	40246	2184	0.01%	110797	1439	0.01%	1626	11.51%	
636	LCKr EN1 00127393	<i>Chrysophyllum ruminatum</i>	juv	Cleveland NF	Quercus	dry	2009	1.16	17.2	10	130414	130419	130420	12.71%	1424	0.82%	110366	127	0.01%	307	847	0.04%	847	219	0.01%	1626	11.51%	
637	LCKr EN1 00127391	<i>Chrysophyllum ruminatum</i>	juv	Cleveland NF	Quercus	dry	2009	0.54	18.1	10	130940	130939	130940	0.90%	1026	0.72%	110366	110	0.01%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
638	LCKr EN1 00127392	<i>Chrysophyllum ruminatum</i>	juv	Cleveland NF	Quercus	dry	2009	1.07	14.8	10	130940	130939	130940	0.90%	1026	0.72%	110366	110	0.01%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
639	LCKr EN1 00127390	<i>Chrysophyllum ruminatum</i>	juv	San Bernardino CA	Quercus	dry	2009	1.5	18	10	131172	131170	131160	7.42%	1429	0.82%	45483	174	0.01%	307	847	0.04%	847	219	0.01%	1626	11.51%	
640	LCKr EN1 00127396	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	2.07	17.2	10	130979	130979	130979	0.93%	1026	0.72%	110366	110	0.01%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
641	LCKr EN1 00127395	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	1.416	14.98	10	130979	130979	130979	0.93%	1026	0.72%	110366	110	0.01%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
642	LCKr EN1 00127397	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	4.44	19.48	10	130278	130279	130280	0.93%	1043	1.04%	21243	110	0.02%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
643	LCKr EN1 00127398	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	4.06	19.48	10	130278	130279	130280	0.93%	1043	1.04%	21243	110	0.02%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
644	LCKr EN1 00127399	<i>Chrysophyllum ruminatum</i>	juv	Tahiti	Quercus	dry	2009	4.26	22.42	10	130484	130482	130485	7.20%	1847	0.81%	847	847	0.04%	847	847	0.04%	847	219	0.01%	1626	11.51%	
645	MANNr PB 00082143	<i>Chrysophyllum ruminatum</i>	juv	Sharda CA	Quercus	dry	1978	1.44	13.6	10	130628	130627	130633	0.93%	1044	0.82%	110113	1151	0.02%	40246	2184	0.01%	110797	1439	0.01%	1626	11.51%	
646	MANNr PB 00082144	<i>Chrysophyllum ruminatum</i>	juv	Sharda CA	Quercus	dry	1978	1.44	13.6	10	130628	130627	130633	0.93%	1044	0.82%	110113	1151	0.02%	40246	2184	0.01%	110797	1439	0.01%	1626	11.51%	
647	MANNr PB 00082145	<i>Chrysophyllum ruminatum</i>	juv	San Bernardino CA	Quercus	dry	1977	0.674	10.12	10	130771	130784	130784	10.37%	1024	0.81%	89248	1373	0.01%	97404	168	0.01%	105404	1730	0.01%	1626	11.51%	
648	LCKr EN1 00127400	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	0.296	18	10	131141	131141	131141	0.92%	1115	0.82%	110366	274	0.01%	307	847	0.04%	847	219	0.01%	1626	11.51%	
649	LCKr EN1 00127401	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	4.3	18.4	10	131187	131186	131185	0.92%	1019	0.91%	103659	1006	0.04%	307	74787	1439	0.01%	118	4329	0.04%	1131	9.43%
650	LCKr EN1 00127402	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	4.36	17.6	10	131187	131186	131185	0.92%	1019	0.91%	103659	1006	0.04%	307	74787	1439	0.01%	118	4329	0.04%	1131	9.43%
651	LCKr EN1 00127403	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	4.36	17.6	10	131187	131186	131185	0.92%	1019	0.91%	103659	1006	0.04%	307	74787	1439	0.01%	118	4329	0.04%	1131	9.43%
652	LCKr EN1 00127404	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	0.307	18.08	10	130284	130284	130284	1.00%	221	0.24%	41244	412	0.02%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
653	LCKr EN1 00127405	<i>Chrysophyllum ruminatum</i>	juv	Raja	Quercus	dry	2009	1.2	14.8	10	130284	130284	130284	1.00%	221	0.24%	41244	412	0.02%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
654	LCKr EN1 00127406	<i>Chrysophyllum ruminatum</i>	juv	Raja	Quercus	dry	2009	1.2	14.8	10	130284	130284	130284	1.00%	221	0.24%	41244	412	0.02%	31235	189	0.01%	24137	1439	0.01%	1626	11.51%	
655	LCKr EN1 00127407	<i>Chrysophyllum ruminatum</i>	juv	San Antonio TX	Quercus	dry	2009	1.15	8.2	10	130847	130845	130850	1.49%	398	0.44%	110366	143	0.01%	307	847	0.04%	847	219	0.01%	1626	11.51%	
656	LCKr EN1 0012																											

620 Figure 1. Procedure flowchart.

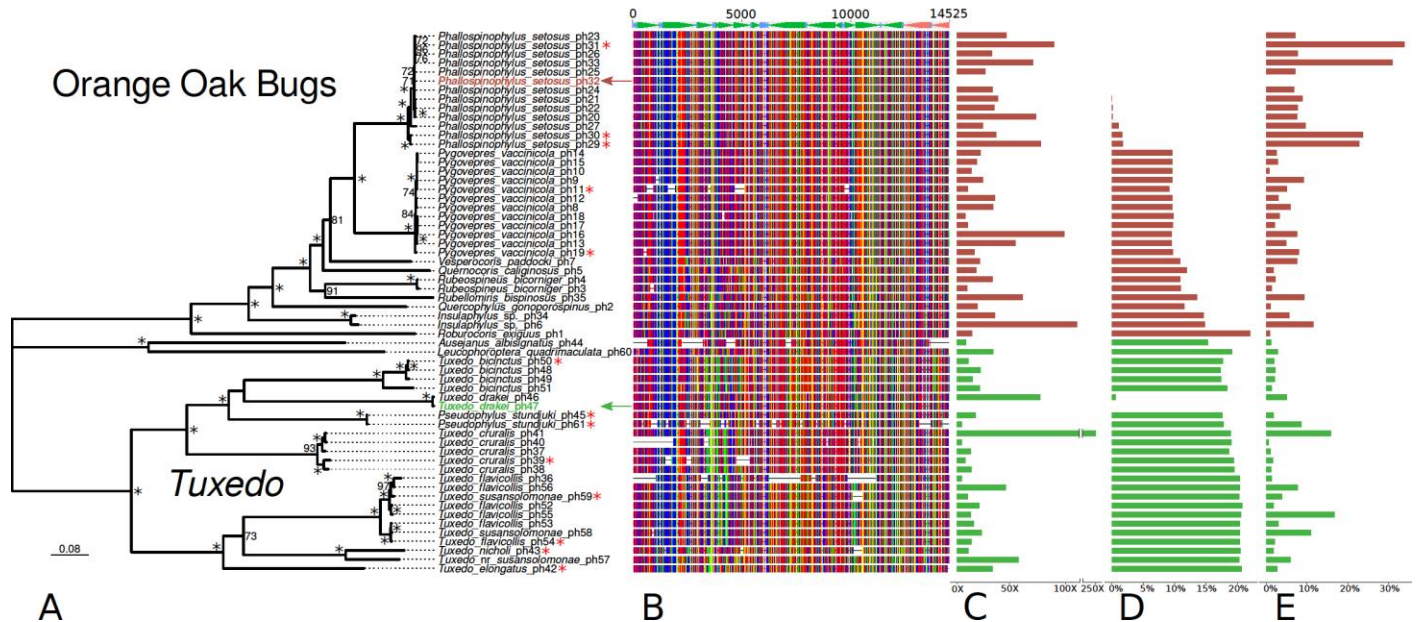


621
622

623 Figure 2. A. Combined phylogeny of the OOB and *Tuxedo* subprojects, generated in RAxML,
 624 values at nodes represent Rapid Bootstrap Support, values below 70 are not shown, asterisks
 625 indicate full support, arrows denote bait samples for the OOB (red) and *Tuxedo* (green)
 626 subprojects, red asterisks denote samples older than 20 years. B. Mitochondrial alignment
 627 completeness, control region excluded. C. Average coverage of mitochondrial contig(s), control
 628 region excluded. D. Pairwise COX1 distances between a bait and a captured sample. E. Total
 629 percent of reads mapping to target including mitochondrial genome, 18S, 28S, and dynein.

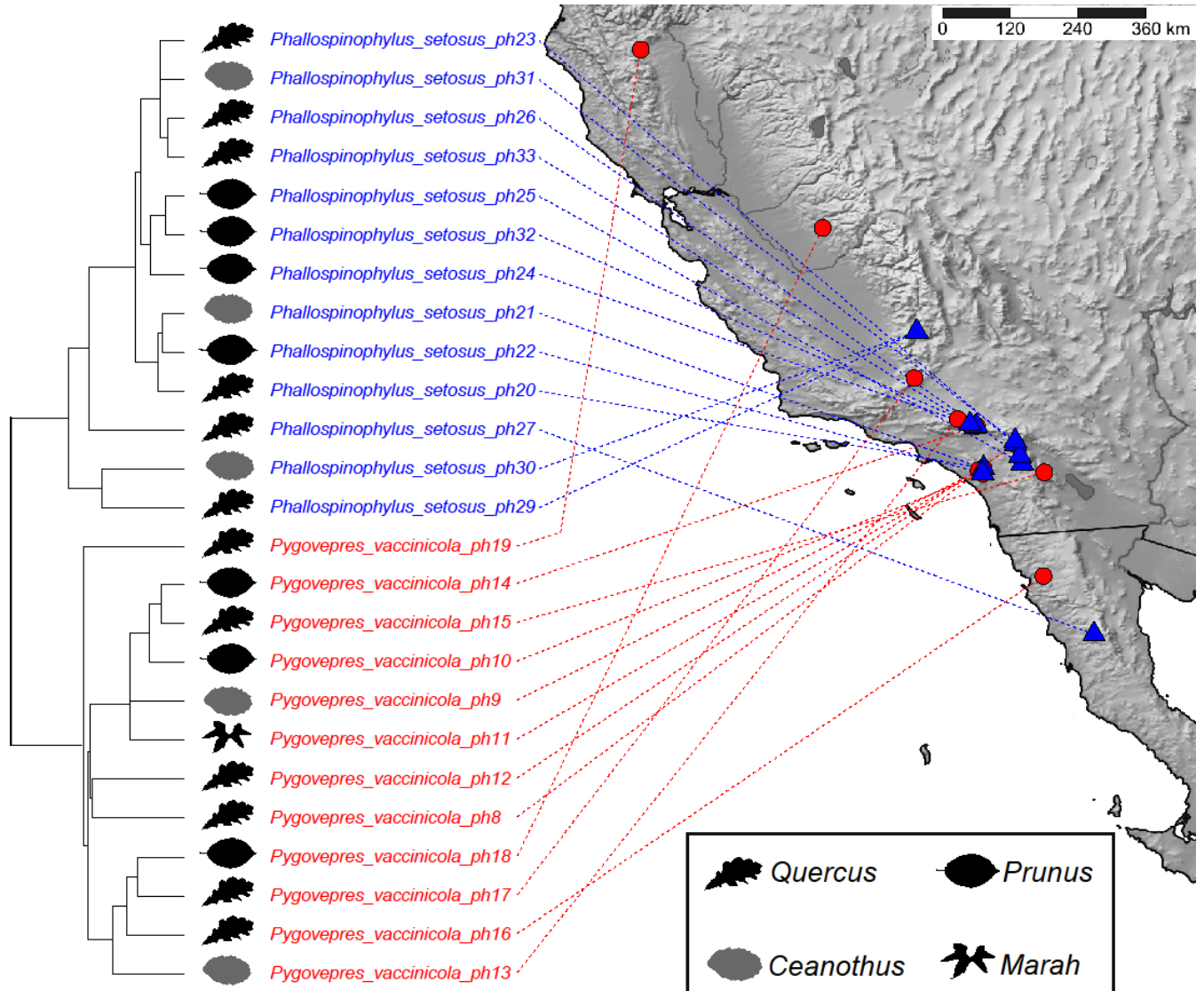
630

631
 632



633 Figure 3. Host and distribution data for the Orange Oak Bug subproject, aligned with phylogeny
634 (branches not to scale) and with host plant of specimens mapped using representative leaf shapes
635 of plant genus.

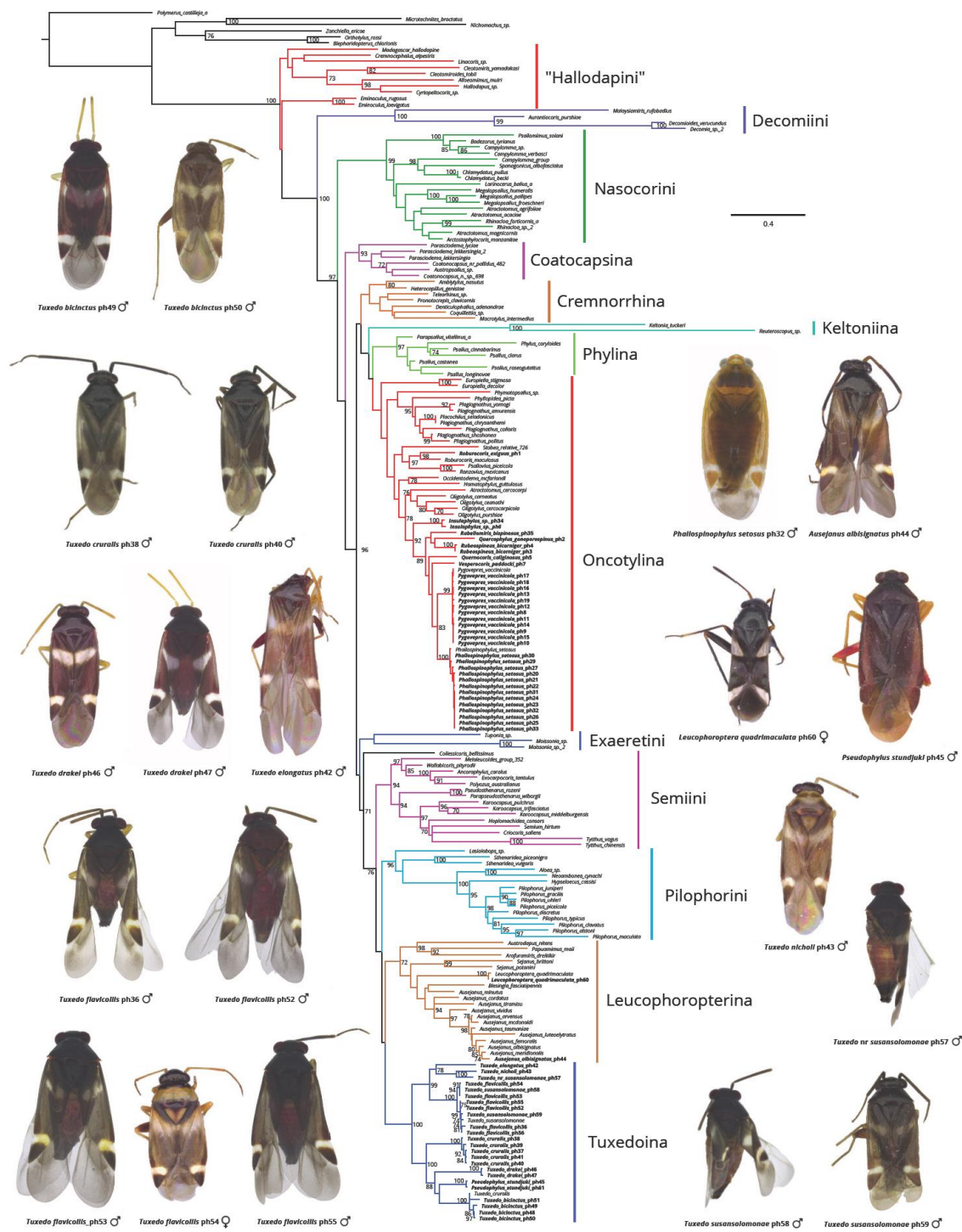
636



637

638

639 Figure 4. Phylogeny of Phylinae, generated in RAXML, with specimens for which new data was
 640 gathered in bold font, values at nodes represent Rapid Bootstrap Support, values below 70 not
 641 shown.



642