

1 **Understanding UCEs: A comprehensive primer on using Ultraconserved Elements for**  
2 **arthropod phylogenomics**

3  
4 **Y. Miles Zhang\***, Jason L. Williams, Andrea Lucky

5 **University of Florida, Department of Entomology & Nematology, Gainesville, FL, 32608**

6 \* **Corresponding author email: [Yuanmeng.zhang@gmail.com](mailto:Yuanmeng.zhang@gmail.com)**

7 **Abstract:**

8 Targeted enrichment of ultraconserved elements (UCE) has emerged as a promising tool for  
9 inferring evolutionary history in many taxa, with utility ranging from phylogenetic and  
10 phylogeographic questions at deep time scales to population level studies at shallow time scales.  
11 However, the methodology can be daunting for beginners. Our goal is to introduce UCE  
12 phylogenomics to a wider audience by summarizing recent advances in arthropod research, and  
13 to familiarize readers with background theory and steps involved. We define terminology used in  
14 association with the UCE approach, evaluate current laboratory and bioinformatic methods and  
15 limitations, and, finally, provide a roadmap of steps in the UCE pipeline to assist  
16 phylogeneticists in making informed decisions as they employ this powerful tool. By facilitating  
17 increased adoption of UCE in phylogenomics studies that deepen our comprehension of the  
18 function of these markers across widely divergent taxa, we aim to ultimately improve  
19 understanding of the arthropod tree of life.

20 **Keywords: Arachnida, Insecta, Phylogenomics Methods, Target Enrichment,**  
21 **Ultraconserved Elements**

22

## 23 **Introduction & Background**

24           The advent of massively-parallel sequencing technology and the subsequent emergence  
25 of the field of phylogenomics has invigorated evolutionary biology in a relatively short time span  
26 (reviewed in Philippe et al. 2011, Jones and Good 2016). This molecular revolution has offered  
27 unprecedented opportunities to generate large-scale datasets, and with the concurrent explosion  
28 of analytic and bioinformatics tools, has made it possible to address previously intractable  
29 challenges due to limited genetic markers. However, the rapidity with which new technologies  
30 have emerged has made it difficult for scientists to stay up to date about useful new tools;  
31 understanding the steps involved in using new methods presents a challenge for researchers.

32           Genome-scale studies are rapidly supplanting the Sanger sequencing-based, multi-locus  
33 molecular phylogenetic methods that dominated from the mid-1990's through the early 2000's;  
34 today, genomic-scale studies dwarf previous approaches in the sheer scale of data they generate  
35 (Bravo et al. 2019). While the cost and scale of whole-genome sequencing are prohibitive for  
36 many researchers, recent advances in sequencing technology and laboratory protocols have made  
37 it possible to generate high quality genomic datasets using a combination of next-generation  
38 sequencing, genomic reduction, and sample multiplexing (Lemmon and Lemmon 2013,  
39 McCormack et al. 2013a). These so-called 'genome reduction' or 'reduced representation'  
40 approaches can rapidly generate datasets with thousands of loci, at relatively low cost, for model  
41 and non-model taxa alike. Methods such as restriction enzyme-associated DNA sequencing  
42 (RADseq; Miller et al. 2007, Baird et al. 2008, Peterson et al. 2012), transcriptomics (Bi et al.  
43 2012), and target enrichment methods such as Anchored Hybrid Enrichment (AHE) (Lemmon et

44 al. 2012) or target capture of Ultraconserved Elements (UCE) (Faircloth 2017) are now widely  
45 used for generating genomic-scale data for phylogenomic studies. These phylogenomics methods  
46 are similar in some respects, but each has strengths and weaknesses which may not be easily  
47 discerned by researchers new to this field. Because of the proliferation of new approaches and  
48 tools in phylogenetics, selecting a method to use in the era of ‘big data’ can be daunting.  
49 Potential users need guidance in choosing methods appropriate to their research questions, and in  
50 navigating confusing terminologies, bioinformatics-heavy data processing, and computationally  
51 intensive analyses.

52 UCE-based phylogenomics continues to develop rapidly, and the lack of comprehensive  
53 review has been a significant challenge for potential users to overcome when exploring this  
54 option. This paper summarizes recent advances in UCE phylogenomics in arthropod research;  
55 we start by familiarizing readers with background theory and terminology, and describing the  
56 steps involved in generating and analyzing UCE data, and then provide quality-control tips to  
57 ensure that data collection and downstream analyses can be performed with confidence.

58

## 59 **What are UCES?**

60 Ultraconserved Elements are highly-conserved regions within the genome that are shared  
61 among evolutionarily distant taxa (Bejerano et al. 2004). The DNA adjacent to each ‘core’ UCE  
62 region, known as flanking DNA, increases in variability with distance from the region (Faircloth  
63 et al. 2012). UCES and flanking regions can be selectively captured, and used to reconstruct the  
64 evolutionary history of taxa at various time scales, from deep to shallow phylogenetic  
65 divergences (Faircloth et al. 2012, McCormack et al. 2012).

66           The UCE approach belongs to the broad category of ‘target enrichment’ phylogenomic  
67 techniques, which involve selective capture of genomic regions from DNA prior to sequencing  
68 (Mamanova et al. 2010). Similar methods include AHE (Anchored Hybrid Enrichment),  
69 BaitFisher (Mayer et al. 2016), and Hyb-Seq (Weitemier et al. 2014). AHE has been the most  
70 widely used method for animal studies, to date, but all target enrichment methods have been  
71 successfully used across a variety of taxa. These techniques universally involve identifying loci  
72 of interest, then designing custom-made molecular probes (also known as baits) which are  
73 hybridized to loci of interest, and ultimately sequencing selected genomic region on a massively-  
74 parallel platform. The main difference between the AHE and UCE approaches is the nature of  
75 the loci targeted; AHE focuses on fewer loci (300-600) that are exclusively exonic, while UCEs  
76 target more loci (>1000) using fewer probes – these may include both exonic and intronic  
77 regions, depending on the organism (Crawford et al. 2012, McCormack et al. 2012, Faircloth et  
78 al. 2015). While AHE can cope with sequence variation at target loci by using a more diverse set  
79 of probes per locus, the details of the methodology are not available for scrutiny as they are, in  
80 part, proprietary (Lemmon et al. 2012). The UCE approach, in contrast, is fully open source,  
81 which has contributed to recent interest in using these markers for arthropod phylogenomics.

82

### 83 **Advantages of UCE Phylogenomics**

84           The UCE approach has become an increasingly popular target enrichment method for  
85 generating phylogenomic data, as it offers advantages over traditional Sanger sequencing  
86 methods in terms of quantity of data generated. UCEs have successfully been used in studies  
87 across a broad array of animal taxa including birds (McCormack et al. 2013b, Musher and  
88 Cracraft 2018), mammals (McCormack et al. 2012, Mclean et al. 2018), fish (Faircloth et al.

89 2013, Alda et al. 2018), amphibians (Newman and Austin 2016, Zarza et al. 2018), reptiles  
90 (Crawford et al. 2012, Streicher and Wiens 2017, Myers et al. 2019), sponges (Ryu et al. 2012),  
91 cnidarians (Quattrini et al. 2018), echinoderms (Ryu et al. 2012), and arthropods (Faircloth et  
92 al. 2015, Baca et al. 2017b, Branstetter et al. 2017c, Hedin et al. 2018b, Kieran et al. 2019).  
93 These studies range widely in evolutionary scale, from phylogenetic and phylogeographic  
94 questions at deep time scales (Faircloth et al. 2013, Smith et al. 2014, Branstetter et al. 2017c) to  
95 population level studies at shallow time scales (Harvey et al. 2016, Manthey et al. 2016, Zarza et  
96 al. 2018, Branstetter and Longino 2019, Myers et al. 2019).

97 Benefits of using UCE include openly shared resources such as probe sets  
98 (<https://www.ultraconserved.org/>), lab protocols (<https://baddna.uga.edu/protocols.html/>), and  
99 bioinformatics tools (<https://phyluce.readthedocs.io/en/latest/>), making it an easy method to learn  
100 and use in comparison to more proprietary alternatives such as AHE. Complete library  
101 preparation for around 100 samples can be completed in approximately two weeks by one  
102 person, or a month if counting DNA extraction and possible troubleshooting. UCE datasets can  
103 be easily standardized, even from multiple studies, by using the same probe set. In this way, data  
104 from studies using the same probe set or with exon/transcriptome data (Bossert et al. 2019,  
105 Kieran et al. 2019) can be combined and can incorporate legacy methods if the probe set includes  
106 Sanger genes (Branstetter et al. 2017a). These are distinct advantages over restriction enzyme-  
107 based methods such as traditional RADseq, which lacks repeatability due to the random nature of  
108 the restriction enzyme digestion that generates random genomic fragments. An additional  
109 advantage of target enrichment methods is the high success rate with degraded or low-quantity  
110 samples; older, dried museum specimens may be unusable in traditional restriction enzyme-  
111 based and transcriptomics studies as they require large quantities of high-quality DNA or RNA

112 from fresh or carefully-preserved tissues (Blaimer et al. 2016a, Lim and Braun 2016, Ruane and  
113 Austin 2017). It is worth noting that newer RAD-based methods such as RADcap (Hoffberg et  
114 al. 2016), Rapture (Ali et al. 2016), and hyRAD (Suchan et al. 2016) address these limitations by  
115 using a combination of restriction enzyme digestion and hybridization capture probes to  
116 overcome traditional RAD-based problems such as allele dropout, and can successfully capture  
117 degraded DNA from older museum samples.

### 118 **UCE and Arthropods**

119 The UCE approach was first demonstrated outside of vertebrates in the insect order  
120 Hymenoptera (See Table 1). To date, two published probe sets exist for Hymenoptera: hym-v2  
121 (31,829 probes for 2,590 UCEs, Branstetter et al. 2017a) includes most of the original hym-v1  
122 probe sets (2,749 probes for 1,510 UCEs, Faircloth et al. 2015), and excludes poorly performing  
123 loci. Other arthropods groups for which published UCE data exist include Arachnida,  
124 Coleoptera, and Hemiptera; as well as multiple upcoming studies for Diptera (E. Buenaventura,  
125 C. Cohen, K. Noble, *pers. comm.*). Psocodea and Lepidoptera probe sets have been developed  
126 but not yet tested *in vitro* (Table 1). The utility of UCEs extends beyond purely phylogenetic and  
127 taxonomic research. For example, UCE-based community phylogenomics has been used to  
128 reveal the importance of bee phylodiversity in agriculture (Grab et al. 2019); UCE-based  
129 phylogeny and geometric morphometrics have been used in combination to explore the evolution  
130 of parasitic wasp body shape (Santos et al. 2019); single nucleotide polymorphisms (SNPs)  
131 generated from UCE data have been used to demonstrate the success of unsupervised machine  
132 learning in species delimitation of harvestmen (Derkarabetian et al. 2019), and UCE phylogenies  
133 have been integrated with environmental niche modeling to examine phylogeographic patterns of  
134 ants across the Brazilian Atlantic Forest (Ströher et al. 2019).

135           This study compiles all currently available UCE-based literature related to arthropods as  
136 of July 2019 (n = 32, Figure 1), but will undoubtedly increase exponentially in the future (Table  
137 1). Our aim is to provide a step-by-step guide to make the UCE research pipeline more  
138 approachable to researchers working across different arthropod groups.

139

## 140 **UCE Phylogenomics Pipeline**

141           The steps in the UCE pipeline are 1) probe selection and design; 2) wet lab work and  
142 sequencing; 3) bioinformatics; and 4) phylogenomic analyses. Below, we visualize the process in  
143 a workflow diagram (Figure 2) and describe the choices a researcher must make at each stage. A  
144 glossary of technical terms is provided as a supplementary document (S1).

145

### 146 **Probe Selection & Design**

147           Probe sets are a collection of oligonucleotides that will bind to specific, conserved  
148 genome regions of interest, often called baits as they can ‘fish’ out the region of interest. These  
149 probes are sometimes interchangeably called ‘baits’ as they are used to fish out target loci from a  
150 ‘pond’ of randomly sheared, adaptor-ligated DNA (Gnirke et al. 2009). However, to avoid  
151 confusion we recommend reserving the term ‘baits’ for the intermediate stage in probe design;  
152 by contrast, ‘probes’ refer to the final products that is synthesized for commercial use (Gustafson  
153 et al. 2019). A probe set functions through a collection of biotinylated oligonucleotides that are  
154 designed to bind with specific genome regions of interest. Probes are combined with denatured  
155 and cooled DNA, allowing for ‘in solution’ hybridization to targets. Streptavidin-coated

156 magnetic beads, which have high affinity for biotin, are added into the solution. The beads then  
157 bind to the probe-DNA hybrids through the biotin on the probe set. Any unwanted DNA  
158 fragments are then washed away, leaving only the desired regions attached to the beads (Gnirke  
159 et al. 2009).

160 Probes are designed based on UCE loci identified from published genomes for each  
161 taxonomic group. Currently probe sets for Arachnida, Coleoptera, Diptera, Hemiptera, and  
162 Hymenoptera are available for purchase through Arbor Biosciences  
163 (<https://arborbiosci.com/products/uces/>). Other taxonomic groups either have no probe sets  
164 available, or have not been tested *in vitro*. Designing new probe sets may prove challenging in  
165 the absence of published genomes for a group of interest. Nevertheless, low coverage genome  
166 sequencing (5x) may be an increasingly affordable appropriate first step (Zhang et al. 2019). The  
167 sequenced genomes selected as the basis for probe design should ideally reflect diversity within  
168 the group of interest. Ideally multiple genomes should be used for probe design, but minimally  
169 probe sets designed based on only two genomes (hym-v1) were shown to be successful in  
170 capturing UCEs across the diverse order Hymenoptera (Faircloth et al. 2015).

171 Whereas probe selection is straightforward (they are either available for the group of  
172 interest or they are not), probe design for new taxonomic groups is a time-consuming process for  
173 any target enrichment method, as the probe sets can differ in number and composition depending  
174 on the target taxa and evolutionary scale. Currently published probe sets for arthropods target  
175 1,100 – 2,700 UCEs loci, and have been made publicly available under public domain license  
176 (CC-0), thus allowing for restriction-free commercial synthesis, testing, use and improvements  
177 by other research groups (<http://ultraconserved.org/#protocols>) (Branstetter et al. 2017a,  
178 Faircloth 2017, Gustafson et al. 2019). A generalized workflow for identifying conserved



179 sequences shared among divergent genomes and enrichment probes design is available (Faircloth  
180 2017), and a new pipeline has been described using low-coverage genome sequencing that can  
181 also be used to design UCE probes (Zhang et al. 2019). In brief, the probe design sequence is 1)  
182 select base genome (s); 2) generate short reads as exemplars of the focal group's diversity and  
183 align to base genome(s); 3) merge approximate reads and find overlapping regions shared among  
184 exemplar taxa and base genome (conserved regions); 4) design temporary bait set from base  
185 genome against conserved regions and align to exemplar genome assemblies to remove  
186 duplicates; 5) design exemplar-specific probes for each locus where temporary baits match  
187 exemplar genome assemblies.

188         How to best optimize the probe design process is an area of active research. Both base  
189 genome choice and initial bait design stringency parameters can greatly affect the number of  
190 resultant probes and, subsequently, the number of loci detected and recovered in Adepagan  
191 beetles (Gustafson et al. 2019). The optimal base genome can be selected by conducting a base  
192 genome experiment by iteratively selecting each taxon as the base genome and finding candidate  
193 loci shared among exemplar taxa, or selected from taxon with the smallest average genetic  
194 distance to the other exemplar taxa through independently generated Sanger markers. Probe sets  
195 can also be modified to incorporate additional loci. The Hymenoptera probe set hym-v1 was  
196 improved by the publication of hym-v2, which included most of the original hym-v1 loci as well  
197 as new loci and probes targeting 16 commonly-sequenced nuclear genes to allow for 'back  
198 compatibility' with Sanger-era data (Branstetter et al. 2017a). The resulting capability of  
199 combining new genomic data with older sequences obtained from 'legacy' markers is vital to  
200 phylogenetic studies, as DNA quality tissue for many rare but vital taxa to phylogenetic studies  
201 may be difficult or impossible to obtain repeatedly. *In silico* tests of existing probe sets

202 demonstrate moderate success with sister outgroups, such as using the Hemiptera probe set to  
203 capture UCEs from thrips (Insecta, Order: Thysanoptera) (Faircloth 2017). Importantly, expense,  
204 time, and computational resources needed should be taken into consideration when designing  
205 new probe sets. The cost of development should be weighed against the potential future use of  
206 the probe set beyond the initial study; UCE probes for larger clades, for example, may be more  
207 likely to be adopted for multiple uses than those designed for species-poor groups.

## 208 **Wet Lab Work & Sequencing**

209 *Specimen selection and DNA extraction.* Selecting appropriate specimens for DNA  
210 extraction is vital to any phylogenetic endeavor. The first major requirement for molecular  
211 phylogenetics is to capture high-quality DNA. DNA capture success rates can be negatively  
212 affected by specimen age and preservation method (Short et al. 2018); Arthropod studies are  
213 often limited by DNA degradation, as most natural history collections have historically preserved  
214 specimens dry (pinned) or stored in 70% ethanol at room temperature which can lead rapidly  
215 deterioration of DNA (Short et al. 2018). Other complications include the number of freeze/thaw  
216 cycles (as few as possible), and the number/frequency of alcohol changes (regular enough to  
217 maintain 95% EtOH concentration and keep specimens submerged).

218 The degraded DNA of older specimens preserved by less-than-ideal methods can,  
219 fortunately, be captured by massive-parallel methods successfully incorporate shorter, more  
220 degraded DNA fragments than can be used for sanger sequencing. One illuminating study  
221 generated nearly 1000 UCEs loci from pinned bee specimens up to 121 years old (Blaimer et al.  
222 2016a). This study demonstrated that pinned specimens less than 20 years old had significantly  
223 higher pre- and post-library concentrations, UCE contig lengths, and locus counts compared to  
224 older specimens. The small size (<5mm), and often corresponding low DNA yield, of many

225 arthropod specimens is another challenge to successful capture of genomic data, and a problem  
226 that may be exacerbated using non-destructive sampling to retain voucher specimens. Total yield  
227 of genetic material can be increased with the use of DNA amplification kits, albeit at a higher  
228 cost (Cruaud et al. 2018). UCE data has been successfully generated from minute, non-  
229 destructively sampled chalcidoid wasps (average DNA input = 25ng, Cruaud et al. 2019). This  
230 study used a modified protocol to maximize DNA yield from a commercially available  
231 extraction kit (Qiagen DNeasy Blood and Tissue Kit, Valencia, CA), by using LoBind tubes  
232 (Eppendorf) and heating the elution buffer for longer periods, while decreasing the number of  
233 purification steps. While no correlation between input DNA quantity and number of UCE loci  
234 captured in this dataset, a large amount of missing data ultimately resulted in a data matrix that  
235 was only 25% complete. Ultimately, in order to ensure high quality DNA generation, using fresh,  
236 well-preserved specimens preserved in 95% EtOH and stored in -80°C or -20°C, is  
237 recommended. Pinned specimens collected within the past 20 years are also suitable. Destructive  
238 sampling will likely generate higher DNA yield, but for rare specimens non-destructive soaking  
239 within lysis buffer will suffice.

240 Careful selection of tissue types can significantly lower the potential of contamination by  
241 non-target organisms. Precautions can be taken by decontaminating specimens using UV light, as  
242 well as separating areas used for DNA extraction from amplification areas (Yeates et al. 2016).  
243 Additional recommendations include removing appendages used by predators to capture prey  
244 (Bossert and Danforth 2018), targeting life stages, such as adults, that are less likely to host  
245 endoparasitoids. Contamination can also be reduced by using either strict bioinformatic  
246 processing parameters, or methods such as the *phyluce\_assembly\_match\_contigs\_to\_barcode*s

247 script in PHYLUCE which extracts the *COI* barcode region, which is used for validating the  
248 presence of a single or multiple species (Bossert and Danforth 2018).

249 *Library Preparation.* Once DNA has been extracted from target organisms, wet lab  
250 protocol for preparing the DNA libraries for sequencing varies little across taxa. Depending on  
251 the quality of DNA or level of DNA degradation, the extracted and quantified genomic DNA  
252 may need to be sheared using sonication or enzymatic digestion to reach the target size of 400–  
253 600bp. The degree of DNA degradation will determine the duration of sonication needed; this  
254 can be assessed using gel electrophoresis, or automated electrophoresis systems such as  
255 TapeStation or Bioanalyzer.

256 At this stage, UCE sample preparation consists of seven main steps: 1) DNA  
257 quantification; 2) adaptor ligation; 3) PCR amplification and initial pooling of specimens; 4)  
258 hybrid enrichment; 5) amplification of enriched libraries; 6) Quantification and final pooling;  
259 and 7) size selection and final quantification (detailed in (Branstetter et al. 2017a)).

260

## 261 **Bioinformatics**

262 Once sequencing is complete it is time to proceed to data analysis. Like other genomic  
263 datasets, one of the advantages of UCEs is the volume of data returned; managing datasets at this  
264 scale also presents a challenge to researchers new to genomics. Processing UCE data involves  
265 three principal steps: 1) demultiplexing, filtering, and trimming the raw Illumina reads; 2) contig  
266 assembly; and 3) UCE processing for phylogenomic analysis. Currently, the most widely-used  
267 bioinformatics pipeline for UCE data processing is PHYLUCE (Faircloth 2015), which includes  
268 a suite of Python wrapper scripts for these steps by calling other programs (detailed below) and

269 batch processing many samples at once. Additional bioinformatic programs not currently  
270 included within PHYLUCE can also be used to process data, as data can easily be imported back  
271 into the pipeline. Alternatively, the SECAPR (Andermann et al. 2018) pipeline also functions  
272 similarly to PHYLUCE and can be used for batch processing of UCE data, while MitoFinder  
273 (Allio et al. 2019) pipeline can be used to extract both UCE and mitogenomic data.

274 1) *Demultiplexing, filtering and trimming of raw Illumina reads.* Analyzing Illumina data always  
275 begins with batch trimming of adapters and low-quality bases of de-multiplexed data. In the  
276 PHYLUCE pipeline this is achieved using Illumiprocessor (Faircloth 2013), which is built  
277 around the Trimmomatic program (Bolger et al. 2014). Alternatively, external trimming  
278 programs such as Trim Galore! (<https://github.com/FelixKrueger/TrimGalore>) can be used  
279 instead of Illumiprocessor.

280 2) *Contig Assembly.* Currently PHYLUCE supports multiple programs such as velvet (Zerbino  
281 and Birney 2008), Trinity (Grabherr et al. 2011), ABySS (Simpson et al. 2009), and SPAdes  
282 (Bankevich et al. 2012) for genome assembly. While Trinity has been the most widely used of  
283 the assembly methods in published papers, updates to PHYLUCE are in the process of  
284 eliminating compatibility with Trinity due to technical issues. Both ABySS and velvet require an  
285 input for  $k$ -mer value, which is as part of the De Bruijn graph assembly algorithm. Smaller  $k$ -  
286 mers result in the assembly of shorter contigs with more connections, while large  $k$ -mers can  
287 result in longer but fewer contigs. However, it is difficult to determine the  $k$ -mer size for UCE  
288 data as the depth of coverage for each locus is variable due to capture efficiency. Therefore,  
289 testing multiple  $k$ -mer values is recommended, starting at the default of 35 and moving up to 55–  
290 65 to find the best trade-off in terms of contig size vs.  $k$ -mer number. Automatic estimation of  $k$ -  
291 mers is possible using SPAdes or the VelvetOptimiser wrapper script along with velvet.

292 Mitogenomic assemblers, such as MetaSPAdes (Nurk et al. 2017), are a promising alternative to  
293 currently used genomic and transcriptomic *de novo* assemblers (Allio et al. 2019). These new  
294 tools are designed to account for variance in sequencing coverage, and are thus capable of  
295 generating larger and more complete supermatrices in a fraction of the time required by Trinity.

296 3) *UCE processing for phylogenomics*. Once assembled, contigs must be processed to determine  
297 which ones represent enriched UCEs loci. Orthologs are identified by aligning the assembled  
298 contigs to a FASTA file of target enrichment baits, and paralogs are subsequently removed  
299 (Faircloth 2015). The output is then screened to identify 1) assembled contigs match by probes  
300 targeting different loci, and 2) different contigs match by probes targeting the same loci. The  
301 latter must be removed from downstream analysis because they will be identified as potentially  
302 paralogous genes by PHYLUCE (Faircloth 2015), which can be problematic if the probes are not  
303 well-designed (see Current and Future Challenges below). The resulting FASTA files are then  
304 aligned using MUSCLE (Edgar 2004) or MAFFT (Kato et al. 2002) within PHYLUCE,  
305 followed by trimming for data matrix completeness using GBlocks (Castresana 2000) or TrimAl  
306 (Capella-Gutiérrez et al. 2009). Finally, the completed data matrices can be exported in a variety  
307 of commonly used formats (e.g. phylip, nexus, etc.) for downstream phylogenomic analyses.

### 308 **Allelic phasing**

309 Allelic phasing is an additional, optional data processing step that extracts SNPs from  
310 UCE loci by separating (phasing) the heterozygous sites into two allele sequences; this approach  
311 can be used to increase resolution for shallow-level phylogenetic or species delimitation studies  
312 (Zarza et al. 2018, Andermann et al. 2019, Derkarabetian et al. 2019). Allelic phasing can be has  
313 been shown to provide more accurate estimation of tree topology and divergence times than  
314 using contig sequences, especially at shallow phylogenetic levels under multispecies coalescent

315 (MSC) models (Andermann et al. 2019), and can be performed in both PHYLUCÉ and  
316 SECAPR. This is in part due to common assembler programs not originally designed for  
317 heterozygous sequences or genomes, and as a result contig sequences generated by these  
318 programs will mask information by eliminating one of the two variants at a heterozygous site  
319 (Bodily et al. 2015). Another benefit of phasing the sequence doubles the sample size, as each  
320 diploid individual will have two strands of DNA sequences (Andermann et al. 2019). While this  
321 isn't always necessary for deep level phylogenomic studies, we recommend performing allelic  
322 phasing for UCE datasets intended for shallow-scale evolutionary studies, such as species  
323 delimitation or population genomics. However, sufficient sequence coverage is needed to ensure  
324 the quality of phased results, as contigs with lower coverage risk being phased inaccurately.

325

## 326 **Phylogenomic Analyses**

327 At this stage, data are nearly ready for use in phylogenetic reconstruction. Before  
328 beginning, however, it is advisable to perform inspection of sequence alignments for each gene,  
329 whether using programs such as GUIDANCE2 (Sela et al. 2015) or custom scripts, rather than  
330 labor-intensive inspection by eye. Preparation of raw data for tree building has become highly  
331 automated in response to the large volumes of data generated by high-throughput methods.  
332 Standardized sequence inspection helps reduce errors and inconsistencies, but can also be  
333 responsible for introducing errors in UCE datasets.

### 334 **Data Filtering**

335 Data filtering is a vital step in quality control of phylogenomic studies, as sequencing  
336 thousands of genes across many samples can lead to missing data in certain taxa. We advise

337 using different filtering criteria to generate multiple datasets and thereby find a balance between  
338 maintaining sequence quantity and quality. For example, GC bias has been demonstrated to be  
339 negatively correlated with topological support in bees (Bossert et al. 2017), and incongruences  
340 among analyses have been found to be exacerbated in studies of ants that used only “high signal”  
341 loci with highest average bootstrap (Borowiec 2019). While there is no current consensus on the  
342 best approach to filtering UCE data, multiple strategies exist. The program BaCoCa (Kuck and  
343 Struck 2014) can be used to filter out genes based on statistical properties such as saturation of  
344 nucleotides, compositional bias and heterogeneity, and proportion of shared missing data. The  
345 program Phylo-MCOA (de Vienne et al. 2012) can be used to detect outlier genes or species that  
346 cause topological incongruences; these can be subsequently filtered out for phylogenetic  
347 reconstruction. Another approach to data filtering is to retain only protein-coding genes rather  
348 than every UCE locus, some of which may include non-coding regions (see Current and Future  
349 Challenges below for more details). Also promising is analysis of protein-coding genes, which  
350 evolve under purifying selection and can be analyzed separately as amino acids; a custom script  
351 is now available for extracting putative protein-coding genes from UCE data (Borowiec 2019).

352

### 353 **Data Partitioning**

354 Approaches to partitioning UCE data can be divided into three strategies: 1) assign all  
355 UCE loci to a single partition; this assumes that every site in the alignment has evolved under a  
356 common evolutionary process; 2) assign each UCE locus to a separate partition; this allows for  
357 variation in rates and patterns of evolution between UCEs but assumes that all sites within each  
358 UCE locus have evolved under the same Markov process; or 3) *k*-means clustering of sites based  
359 on evolutionary rates (Frandsen et al. 2015), which subdivides data into partitions based on



360 evolutionary rates, thus avoiding *a priori* partitioning by the user. All three are used, however,  
361 recent studies have shown the *k*-means algorithm could be unreliable for UCE data, as it  
362 generates a partition comprised of all the invariant sites in the dataset, possibly misleading  
363 phylogenetic inference methods (Baca et al. 2017a). A promising new method for partitioning  
364 UCE data is the Sliding-Window Site Characteristics (SWSC, Tagliacollo and Lanfear 2018),  
365 which divides each UCE locus into three data blocks (right flank, core, and left flank) as the  
366 UCE core regions are conserved, while the two flanking regions become increasingly more  
367 variable (Faircloth et al. 2012). Different methods can be used by SWSC to evaluate sites, but  
368 the site entropies (EN), in particular, have been shown to most accurately account for within-  
369 UCE heterogeneity (Tagliacollo and Lanfear 2018). Using the SWSC-EN partitioning schemes  
370 account for within-UCE heterogeneity and leads to an increase in model fit (Tagliacollo and  
371 Lanfear 2018, Branstetter and Longino 2019).

372

### 373 **Tree Building**

374 Once datasets have been generated, downstream analyses on UCE data are similar to  
375 phylogenetic analyses performed on most other data types (e.g. Sanger sequencing, SNPs, etc.).  
376 A variety of tree-building methods (Figure 3) can be used for reconstructing phylogeny from  
377 UCE datasets, including Maximum Likelihood (ML), Bayesian Inference (BI), or Multispecies  
378 Coalescent/Species Tree (MSC). While the intricacies of phylogenetic analyses are beyond the  
379 scope of this paper, excellent and detailed overviews – both theoretical and practical – are  
380 available (Yang and Rannala 2012, Liu et al. 2015, Bromham et al. 2018, Bravo et al. 2019).

### 381 **Maximum Likelihood**

382           Maximum likelihood (ML) is a statistical methodology for estimating unknown  
383 parameters in a model. ML is widely used in phylogenetic studies due to its use of complex  
384 substitution models and its robustness to many violations to the assumptions of these models  
385 (Yang and Rannala 2012). The most widely used programs for phylogenetic reconstruction in the  
386 ML framework includes RAxML (Stamatakis 2006, Kozlov et al. 2018) and IQ-TREE (Nguyen  
387 et al. 2014). One advantage of these programs is their speed, with the former being the dominant  
388 method within UCE literature despite having very limited evolutionary model choices. IQ-TREE  
389 has gained momentum in recent years for its ability to produce accurate trees without sacrificing  
390 speed (Zhou et al. 2018). It includes functions such as ModelFinder for finding appropriate  
391 evolutionary models (Kalyaanamoorthy et al. 2017); approximation-based methods such as  
392 ultrafast bootstrap (UFBoot) and Shimodaira-Hasegawa like approximate likelihood ratio test  
393 (SH-aLRT), which greatly decreases computational time compared to traditional nonparametric  
394 bootstrap methods (Guindon et al. 2010, Hoang et al. 2017); and gene/site concordance factors as  
395 alternative support measures to illustrate disagreement among loci and sites (Minh et al. 2018).

### 396 **Bayesian Inference and Divergence Dating**

397           Like ML, Bayesian inference (BI) is also a general methodology of statistical inference  
398 that has been widely adopted for phylogenetic analyses. Bayesian inference differs from ML in  
399 that parameters in the models are considered to be random variables within statistical  
400 distributions rather than unknown fixed constants (Yang and Rannala 2012). Today BI using the  
401 Markov chain Monte Carlo (MCMC) sampling is a widely adopted method used for  
402 phylogenetic analysis, as the incorporation of prior knowledge into the analysis offers an  
403 appealing alternative to ML even at the cost of slower computational speed (Nascimento et al.  
404 2017). The commonly used Bayesian programs for phylogenomics data include BEAST

405 (Drummond and Rambaut 2007), BEAST 2 (Bouckaert et al. 2014), and ExaBayes (Aberer et al.  
406 2014). Bayesian analyses are extremely sensitive to prior probabilities set by users, and often  
407 default priors may not be appropriate for the data being analyzed as they can affect resulting  
408 topologies (Nascimento et al. 2017). Because setting priors can be daunting for beginners, we  
409 advise users to resist the all-too-common tendency to employ default settings and instead urge  
410 users to follow steps outlined in Bromham et al. (2018) to make informed choices when setting  
411 up Bayesian analyses. Running an ‘empty’ analysis without data to allow MCMC algorithm  
412 sampling from the prior is a good way of checking whether the data were informative enough to  
413 return posterior distributions different from the marginal priors, and to assess for good  
414 convergence and mixing of the MCMC chains (Nascimento et al. 2017, Blaimer et al. 2018b).

415 Divergence time estimation analyses can also be implemented for UCE phylogenies to  
416 generate dated chronograms, using carefully selected fossils as calibration points. The commonly  
417 used node-dating approach assigns the oldest fossil that can be confidently identified to the  
418 youngest internal node, imposing the age of the fossil a minimum age constraint (Arcila et al.  
419 2015). An alternative method called total-evidence, or tip-dating methods can include all  
420 available paleontological information, ameliorating fossil-placement uncertainty while  
421 simultaneously incorporating fossil ages into the analysis (Ronquist et al. 2012a). Both of these  
422 methods can be implemented for divergence date estimation using Bayesian inference programs  
423 such as MrBayes (Ronquist et al. 2012b), BEAST/BEAST2, and the MCMCTree package in  
424 PAML (Yang 2007). While MCMCTree is faster computationally, the setup for prior  
425 distributions on fossil calibrations is less intuitive. BEAST/BEAST2, by comparison, is easier to  
426 understand and offers more analytical options such as the incorporation of fossils directly into  
427 the phylogeny with the newly developed node-dating method using the fossilized birth-death

428 model (Heath et al. 2014). The fossilized birth-death model offers an advantage over other  
429 methods by combining morphological and molecular data as well as stratigraphic range data  
430 from the fossil record, and can be implemented directly in RevBayes (Höhna et al. 2016), or in  
431 BEAST2 with add on package sampled-ancestor (Gavryushkina et al. 2014). In general, large  
432 data volumes associated with UCEs makes most Bayesian analyses too computationally  
433 intensive to be practical. To overcome this limitation, many studies reduce data size by removing  
434 taxa or loci in order to reduce the analysis time (Blaimer et al. 2018b, Borowiec 2019). It is also  
435 worth noting that tip-dating models have been shown to recover older ages than traditional node-  
436 dating models, and might produce inaccurate date estimations (Arcila et al. 2015). Regardless of  
437 the approach, the resulting dated chronogram can be then used as input for additional analyses  
438 such as ancestral state reconstruction, historical biogeographic analysis, or diversification rates  
439 estimation.

#### 440 **Multispecies Coalescent/Species Tree**

441 One key advance in molecular phylogenetics has been the acknowledgement that high  
442 levels of incomplete lineage sorting (ILS) or other stochastic errors can yield misleading results  
443 for traditional methods concatenation methods (Liu et al. 2015, Bravo et al. 2019). Incorporation  
444 of discordance between gene trees and species trees as a result of high incomplete lineage sorting  
445 (ILS), under the MSC model (Heled and Drummond 2009) can alleviate this problem.  
446 Commonly used MSC tree summary-based methods such as ASTRAL (Mirarab et al. 2014,  
447 Mirarab and Warnow 2015, Zhang et al. 2018) and MP-EST (Liu et al. 2010) are performed in  
448 two steps, wherein gene trees are estimated first and separately, then used as input to generate a  
449 species tree based on various summaries of coalescent process (Bravo et al. 2019). Because the  
450 accuracy of the individual input gene trees directly affects the resulting species tree, these

451 summary-based methods are especially susceptible to gene tree estimation errors (Molloy and  
452 Warnow 2018). Therefore, checking individual gene trees for incongruences is advised to ensure  
453 species tree accuracy in summary-based methods. Methods such as concordance factors (Ané et  
454 al. 2006, Minh et al. 2018) should be used to provide insight into the influence of ILS versus  
455 other factors such as introgression on the resulting topology. Alternatively, site-based coalescent  
456 method such as SVDquartets (Chifman and Kubatko 2014) and SVDquest (Vachaspati and  
457 Warnow 2018) bypass gene tree estimation, and is comparable or even more accurate than  
458 summary methods in cases of high ILS (Chou et al. 2015, Molloy and Warnow 2018). Finally,  
459 the newly developed StarBEAST2 (Ogilvie et al. 2017) package for BEAST2 is a promising  
460 implementation of the full MSC model which can jointly infer gene trees and species trees, but  
461 the current version is too computationally intensive to use on large UCE datasets.

462

### 463 **Resources and Costs**

464 Most steps of the wet lab protocol can be performed in standard molecular labs that have  
465 access to equipment such as a centrifuge and thermocycler. More specialized equipment such as  
466 a sonicator for shearing DNA, TapeStation/Bioanalyzer for quantifying DNA, and  
467 BluePippin/PippinHT for size selection can all be substituted with cheaper, albeit less accurate  
468 alternatives such as restriction enzymes, gel electrophoresis, and magnetic beads.

469 Illumina platforms (HiSeq, NextSeq, NovaSeq) are generally used for UCE studies due to  
470 their high throughput and low cost per base pair. The current estimated cost per specimen is  
471 approximately \$30 – 40 USD, accounting for costs of all reagents in library preparation and  
472 paired-end Illumina run (See Supp Table 1 for sample cost breakdown). Some commercial

473 laboratories (e.g. RAPID Genomics, Gainesville, FL, USA) also offer UCE enrichment services,  
474 handling all library preparation, enrichment, and sequencing; customers simply submit DNA  
475 extracts and then receive sequence data. Costs associated with such ‘concierge service’ are  
476 considerably higher (approximately ~\$120 per specimen), but this may be an attractive option for  
477 researchers lacking the infrastructure or personnel to undertake wet lab protocols.

478         Having access to high performance computing (HPC) greatly expedites bioinformatic and  
479 phylogenomic analyses, especially when processing large batches of samples. While PHYLUCE  
480 and many associated data-processing programs can be run in local Linux/Unix environment, the  
481 parallelization using HPC will reduce execution time in computationally intensive steps such as  
482 demultiplexing and assembly. Similarly, many phylogenomics programs discussed above can  
483 also be expedited through this process.

484

### 485 **Data Availability Recommendations**

486         One hallmark feature of UCE data is its open source nature, probe sets, protocols, and  
487 previously published data are made publicly available, ensuring repeatability – the foundation of  
488 open scientific research. To this end, untrimmed raw Illumina reads should be uploaded to public  
489 database such as Sequence Read Archive (SRA) once studies are published, giving interested  
490 readers the full ability to download and process the data using different trimming settings. All  
491 analytical methods such as software and code used to process data should also be made publicly  
492 available on repositories such as Dryad or GitHub. UCE contigs can be uploaded to GenBank as  
493 targeted locus studies, making the data available for BLAST. The voucher specimens from  
494 which DNA was extracted should be deposited in recognized scientific collections and museums;

495 associated information such as collection locality, identification, etc., should be included as  
496 metadata with all molecular sequence (Bravo et al. 2019).

497

## 498 **Current and Future Challenges**

499 UCE and similar methods offer the ability to generate massive amounts of data from  
500 many loci, and yet, despite the increase in data volume, the same concerns that have long  
501 plagued phylogenetic analyses remain as relevant as ever: taxon sampling, choice of alignment  
502 methods, and composition bias (Bossert et al. 2017, Mclean et al. 2018). Recent research also  
503 suggests phylogenomic results can be strongly affected by a tiny proportion of highly biased loci  
504 or sites (Shen et al. 2017), and reduction of phylogenetic noise resulting from compositional  
505 heterogeneity and saturation can increase congruence among different analytic methods  
506 (Borowiec 2019). With that in mind, we strongly encourage performing sensitivity analyses to  
507 test the robustness of results when interpreting phylogenomic data (Borowiec 2019, Camacho et  
508 al. 2019). As these large datasets are less prone to uncertainty, and instead may give strongly  
509 supported wrong results if model violations are not carefully evaluated (Borowiec et al. 2019).

510 The fact that the function of UCEs remains largely unknown is the basis of active  
511 research and a current challenge for identifying and modeling UCEs in a phylogenomic context  
512 (Bejerano et al. 2004). Vertebrate UCEs are characterized as predominantly non-coding  
513 sequences, non-randomly distributed across chromosomes and acting as regulators and/or  
514 enhancers of gene expression (Baira et al. 2008, Polychronopoulos et al. 2017). By contrast,  
515 studies of invertebrate UCEs reveal that most flanking regions captured include exons  
516 (Branstetter et al. 2017a), with the most widely shared loci being either exclusively conserved

517 exons or partially exonic regions in Hymenoptera and Arachnida (Bossert and Danforth 2018,  
518 Hedin et al. 2019). This is an exciting discovery, as the exonic flanking regions captured by the  
519 UCE process and transcriptome sequence data within these groups can be meaningfully  
520 combined, without the need to design specific probe sets to target them, as demonstrated in  
521 Apidae (Bossert et al. 2019). However, since the genomic landscapes of different animal taxa  
522 can differ substantially, the wider application of combining transcriptomic data with UCEs in  
523 other taxonomic groups still needs to be tested. Currently, it appears that the function of UCEs is  
524 highly variable, with flanking regions containing exons and introns; whether this variability will  
525 affect downstream analyses remains to be seen.

526         Continued refinement of existing probe sets is needed to increase capture success while  
527 minimizing duplicates and paralogous loci. It has been shown in the arachnid probe set, different  
528 UCE probes sometimes target regions of the same protein, or include non-orthologous sequences  
529 (Hedin et al. 2019). This is unsurprising given the wide phylogenetic depth of the probe set,  
530 which was designed to target all arachnids, but given that PHYLUCE cannot detect these non-  
531 orthologous sequences as the program only removes different contigs hit by probes targeting the  
532 same loci (Faircloth 2015), additional manual filtering is needed to ensure the exclusion of  
533 misleading paralogous sequences into the final data matrix (Hedin et al. 2019).

534

## 535 **Conclusion**

536         Ultraconserved elements-based phylogenomic studies have been rapidly adopted by  
537 researchers working on arthropod taxa since their introduction by Faircloth et al. (2012). This  
538 review described the versatility of UCE data at both deep and shallow evolutionary scale, and



539 provided a step-by-step guide to generating and analyzing UCEs; we then summarized current  
540 practices, challenges, and unresolved questions that surround this active field. Our hope is to  
541 make UCE-based phylogenomic studies more accessible to users with diverse taxonomic  
542 interests, and thereby deepen our collective understanding of the roles and functions of UCEs  
543 across widely divergent taxa. As our understanding of UCEs develops through studies of  
544 different organisms, identifying individual genes and incorporation of functional genomics will  
545 yield interesting comparative studies across deeply divergent taxonomic groups and provide new  
546 insights in the continued pursuit of building the tree of life.

547

548

## 549 **Acknowledgements**

550 We would like to thank Michael Branstetter and two anonymous reviewers for comments  
551 that greatly improved the manuscript, Eliana Buenaventura, Chris Cohen, Shahan Derkarabetian,  
552 Grey Gustafson, Katherine Noble, and Matthew van Dam for informative discussions on UCE  
553 research in their perspective groups. We thank Rachel Atchison and Suzy Rodriguez for graphic  
554 design assistance.

555

## 556 **References**

557 **Aberer, A. J., K. Kobert, and A. Stamatakis. 2014.** ExaBayes: massively parallel Bayesian  
558 tree inference for the whole-genome era. *Molecular Biology and Evolution* 31: 2553-2556.  
559 **Alda, F., V. A. Tagliacollo, M. J. Bernt, B. T. Waltz, W. B. Ludt, B. C. Faircloth, M. E.**  
560 **Alfaro, J. S. Albert, and P. Chakrabarty. 2018.** Resolving deep nodes in an ancient radiation

- 561 of neotropical fishes in the presence of conflicting signals from incomplete lineage sorting.  
562 *Systematic Biology* 68: 573-593.
- 563 **Ali, O. A., S. M. O'Rourke, S. J. Amish, M. H. Meek, G. Luikart, C. Jeffres, and M. R.**  
564 **Miller. 2016.** RAD Capture (Rapture): Flexible and efficient sequence-based genotyping.  
565 *Genetics* 202: 389-400.
- 566 **Allio, R., A. Schomaker-Bastos, J. Romiguier, F. Prosdocimi, B. Nabholz, and F. Delsuc.**  
567 **2019.** MitoFinder: efficient automated large-scale extraction of mitogenomic data in target  
568 enrichment phylogenomics. *bioRxiv* 685412.
- 569 **Andermann, T., A. Cano, A. Zizka, C. Bacon, and A. Antonelli. 2018.** SECAPR-a  
570 bioinformatics pipeline for the rapid and user-friendly processing of targeted enriched Illumina  
571 sequences, from raw reads to alignments. *PeerJ* 6: e5175.
- 572 **Andermann, T., A. M. Fernandes, U. Olsson, M. Topel, B. Pfeil, B. Oxelman, A. Aleixo, B.**  
573 **C. Faircloth, and A. Antonelli. 2019.** Allele phasing greatly improves the phylogenetic utility  
574 of ultraconserved elements. *Systematic Biology* 68: 32-46.
- 575 **Ané, C., B. Larget, D. A. Baum, S. D. Smith, and A. Rokas. 2006.** Bayesian Estimation of  
576 Concordance among Gene Trees. *Molecular Biology and Evolution* 24: 412-426.
- 577 **Arcila, D., R. Alexander Pyron, J. C. Tyler, G. Orti, and R. R. Betancur. 2015.** An  
578 evaluation of fossil tip-dating versus node-age calibrations in tetraodontiform fishes (Teleostei:  
579 Percomorphaceae). *Molecular Phylogenetics and Evolution* 82 Pt A: 131-145.
- 580 **Baca, S. M., E. F. A. Toussaint, K. B. Miller, and A. E. Z. Short. 2017a.** Molecular  
581 phylogeny of the aquatic beetle family Noteridae (Coleoptera: Adephaga) with an emphasis on  
582 data partitioning strategies. *Molecular Phylogenetics and Evolution* 107: 282-292.
- 583 **Baca, S. M., A. Alexander, G. T. Gustafson, and A. E. Z. Short. 2017b.** Ultraconserved  
584 elements show utility in phylogenetic inference of Adephaga (Coleoptera) and suggest paraphyly  
585 of 'Hydradephaga'. *Systematic Entomology* 42: 786-795.
- 586 **Baira, E., J. Greshock, G. Coukos, and L. Zhang. 2008.** Ultraconserved elements: genomics,  
587 function and disease. *RNA Biology* 5: 132-134.
- 588 **Baird, N. A., P. D. Etter, T. S. Atwood, M. C. Currey, A. L. Shiver, Z. A. Lewis, E. U.**  
589 **Selker, W. A. Cresko, and E. A. Johnson. 2008.** Rapid SNP discovery and genetic mapping  
590 using sequenced RAD markers. *PLoS One* 3: e3376.
- 591 **Bankevich, A., S. Nurk, D. Antipov, A. A. Gurevich, M. Dvorkin, A. S. Kulikov, V. M.**  
592 **Lesin, S. I. Nikolenko, S. Pham, and A. D. Prjibelski. 2012.** SPAdes: a new genome assembly  
593 algorithm and its applications to single-cell sequencing. *Journal of Computational Biology* 19:  
594 455-477.
- 595 **Bejerano, G., M. Pheasant, I. Makunin, S. Stephen, W. J. Kent, J. S. Mattick, and D.**  
596 **Haussler. 2004.** Ultraconserved elements in the human genome. *Science* 304: 1321-1325.
- 597 **Bi, K., D. Vanderpool, S. Singhal, T. Linderoth, C. Moritz, and J. M. Good. 2012.**  
598 Transcriptome-based exon capture enables highly cost-effective comparative genomic data  
599 collection at moderate evolutionary scales. *BMC Genomics* 13: 403.
- 600 **Blaimer, B. B., J. R. Mawdsley, and S. G. Brady. 2018a.** Multiple origins of sexual  
601 dichromatism and aposematism within large carpenter bees. *Evolution* 72: 1874-1889.
- 602 **Blaimer, B. B., M. W. Lloyd, W. X. Guillory, and S. G. Brady. 2016a.** Sequence capture and  
603 phylogenetic utility of genomic ultraconserved elements obtained from pinned insect specimens.  
604 *PLoS One* 11: e0161531.

- 605 **Blaimer, B. B., J. S. LaPolla, M. G. Branstetter, M. W. Lloyd, and S. G. Brady. 2016b.**  
606 Phylogenomics, biogeography and diversification of obligate mealybug-tending ants in the genus  
607 *Acropyga*. *Molecular Phylogenetics and Evolution* 102: 20-29.
- 608 **Blaimer, B. B., P. S. Ward, T. R. Schultz, B. L. Fisher, and S. G. Brady. 2018b.**  
609 Paleotropical diversification dominates the evolution of the hyperdiverse ant tribe  
610 Crematogastrini (Hymenoptera: Formicidae). *Insect Systematics and Diversity* 2: 3.
- 611 **Blaimer, B. B., S. G. Brady, T. R. Schultz, M. W. Lloyd, B. L. Fisher, and P. S. Ward. 2015.**  
612 Phylogenomic methods outperform traditional multi-locus approaches in resolving deep  
613 evolutionary history: a case study of formicine ants. *BMC Evolutionary Biology* 15: 271.
- 614 **Bodily, P. M., M. S. Fujimoto, C. Ortega, N. Okuda, J. C. Price, M. J. Clement, and Q.**  
615 **Snell. 2015.** Heterozygous genome assembly via binary classification of homologous sequence.  
616 *BMC Bioinformatics* 16: S5.
- 617 **Bolger, A. M., M. Lohse, and B. Usadel. 2014.** Trimmomatic: a flexible trimmer for Illumina  
618 sequence data. *Bioinformatics* 30: 2114-2120.
- 619 **Borowiec, M. L. 2019.** Convergent evolution of the army ant syndrome and congruence in big-  
620 data phylogenetics. *Systematic Biology* 68: 642-656.
- 621 **Borowiec, M. L., C. Rabeling, S. G. Brady, B. L. Fisher, T. R. Schultz, and P. S. Ward.**  
622 **2019.** Compositional heterogeneity and outgroup choice influence the internal phylogeny of the  
623 ants. *Molecular Phylogenetics and Evolution* 134: 111-121.
- 624 **Bossert, S., and B. N. Danforth. 2018.** On the universality of target-enrichment baits for  
625 phylogenomic research. *Methods in Ecology and Evolution* 9: 1453-1460.
- 626 **Bossert, S., E. A. Murray, B. B. Blaimer, and B. N. Danforth. 2017.** The impact of GC bias  
627 on phylogenetic accuracy using targeted enrichment phylogenomic data. *Molecular*  
628 *Phylogenetics and Evolution* 111: 149-157.
- 629 **Bossert, S., E. A. Murray, E. A. B. Almeida, S. G. Brady, B. B. Blaimer, and B. N.**  
630 **Danforth. 2019.** Combining transcriptomes and ultraconserved elements to illuminate the  
631 phylogeny of Apidae. *Molecular Phylogenetics and Evolution* 130: 121-131.
- 632 **Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard, A.**  
633 **Rambaut, and A. J. Drummond. 2014.** BEAST 2: a software platform for Bayesian  
634 evolutionary analysis. *PLoS Computational Biology* 10: e1003537.
- 635 **Branstetter, M. G., and J. T. Longino. 2019.** Ultra-conserved element phylogenomics of New  
636 World *Ponera* (Hymenoptera: Formicidae) illuminates the origin and phylogeographic history of  
637 the Endemic Exotic Ant *Ponera exotica*. *Insect Systematics and Diversity* 3: 1.
- 638 **Branstetter, M. G., J. T. Longino, P. S. Ward, B. C. Faircloth, and S. Price. 2017a.**  
639 Enriching the ant tree of life: enhanced UCE bait set for genome-scale phylogenetics of ants and  
640 other Hymenoptera. *Methods in Ecology and Evolution* 8: 768-776.
- 641 **Branstetter, M. G., A. Jesovnik, J. Sosa-Calvo, M. W. Lloyd, B. C. Faircloth, S. G. Brady,**  
642 **and T. R. Schultz. 2017b.** Dry habitats were crucibles of domestication in the evolution of  
643 agriculture in ants. *Proceedings of the Royal Society of London B: Biological Sciences* 284.
- 644 **Branstetter, M. G., B. N. Danforth, J. P. Pitts, B. C. Faircloth, P. S. Ward, M. L.**  
645 **Buffington, M. W. Gates, R. R. Kula, and S. G. Brady. 2017c.** Phylogenomic Insights into the  
646 Evolution of Stinging Wasps and the Origins of Ants and Bees. *Current Biology* 27: 1019-1025.
- 647 **Bravo, G. A., A. Antonelli, C. D. Bacon, K. Bartoszek, M. P. K. Blom, S. Huynh, G. Jones,**  
648 **L. L. Knowles, S. Lamichhaney, T. Marcussen, H. Morlon, L. K. Nakhleh, B. Oxelman, B.**  
649 **Pfeil, A. Schliep, N. Wahlberg, F. P. Werneck, J. Wiedenhoeft, S. Willows-Munro, and S.**

- 650 **V. Edwards. 2019.** Embracing heterogeneity: coalescing the Tree of Life and the future of  
651 phylogenomics. *PeerJ* 7: e6399.
- 652 **Bromham, L., S. Duchene, X. Hua, A. M. Ritchie, D. A. Duchene, and S. Y. W. Ho. 2018.**  
653 Bayesian molecular dating: opening up the black box. *Biological Reviews* 93: 1165-1191.
- 654 **Camacho, G. P., M. R. Pie, R. M. Feitosa, and M. S. Barbeitos. 2019.** Exploring gene tree  
655 incongruence at the origin of ants and bees (Hymenoptera). *Zoologica Scripta* 48: 215-225.
- 656 **Capella-Gutiérrez, S., J. M. Silla-Martínez, and T. Gabaldón. 2009.** trimAl: a tool for  
657 automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25: 1972-  
658 1973.
- 659 **Castresana, J. 2000.** Selection of conserved blocks from multiple alignments for their use in  
660 phylogenetic analysis. *Molecular Biology and Evolution* 17: 540-552.
- 661 **Chifman, J., and L. Kubatko. 2014.** Quartet inference from SNP data under the coalescent  
662 model. *Bioinformatics* 30: 3317-3324.
- 663 **Chou, J., A. Gupta, S. Yaduvanshi, R. Davidson, M. Nute, S. Mirarab, and T. Warnow.**  
664 **2015.** A comparative study of SVDquartets and other coalescent-based species tree estimation  
665 methods. *BMC Genomics* 16: S2.
- 666 **Cooke, C. 2018.** Forest Micro-Hymenoptera, including those attacking trees (Cynipidae: Oak  
667 Gall Wasps) and those potentially defending them (Parasitic Pteromalidae). PhD, University of  
668 Maryland, College Park, MD, USA.
- 669 **Crawford, N. G., B. C. Faircloth, J. E. McCormack, R. T. Brumfield, K. Winker, and T. C.**  
670 **Glenn. 2012.** More than 1000 ultraconserved elements provide evidence that turtles are the sister  
671 group of archosaurs. *Biology Letters* rsbl20120331.
- 672 **Cruaud, A., G. Groussier, G. Genson, L. Saune, A. Polaszek, and J. Y. Rasplus. 2018.**  
673 Pushing the limits of whole genome amplification: successful sequencing of RADseq library  
674 from a single microhymenopteran (Chalcidoidea, *Trichogramma*). *PeerJ* 6: e5640.
- 675 **Cruaud, A., S. Nidelet, P. Arnal, A. Weber, L. Fusu, A. Gumovsky, J. Huber, A. Polaszek,**  
676 **and J. Y. Rasplus. 2019.** Optimized DNA extraction and library preparation for minute  
677 arthropods: Application to target enrichment in chalcid wasps used for biocontrol. *Molecular*  
678 *Ecology Resources*.
- 679 **de Vienne, D. M., S. Ollier, and G. Aguilera. 2012.** Phylo-MCOA: a fast and efficient method  
680 to detect outlier genes and species in phylogenomics using multiple co-inertia analysis.  
681 *Molecular Biology and Evolution* 29: 1587-1598.
- 682 **Derkarabetian, S., S. Castillo, P. K. Koo, S. Ovchinnikov, and M. Hedin. 2019.** A  
683 demonstration of unsupervised machine learning in species delimitation. *Molecular*  
684 *Phylogenetics and Evolution* 139.
- 685 **Derkarabetian, S., J. Starrett, N. Tsurusaki, D. Ubick, S. Castillo, and M. Hedin. 2018.** A  
686 stable phylogenomic classification of Travunioidea (Arachnida, Opiliones, Laniatores) based on  
687 sequence capture of ultraconserved elements. *ZooKeys* 1-36.
- 688 **Drummond, A. J., and A. Rambaut. 2007.** BEAST: Bayesian evolutionary analysis by  
689 sampling trees. *BMC Evolutionary Biology* 7: 214.
- 690 **Edgar, R. C. 2004.** MUSCLE: multiple sequence alignment with high accuracy and high  
691 throughput. *Nucleic Acids Research* 32: 1792-1797.
- 692 **Faircloth, B. 2013.** Illumiprocessor: a trimmomatic wrapper for parallel adapter and quality  
693 trimming.
- 694 **Faircloth, B. C. 2015.** PHYLUCES is a software package for the analysis of conserved genomic  
695 loci. *Bioinformatics* 32: 786-788.

- 696 **Faircloth, B. C. 2017.** Identifying conserved genomic elements and designing universal bait sets  
697 to enrich them. *Methods in Ecology and Evolution* 8: 1103-1112.
- 698 **Faircloth, B. C., L. Sorenson, F. Santini, and M. E. Alfaro. 2013.** A phylogenomic  
699 perspective on the radiation of ray-finned fishes based upon targeted sequencing of  
700 ultraconserved elements (UCEs). *PLoS One* 8: e65923.
- 701 **Faircloth, B. C., M. G. Branstetter, N. D. White, and S. G. Brady. 2015.** Target enrichment  
702 of ultraconserved elements from arthropods provides a genomic perspective on relationships  
703 among Hymenoptera. *Molecular Ecology Resources* 15: 489-501.
- 704 **Faircloth, B. C., J. E. McCormack, N. G. Crawford, M. G. Harvey, R. T. Brumfield, and T.  
705 C. Glenn. 2012.** Ultraconserved elements anchor thousands of genetic markers spanning  
706 multiple evolutionary timescales. *Systematic Biology* 61: 717-726.
- 707 **Forthman, M., C. W. Miller, and R. T. Kimball. 2019.** Phylogenomic analysis suggests  
708 Coreidae and Alydidae (Hemiptera: Heteroptera) are not monophyletic. *Zoologica Scripta* 48:  
709 520-534.
- 710 **Frandsen, P. B., B. Calcott, C. Mayer, and R. Lanfear. 2015.** Automatic selection of  
711 partitioning schemes for phylogenetic analyses using iterative k-means clustering of site rates.  
712 *BMC Evolutionary Biology* 15: 13.
- 713 **Gavryushkina, A., D. Welch, T. Stadler, and A. J. Drummond. 2014.** Bayesian inference of  
714 sampled ancestor trees for epidemiology and fossil calibration. *PLoS Computational Biology* 10:  
715 e1003919.
- 716 **Gnrke, A., A. Melnikov, J. Maguire, P. Rogov, E. M. LeProust, W. Brockman, T. Fennell,  
717 G. Giannoukos, S. Fisher, and C. Russ. 2009.** Solution hybrid selection with ultra-long  
718 oligonucleotides for massively parallel targeted sequencing. *Nature Biotechnology* 27: 182.
- 719 **Grab, H., M. G. Branstetter, N. Amon, K. R. Urban-Mead, M. G. Park, J. Gibbs, E. J.  
720 Blitzer, K. Poveda, G. Loeb, and B. N. Danforth. 2019.** Agriculturally dominated landscapes  
721 reduce bee phylogenetic diversity and pollination services. *Science* 363: 282-284.
- 722 **Grabherr, M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson, I. Amit, X.  
723 Adiconis, L. Fan, R. Raychowdhury, and Q. Zeng. 2011.** Full-length transcriptome assembly  
724 from RNA-Seq data without a reference genome. *Nature Biotechnology* 29: 644.
- 725 **Guindon, S., J.-F. Dufayard, V. Lefort, M. Anisimova, W. Hordijk, and O. Gascuel. 2010.**  
726 New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the  
727 performance of PhyML 3.0. *Systematic Biology* 59: 307-321.
- 728 **Gustafson, G. T., A. Alexander, J. S. Sproul, J. M. Pflug, D. R. Maddison, and A. E. Z.  
729 Short. 2019.** Ultraconserved element (UCE) probe set design: Base genome and initial design  
730 parameters critical for optimization. *Ecology and Evolution* 9: 6933-6948.
- 731 **Harvey, M. G., B. T. Smith, T. C. Glenn, B. C. Faircloth, and R. T. Brumfield. 2016.**  
732 Sequence capture versus restriction site associated DNA sequencing for shallow systematics.  
733 *Systematic Biology* 65: 910-924.
- 734 **Heath, T. A., J. P. Huelsenbeck, and T. Stadler. 2014.** The fossilized birth-death process for  
735 coherent calibration of divergence-time estimates. *Proceedings of the National Academy of  
736 Sciences of the United States of America* 111: E2957-2966.
- 737 **Hedin, M., S. Derkarabetian, J. Blair, and P. Paquin. 2018a.** Sequence capture  
738 phylogenomics of eyeless *Cicurina* spiders from Texas caves, with emphasis on US federally-  
739 endangered species from Bexar County (Araneae, Hahniidae). *ZooKeys* 49-76.

- 740 **Hedin, M., S. Derkarabetian, M. J. Ramírez, C. Vink, and J. E. Bond. 2018b.** Phylogenomic  
741 reclassification of the world's most venomous spiders (Mygalomorphae, Atracinae), with  
742 implications for venom evolution. *Scientific Reports* 8: 1636.
- 743 **Hedin, M., S. Derkarabetian, A. Alfaro, M. J. Ramírez, and J. E. Bond. 2019.** Phylogenomic  
744 analysis and revised classification of atypoid mygalomorph spiders (Araneae, Mygalomorphae),  
745 with notes on arachnid ultraconserved element loci. *PeerJ* 7: e6864.
- 746 **Heled, J., and A. J. Drummond. 2009.** Bayesian inference of species trees from multilocus  
747 data. *Molecular Biology and Evolution* 27: 570-580.
- 748 **Hoang, D. T., O. Chernomor, A. von Haeseler, B. Q. Minh, and S. V. Le. 2017.** UFBoot2:  
749 Improving the ultrafast bootstrap approximation. *Molecular Biology and Evolution* msx281.
- 750 **Hoffberg, S. L., T. J. Kieran, J. M. Catchen, A. Devault, B. C. Faircloth, R. Mauricio, and**  
751 **T. C. Glenn. 2016.** RAD cap: sequence capture of dual-digest RAD seq libraries with  
752 identifiable duplicates and reduced missing data. *Molecular Ecology Resources* 16: 1264-1278.
- 753 **Höhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P.**  
754 **Huelsbeck, and F. Ronquist. 2016.** RevBayes: Bayesian phylogenetic inference using  
755 graphical models and an interactive model-specification language. *Systematic Biology* 65: 726-  
756 736.
- 757 **Ješovnik, A. N. A., J. Sosa-Calvo, M. W. Lloyd, M. G. Branstetter, F. FernÁNdez, and T.**  
758 **R. Schultz. 2017.** Phylogenomic species delimitation and host-symbiont coevolution in the  
759 fungus-farming ant genus *Sericomyrmex* Mayr (Hymenoptera: Formicidae): ultraconserved  
760 elements (UCEs) resolve a recent radiation. *Systematic Entomology* 42: 523-542.
- 761 **Jones, M. R., and J. M. Good. 2016.** Targeted capture in evolutionary and ecological genomics.  
762 *Molecular Ecology* 25: 185-202.
- 763 **Kalyaanamoorthy, S., B. Q. Minh, T. K. Wong, A. von Haeseler, and L. S. Jermiin. 2017.**  
764 **ModelFinder: fast model selection for accurate phylogenetic estimates.** *Nature Methods* 14: 587.
- 765 **Katoh, K., K. Misawa, K. i. Kuma, and T. Miyata. 2002.** MAFFT: a novel method for rapid  
766 multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Research* 30: 3059-  
767 3066.
- 768 **Kieran, T. J., E. R. Gordon, M. Forthman, R. Hoey-Chamberlain, R. T. Kimball, B. C.**  
769 **Faircloth, C. Weirauch, and T. C. Glenn. 2019.** Insight from an ultraconserved element bait  
770 set designed for hemipteran phylogenetics integrated with genomic resources. *Molecular*  
771 *Phylogenetics and Evolution* 130: 297-303.
- 772 **Kozlov, A., D. Darriba, T. Flouri, B. Morel, and A. Stamatakis. 2018.** RAxML-NG: A fast,  
773 scalable, and user-friendly tool for maximum likelihood phylogenetic inference. *bioRxiv*  
774 447110.
- 775 **Kuck, P., and T. H. Struck. 2014.** BaCoCa--a heuristic software tool for the parallel assessment  
776 of sequence biases in hundreds of gene and taxon partitions. *Molecular Phylogenetics and*  
777 *Evolution* 70: 94-98.
- 778 **Lemmon, A. R., S. A. Emme, and E. M. Lemmon. 2012.** Anchored hybrid enrichment for  
779 massively high-throughput phylogenomics. *Systematic Biology* 61: 727-744.
- 780 **Lemmon, E. M., and A. R. Lemmon. 2013.** High-throughput genomic data in systematics and  
781 phylogenetics. *Annual Review of Ecology, Evolution, and Systematics* 44: 99-121.
- 782 **Lim, H. C., and M. J. Braun. 2016.** High-throughput SNP genotyping of historical and modern  
783 samples of five bird species via sequence capture of ultraconserved elements. *Molecular Ecology*  
784 *Resources* 16: 1204-1223.

- 785 **Liu, L., L. Yu, and S. V. Edwards. 2010.** A maximum pseudo-likelihood approach for  
786 estimating species trees under the coalescent model. *BMC Evolutionary Biology* 10: 302.
- 787 **Liu, L., Z. Xi, S. Wu, C. C. Davis, and S. V. Edwards. 2015.** Estimating phylogenetic trees  
788 from genome-scale data. *Annals of the New York Academy of Sciences* 1360: 36-53.
- 789 **Mamanova, L., A. J. Coffey, C. E. Scott, I. Kozarewa, E. H. Turner, A. Kumar, E. Howard,  
790 J. Shendure, and D. J. Turner. 2010.** Target-enrichment strategies for next-generation  
791 sequencing. *Nature Methods* 7: 111.
- 792 **Manthey, J. D., L. C. Campillo, K. J. Burns, and R. G. Moyle. 2016.** Comparison of target-  
793 capture and restriction-site associated DNA sequencing for phylogenomics: a test in cardinalid  
794 tanagers (Aves, Genus: *Piranga*). *Systematic Biology* 65: 640-650.
- 795 **Mayer, C., M. Sann, A. Donath, M. Meixner, L. Podsiadlowski, R. S. Peters, M. Petersen,  
796 K. Meusemann, K. Liere, and J.-W. Wägele. 2016.** BaitFisher: a software package for  
797 multispecies target DNA enrichment probe design. *Molecular Biology and Evolution* 33: 1875-  
798 1886.
- 799 **McCormack, J. E., S. M. Hird, A. J. Zellmer, B. C. Carstens, and R. T. Brumfield. 2013a.**  
800 Applications of next-generation sequencing to phylogeography and phylogenetics. *Molecular*  
801 *Phylogenetics and Evolution* 66: 526-538.
- 802 **McCormack, J. E., B. C. Faircloth, N. G. Crawford, P. A. Gowaty, R. T. Brumfield, and T.  
803 C. Glenn. 2012.** Ultraconserved elements are novel phylogenomic markers that resolve placental  
804 mammal phylogeny when combined with species-tree analysis. *Genome Research* 22: 746-754.
- 805 **McCormack, J. E., M. G. Harvey, B. C. Faircloth, N. G. Crawford, T. C. Glenn, and R. T.  
806 Brumfield. 2013b.** A phylogeny of birds based on over 1,500 loci collected by target enrichment  
807 and high-throughput sequencing. *PLoS One* 8: e54848.
- 808 **Mclean, B. S., K. C. Bell, J. M. Allen, K. M. Helgen, and J. A. Cook. 2018.** Impacts of  
809 inference method and data set filtering on phylogenomic resolution in a rapid radiation of ground  
810 squirrels (Xerinae: Marmotini). *Systematic biology* 68: 298-316.
- 811 **Miller, M. R., J. P. Dunham, A. Amores, W. A. Cresko, and E. A. Johnson. 2007.** Rapid and  
812 cost-effective polymorphism identification and genotyping using restriction site associated DNA  
813 (RAD) markers. *Genome Research* 17: 240-248.
- 814 **Minh, B. Q., M. Hahn, and R. Lanfear. 2018.** New methods to calculate concordance factors  
815 for phylogenomic datasets. *bioRxiv* 487801.
- 816 **Mirarab, S., and T. Warnow. 2015.** ASTRAL-II: coalescent-based species tree estimation with  
817 many hundreds of taxa and thousands of genes. *Bioinformatics* 31: i44-i52.
- 818 **Mirarab, S., R. Reaz, M. S. Bayzid, T. Zimmermann, M. S. Swenson, and T. Warnow.  
819 2014.** ASTRAL: genome-scale coalescent-based species tree estimation. *Bioinformatics* 30:  
820 i541-i548.
- 821 **Molloy, E. K., and T. Warnow. 2018.** To Include or Not to Include: The Impact of Gene  
822 Filtering on Species Tree Estimation Methods. *Systematic Biology* 67: 285-303.
- 823 **Musher, L. J., and J. Cracraft. 2018.** Phylogenomics and species delimitation of a complex  
824 radiation of Neotropical suboscine birds (*Pachyramphus*). *Molecular Phylogenetics and*  
825 *Evolution* 118: 204-221.
- 826 **Myers, E. A., R. W. Bryson Jr, R. W. Hansen, M. L. Aardema, D. Lazcano, and F. T.  
827 Burbrink. 2019.** Exploring Chihuahuan Desert diversification in the gray-banded kingsnake,  
828 *Lampropeltis alterna* (Serpentes: Colubridae). *Molecular Phylogenetics and Evolution* 131: 211-  
829 218.

- 830 **Nascimento, F. F., M. dos Reis, and Z. Yang. 2017.** A biologist's guide to Bayesian  
831 phylogenetic analysis. *Nature Ecology & Evolution* 1: 1446-1454.
- 832 **Newman, C. E., and C. C. Austin. 2016.** Sequence capture and next-generation sequencing of  
833 ultraconserved elements in a large-genome salamander. *Molecular Ecology* 25: 6162-6174.
- 834 **Nguyen, L.-T., H. A. Schmidt, A. von Haeseler, and B. Q. Minh. 2014.** IQ-TREE: a fast and  
835 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular*  
836 *Biology and Evolution* 32: 268-274.
- 837 **Nurk, S., D. Meleshko, A. Korobeynikov, and P. A. Pevzner. 2017.** metaSPAdes: a new  
838 versatile metagenomic assembler. *Genome Research* 27: 824-834.
- 839 **Ogilvie, H. A., R. R. Bouckaert, and A. J. Drummond. 2017.** StarBEAST2 brings faster  
840 species tree inference and accurate estimates of substitution rates. *Molecular Biology and*  
841 *Evolution* 34: 2101-2114.
- 842 **Peterson, B. K., J. N. Weber, E. H. Kay, H. S. Fisher, and H. E. Hoekstra. 2012.** Double  
843 digest RADseq: an inexpensive method for de novo SNP discovery and genotyping in model and  
844 non-model species. *PLoS One* 7: e37135.
- 845 **Philippe, H., H. Brinkmann, D. V. Lavrov, D. T. Littlewood, M. Manuel, G. Worheide, and**  
846 **D. Baurain. 2011.** Resolving difficult phylogenetic questions: why more sequences are not  
847 enough. *PLoS Biology* 9: e1000602.
- 848 **Pierce, M. P., M. G. Branstetter, and J. T. Longino. 2017.** Integrative taxonomy reveals  
849 multiple cryptic species within Central American *Hylomyrma* Forel, 1912 (Hymenoptera:  
850 Formicidae). *Myrmecological News* 25: 131-143.
- 851 **Polychronopoulos, D., J. W. King, A. J. Nash, G. Tan, and B. Lenhard. 2017.** Conserved  
852 non-coding elements: developmental gene regulation meets genome organization. *Nucleic Acids*  
853 *Research* 45: 12611-12624.
- 854 **Prebus, M. 2017.** Insights into the evolution, biogeography and natural history of the acorn ants,  
855 genus *Temnothorax* Mayr (hymenoptera: Formicidae). *BMC Evolutionary Biology* 17: 250.
- 856 **Quattrini, A. M., B. C. Faircloth, L. F. Dueñas, T. C. Bridge, M. R. Brugler, I. F. Calixto-**  
857 **Botía, D. M. DeLeo, S. Foret, S. Herrera, and S. M. Lee. 2018.** Universal target-enrichment  
858 baits for anthozoan (Cnidaria) phylogenomics: New approaches to long-standing problems.  
859 *Molecular Ecology Resources* 18: 281-295.
- 860 **Ronquist, F., S. Klopfstein, L. Vilhelmsen, S. Schulmeister, D. L. Murray, and A. P.**  
861 **Rasnitsyn. 2012a.** A total-evidence approach to dating with fossils, applied to the early radiation  
862 of the Hymenoptera. *Systematic Biology* 61: 973-999.
- 863 **Ronquist, F., M. Teslenko, P. Van Der Mark, D. L. Ayres, A. Darling, S. Höhna, B. Larget,**  
864 **L. Liu, M. A. Suchard, and J. P. Huelsenbeck. 2012b.** MrBayes 3.2: efficient Bayesian  
865 phylogenetic inference and model choice across a large model space. *Systematic Biology* 61:  
866 539-542.
- 867 **Ruane, S., and C. C. Austin. 2017.** Phylogenomics using formalin-fixed and 100+ year-old  
868 intractable natural history specimens. *Molecular Ecology Resources* 17: 1003-1008.
- 869 **Ryu, T., L. Seridi, and T. Ravasi. 2012.** The evolution of ultraconserved elements with  
870 different phylogenetic origins. *BMC Evolutionary Biology* 12: 236.
- 871 **Santos, B. F., A. Perrard, and S. G. Brady. 2019.** Running in circles in phylomorphospace:  
872 host environment constrains morphological diversification in parasitic wasps. *Proceedings of the*  
873 *Royal Society of London B: Biological Sciences* 286.



- 874 **Sela, I., H. Ashkenazy, K. Katoh, and T. Pupko. 2015.** GUIDANCE2: accurate detection of  
875 unreliable alignment regions accounting for the uncertainty of multiple parameters. *Nucleic*  
876 *Acids Research* 43: W7-W14.
- 877 **Shen, X. X., C. T. Hittinger, and A. Rokas. 2017.** Contentious relationships in phylogenomic  
878 studies can be driven by a handful of genes. *Nature Ecology & Evolution* 1: 126.
- 879 **Short, A. E. Z., T. Dikow, and C. S. Moreau. 2018.** Entomological Collections in the Age of  
880 Big Data. *Annual Review of Entomology* 63: 513-530.
- 881 **Simpson, J. T., K. Wong, S. D. Jackman, J. E. Schein, S. J. Jones, and I. Birol. 2009.**  
882 ABySS: a parallel assembler for short read sequence data. *Genome Research* 19: 1117-1123.
- 883 **Smith, B. T., M. G. Harvey, B. C. Faircloth, T. C. Glenn, and R. T. Brumfield. 2014.** Target  
884 capture and massively parallel sequencing of ultraconserved elements for comparative studies at  
885 shallow evolutionary time scales. *Systematic Biology* 63: 83-95.
- 886 **Stamatakis, A. 2006.** RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with  
887 thousands of taxa and mixed models. *Bioinformatics* 22: 2688-2690.
- 888 **Starrett, J., S. Derkarabetian, M. Hedin, R. W. Bryson Jr, J. E. McCormack, and B. C.**  
889 **Faircloth. 2017.** High phylogenetic utility of an ultraconserved element probe set designed for  
890 Arachnida. *Molecular Ecology Resources* 17: 812-823.
- 891 **Streicher, J. W., and J. J. Wiens. 2017.** Phylogenomic analyses of more than 4000 nuclear loci  
892 resolve the origin of snakes among lizard families. *Biology Letters* 13: 20170393.
- 893 **Ströher, P. R., A. L. S. Meyer, E. Zarza, W. L. E. Tsai, J. E. McCormack, and M. R. Pie.**  
894 **2019.** Phylogeography of ants from the Brazilian Atlantic Forest. *Organisms Diversity &*  
895 *Evolution*.
- 896 **Suchan, T., C. Pitteloud, N. S. Gerasimova, A. Kostikova, S. Schmid, N. Arrigo, M.**  
897 **Pajkovic, M. Ronikier, and N. Alvarez. 2016.** Hybridization capture using RAD probes  
898 (hyRAD), a new tool for performing genomic analyses on collection specimens. *PLoS One* 11:  
899 e0151651.
- 900 **Tagliacollo, V. A., and R. Lanfear. 2018.** Estimating improved partitioning schemes for  
901 ultraconserved elements. *Molecular Biology and Evolution* 35: 1798-1811.
- 902 **Vachaspati, P., and T. Warnow. 2018.** SVDquest: Improving SVDquartets species tree  
903 estimation using exact optimization within a constrained search space. *Molecular Phylogenetics*  
904 *and Evolution* 124: 122-136.
- 905 **Van Dam, M. H., M. Trautwein, G. S. Spicer, and L. Esposito. 2019.** Advancing mite  
906 phylogenomics: Designing ultraconserved elements for Acari phylogeny. *Molecular Ecology*  
907 *Resources* 19: 465-475.
- 908 **Van Dam, M. H., A. W. Lam, K. Sagata, B. Gewa, R. Laufa, M. Balke, B. C. Faircloth, and**  
909 **A. Riedel. 2017.** Ultraconserved elements (UCEs) resolve the phylogeny of Australasian smurf-  
910 weevils. *PLoS One* 12: e0188044.
- 911 **Ward, P. S., and M. G. Branstetter. 2017.** The acacia ants revisited: convergent evolution and  
912 biogeographic context in an iconic ant/plant mutualism. *Proceedings of the Royal Society of*  
913 *London B: Biological Sciences* 284.
- 914 **Weitemier, K., S. C. Straub, R. C. Cronn, M. Fishbein, R. Schmickl, A. McDonnell, and A.**  
915 **Liston. 2014.** Hyb-Seq: Combining target enrichment and genome skimming for plant  
916 phylogenomics. *Applications in Plant Sciences* 2: 1400042.
- 917 **Wood, H. M., V. L. Gonzalez, M. Lloyd, J. Coddington, and N. Scharff. 2018.** Next-  
918 generation museum genomics: Phylogenetic relationships among palpimanoid spiders using

- 919 sequence capture techniques (Araneae: Palpimanoidea). *Molecular Phylogenetics and Evolution*  
920 127: 907-918.
- 921 **Yang, Z. 2007.** PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and*  
922 *Evolution* 24: 1586-1591.
- 923 **Yang, Z., and B. Rannala. 2012.** Molecular phylogenetics: principles and practice. *Nature*  
924 *Reviews Genetics* 13: 303-314.
- 925 **Yeates, D. K., A. Zwick, and A. S. Mikheyev. 2016.** Museums are biobanks: unlocking the  
926 genetic potential of the three billion specimens in the world's biological collections. *Current*  
927 *opinion in insect science* 18: 83-88.
- 928 **Zarza, E., E. M. Connors, J. M. Maley, W. L. E. Tsai, P. Heimes, M. Kaplan, and J. E.**  
929 **McCormack. 2018.** Combining ultraconserved elements and mtDNA data to uncover lineage  
930 diversity in a Mexican highland frog (*Sarcohyla*; Hylidae). *PeerJ* 6.
- 931 **Zerbino, D., and E. Birney. 2008.** Velvet: algorithms for *de novo* short read assembly using de  
932 Bruijn graphs. *Genome Research* 18: 821-829.
- 933 **Zhang, C., M. Rabiee, E. Sayyari, and S. Mirarab. 2018.** ASTRAL-III: polynomial time  
934 species tree reconstruction from partially resolved gene trees. *BMC Bioinformatics* 19: 153.
- 935 **Zhang, F., Y. Ding, C. Zhu, X. Zhou, M. C. Orr, S. Scheu, and Y.-X. Luan. 2019.**  
936 Phylogenomics from Low-coverage Whole-genome Sequencing. *Methods in Ecology and*  
937 *Evolution* 10: 507-517.
- 938 **Zhou, X., X. X. Shen, C. T. Hittinger, and A. Rokas. 2018.** Evaluating Fast Maximum  
939 Likelihood-Based Phylogenetic Programs Using Empirical Phylogenomic Data Sets. *Molecular*  
940 *Biology and Evolution* 35: 486-503.

941

942

943

944 **Table 1.** Published studies using arthropod UCE datasets as of July 2019.

<b>Taxonomic Group</b>	<b>Probe Sets</b>	<b>References</b>
<b>Hymenoptera</b>	Hymenoptera 1.5Kv1 Hymenoptera 2.5Kv2 Full Hymenoptera 2.5Kv2 Ant-Specific Hymenoptera 2.5Kv2 Bee-Ant-Specific	(Blaimer et al. 2015, Faircloth et al. 2015, Blaimer et al. 2016b, Blaimer et al. 2016a, Branstetter et al. 2017c, Branstetter et al. 2017b, Branstetter et al. 2017a, Ješovnik et al. 2017, Pierce et al. 2017, Prebus 2017, Ward and Branstetter 2017, Blaimer et al. 2018a, Blaimer et al. 2018b, Cooke 2018, Borowiec 2019, Bossert et al. 2019, Branstetter and Longino 2019, Cruaud et al. 2019, Grab et al. 2019, Santos et al. 2019, Ströher et al. 2019)
<b>Arachnida</b>	Arachnida 1.1Kv1 Mite-v2	(Faircloth 2017, Starrett et al. 2017, Derkarabetian et al. 2018, Hedin et al. 2018a, Hedin et al. 2018b, Wood et al. 2018, Derkarabetian et al. 2019, Hedin et al. 2019, Van Dam et al. 2019)
<b>Coleoptera</b>	Coleoptera 1.1Kv1 Adephaga_2.9Kv1	(Baca et al. 2017b, Faircloth 2017, Van Dam et al. 2017, Gustafson et al. 2019)
<b>Hemiptera</b>	Hemiptera 2.7Kv1	(Faircloth 2017, Forthman et al. 2019, Kieran et al. 2019)
<b>Diptera</b>	Diptera 2.7Kv1	(Faircloth 2017)
<b>Lepidoptera</b>	( <i>in silico</i> only)	(Faircloth 2017)
<b>Psocodea</b>	Phthiraptera-2.8Kv1 ( <i>in silico</i> only)	(Zhang et al. 2019)

945

946

947

948

949

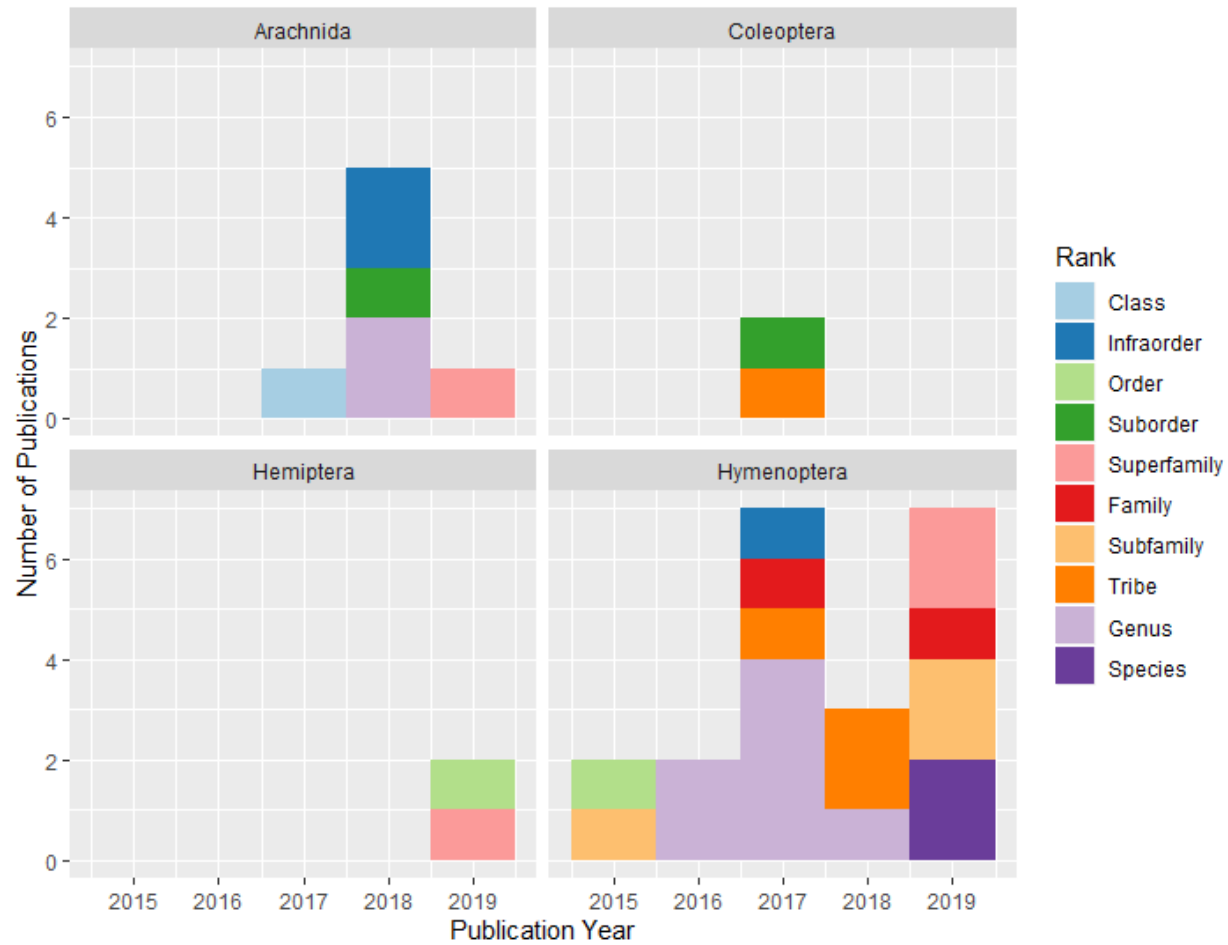
950 **Supplemental Table 1.** Sample cost breakdown of DNA extraction, library preparation, UCE enrichment, and Illumina sequencing  
 951 estimated at around \$40 USD per specimen for 96 samples. (as of July 2019).

Reagent / Kit	Unit of Total Kit	Total Kit Cost	Cost Per Prep	Number Preps	Total Cost
DNeasy Blood and Tissue Kit	250 rxns	\$644.00	\$2.58	96	\$247.30
Library Prep Kits	96 reactions	\$2,496.00	\$6.50	96	\$624.00
HotStart ReadyMix	500 x 25 µL reactions	\$525.00	\$1.05	108	\$113.40
Magnetic SpeedBeads	15 mL	\$385.50	\$25.70	1	\$25.70
			<b>\$35.83</b>		<b>\$1,010.40</b>
Dynabeads MyOne Streptavidin T1	10000 µL	\$1,624.00	\$8.12	12	\$97.44
MYbaits-1	12 reactions - 5.5µL per reaction	\$2,400.00	\$32.73	12	\$392.73
Enrichment Reagents	Misc. Chemicals / Reagents	NA	\$15.00	12	\$180.00
			<b>\$55.85</b>		<b>\$670.17</b>
qPCR Library Quant	500 x 20 µL reactions	\$604.00	\$0.60	90	\$54.36
Gel Cassettes, BluePippin	Cassette for 5 samples	\$450.00	\$50.00	1	\$50.00
High Sensitivity D1000	Tape for 16 samples	\$513.51	\$50.00	1	\$50.00
			<b>\$114.50</b>		<b>\$154.36</b>
HiSeq 125 cycle paired-end	1 lane	\$2,140.00	\$2,140.00	1	<b>\$2,140.00</b>
					<b>\$3,974.92</b>

952

953 **Supplementary Table 2.** Breakdown of methods used by used in 32 UCE-based arthropod studies (as of July 2019).

954 **Figure 1.** Breakdown of the number of arthropod UCEs-based publications per year (as of July  
 955 2019) by taxonomic group and taxonomic hierarchy.



956

957

958

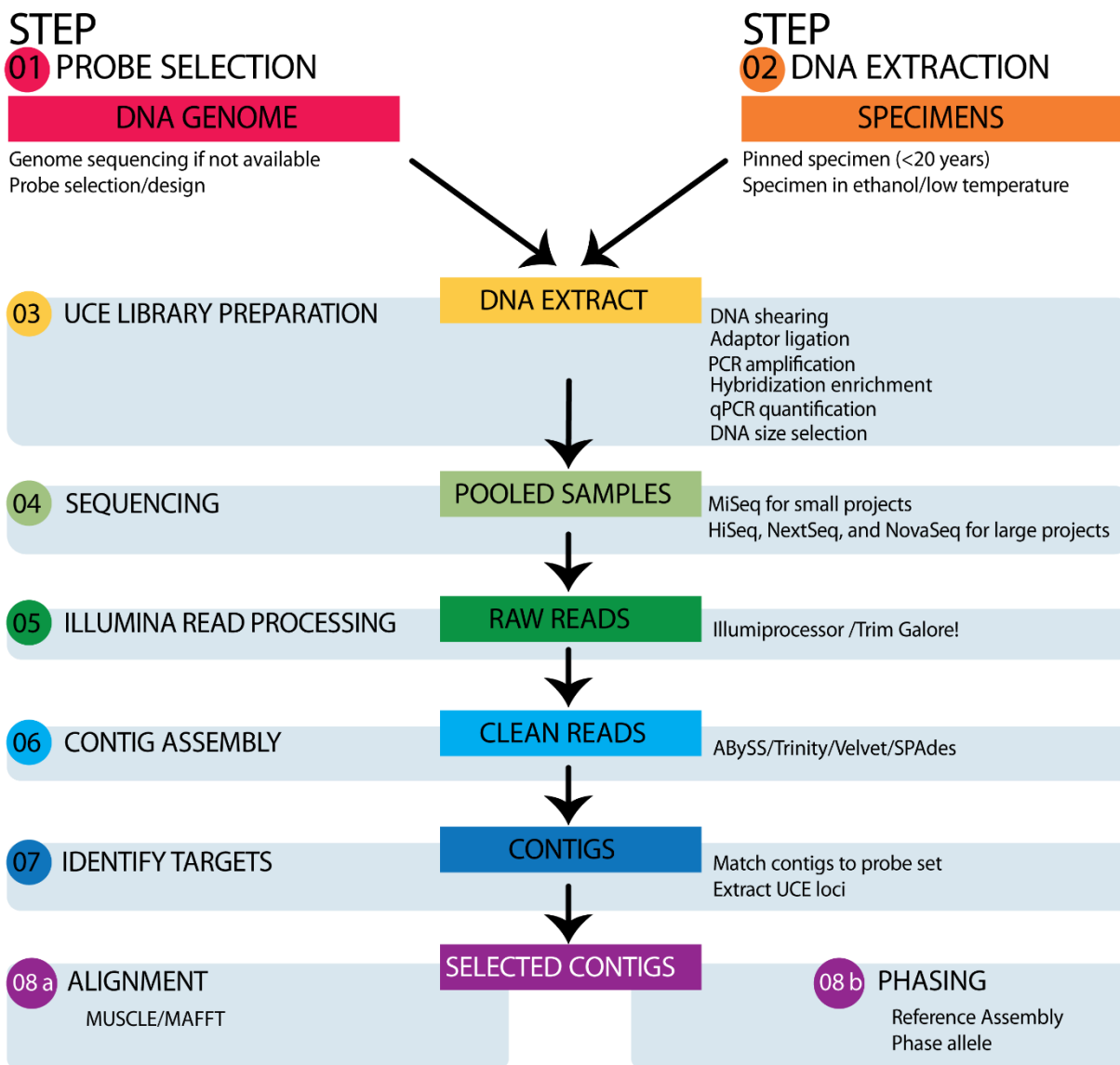
959

960

961

962

963 **Figure 2.** Generalized workflow of the UCE pipeline.



964

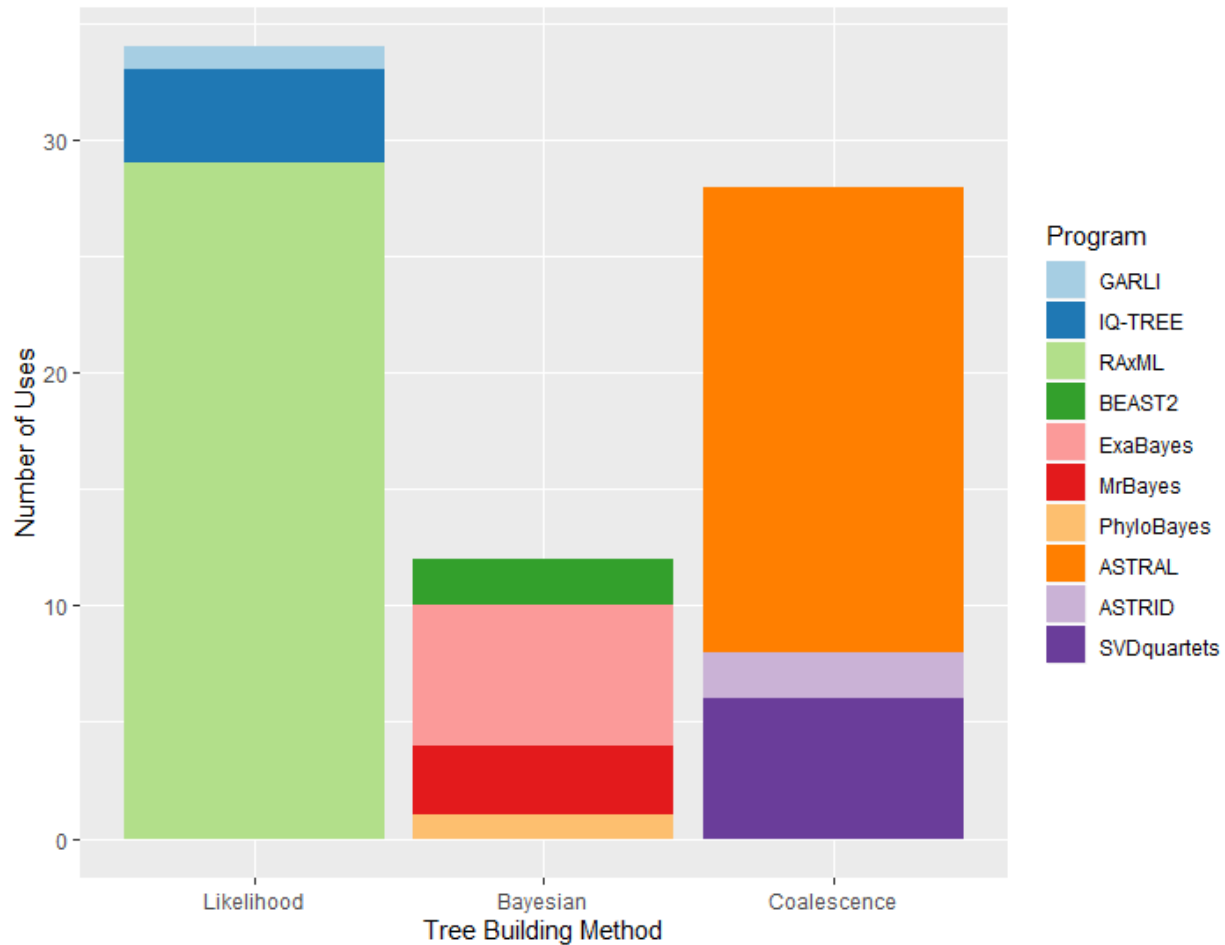
965

966

967

968

969 **Figure 3.** Breakdown of the phylogenetic programs used by arthropod UCE-based publications  
 970 (as of July 2019).



971

972

973 **Supplementary Data. Glossary of commonly encountered terminology in UCE generation**  
 974 **and analyses.**

975

976 **AHE:** Anchored Hybrid Enrichment = A target enrichment method that uses in solution  
 977 hybridization to capture exonic genes for phylogenomic studies.

978

979 **Adapter:** A short piece of known DNA attached to the genomic DNA of interest to identify the  
 980 sample once mixed with other samples.

981

982 **Assembly:** Assembly of fragment sequences into higher order structures based on their overlap  
 983 and reference sequence, where appropriate.

984

985 **Biotinylation:** The process of attaching biotin to proteins and other macromolecules, in this case  
986 to bind the DNA regions of interest to the streptavidin magnet beads during the in-solution  
987 hybridization process. Streptavidin has a high affinity for biotin, being one of the strongest non-  
988 covalent interactions known in nature.

989

990 **Contig:** A contiguous stretch of DNA sequence that is the result of assembly of multiple  
991 overlapping sequence reads into a single consensus sequence.

992

993 **de Bruijn Graph:** A graph theory method for assembling a long sequence from overlapping  
994 fragments. The de Bruijn graph is a set of unique substrings (words) of a fixed length (a  $k$ -mer)  
995 that contain all possible words in the data set exactly once. The sequence reads are split into all  
996 possible  $k$ -mers, and overlapping  $k$ -mers are linked by edges in the graph. Reads are then  
997 mapped onto the graph of overlapping  $k$ -mers in a single pass, greatly reducing the  
998 computational complexity of genome assembly.

999

1000 **de novo Assembly:** Assembly of contigs without a reference genome.

1001

1002 **Exon:** A portion of a gene that is transcribed and spliced to form the final messenger RNA  
1003 (mRNA). Exons contain protein-coding sequence and untranslated upstream and downstream  
1004 regions (3' UTR and 5' UTR). Exons are separated by introns, which are sequences that are  
1005 transcribed by RNA polymerase, but spliced out after transcription and not included in the  
1006 mature mRNA.

1007

1008 **Flanking regions:** The areas immediately to the left and right of the UCE core, which are  
1009 variable and, therefore, the target of UCE capture.

1010

1011 **K-mer:** A motif (or a small word) of length  $k$  observed more than once in a genomic or  
1012 sequenced sequence. E.g., a dinucleotide is a  $k$ -mer where  $k=2$ .

1013

1014 **In Solution Hybridization:** Binding of biotinylated probes with denatured genomic regions of  
1015 interest in the process of several hours in liquid.

1016

1017 **Library:** A set of nucleic acid fragments which has undergone all processing steps and is ready  
1018 for actual sequencing.

1019

1020 **Multiplex:** A library containing various samples labelled with adapters.

1021



1022 **Paired-End Read:** A technology that obtains sequence reads from both ends of a DNA fragment  
1023 template. The use of paired-end sequencing can greatly improve de novo sequencing applications  
1024 by allowing contigs to be joined when they contain read pairs from a single template fragment,  
1025 even if no reads overlap.

1026  
1027 **Probe/Bait:** A collection of oligonucleotides that will bind to specific, conserved genome  
1028 regions of interest, often called baits as they can ‘fish’ out the region of interest. In our review  
1029 we refer to baits as the temporary oligonucleotides used during the probe design process,  
1030 whereas probes are the final product synthesized to capture UCE loci during library prep.

1031  
1032 **Read:** Data output from the analysis of a single fragment (sequence).

1033  
1034  
1035 **SNP:** Single-Nucleotide Polymorphism = sequence divergence in the range of a single base.

1036  
1037 **Target Enrichment:** Capturing genomic regions of interest by hybridization to target-specific  
1038 biotinylated probes, which are then isolated by magnetic beads.

1039  
1040 **UCE:** Ultraconserved Elements = highly-conserved regions within the genome that are shared  
1041 among evolutionarily distant taxa.

1042

1043