

1 **Paralogs and off-target sequences improve phylogenetic resolution in a densely-sampled**
2 **study of the breadfruit genus (*Artocarpus*, Moraceae)**

3

4 Elliot M. Gardner^{1,2,3,4,a,b}, Matthew G. Johnson^{5,a}, Joan T. Pereira⁶, Aida Shafreena Ahmad

5 Puad⁷, Deby Arifianis⁸, Sahromi⁹, Norman J. Wickett¹, and Nyree J.C. Zerega^{1,2,b}

6 ¹ Chicago Botanic Garden, Negaunee Institute for Plant Conservation Science and

7 Action, 1000 Lake Cook Road, Glencoe, IL, 60022, USA

8 ² Northwestern University, Plant Biology and Conservation Program, 2205 Tech Dr.,

9 Evanston, IL, 60208, USA

10 ³ The Morton Arboretum, 4100 IL-53, Lisle, Illinois 60532, USA (current affiliation)

11 ⁴ Case Western Reserve University, Department of Biology, 2080 Adelbert Road,

12 Cleveland, Ohio 44106, USA (current affiliation)

13 ⁵ Texas Tech University, Department of Biological Sciences, 2901 Main Street, Lubbock,

14 TX, 79409-3131, USA

15 ⁶ Forest Research Centre, Sabah Forestry Department, P.O. Box 1407, 90715 Sandakan,

16 Sabah, Malaysia

17 ⁷ Universiti Malaysia Sarawak, Kuching, Sarawak, Malaysia

18 ⁸ Herbarium Bogoriense, Research Center for Biology, Indonesian Institute of Sciences,

19 Cibinong, Jawa Barat, Indonesia

20 ⁹ Center for Plant Conservation Botanic Gardens, Indonesian Institute Of Sciences,

21 Bogor, Jawa Barat, Indonesia

22 ^a Co-first authors

23 ^b For correspondence, email elliott.gardner@case.edu or n-zerega@northwestern.edu

24

25 **Abstract**

26 We present a 517-gene phylogenetic framework for the breadfruit genus *Artocarpus* (ca. 70 spp.,
27 Moraceae), making use of silica-dried leaves from recent fieldwork and herbarium specimens
28 (some up to 106 years old) to achieve 96% taxon sampling. We explore issues relating to
29 assembly, paralogous loci, partitions, and analysis method to reconstruct a phylogeny that is
30 robust to variation in data and available tools. While codon partitioning did not result in any
31 substantial topological differences, the inclusion of flanking non-coding sequence in analyses
32 significantly increased the resolution of gene trees. We also found that increasing the size of
33 datasets increased convergence between analysis methods but did not reduce gene tree conflict.
34 We optimized the HybPiper targeted-enrichment sequence assembly pipeline for short sequences
35 derived from degraded DNA extracted from museum specimens. While the subgenera of
36 *Artocarpus* were monophyletic, revision is required at finer scales, particularly with respect to
37 widespread species. We expect our results to provide a basis for further studies in *Artocarpus*
38 and provide guidelines for future analyses of datasets based on target enrichment data,
39 particularly those using sequences from both fresh and museum material, counseling careful
40 attention to the potential of off-target sequences to improve resolution.

41

42 **Key words:** Phylogenomics, target enrichment, non-coding sequences, Moraceae, *Artocarpus*

43

44

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

45 Despite the increasing availability of whole-genome sequencing, costs and computational
46 limitations still make it impractical, and often unnecessary, for large phylogenetic projects.
47 Reduced-representation methods such as target enrichment sequencing (e.g., HybSeq) have
48 become important tools for phylogenetic studies, enabling high-throughput and cost-effective
49 sequencing of hundreds of informative loci (Faircloth et al. 2012; Mandel et al. 2014; Weitemier
50 et al. 2014). HybSeq involves hybridizing a randomly-sheared sequencing library (e.g., Illumina)
51 to oligonucleotide bait sequences, typically exonic sequences from one or more taxa within or
52 near the target clade. The resulting sequence data include exons and flanking non-coding
53 sequences (e.g., introns, UTRs). While HybSeq recovers fewer loci than the tens of thousands
54 available from RAD-seq, it is more repeatable, recovers longer loci, and has much lower rates of
55 missing data (Weitemier et al. 2014). Accordingly, researchers have successfully employed
56 HybSeq in studies ranging from deep phylogenetics (Prum et al. 2015; Liu et al. 2019) to within-
57 species phylogeography (Villaverde et al. 2018). Like many high throughput sequencing
58 methods, including transcriptomics and RAD-seq, HybSeq can provide hundreds of thousands of
59 characters for phylogenetic reconstruction. Making the most of these large datasets requires
60 careful attention to assembly and analysis methods, particularly when dealing with degraded
61 DNA extracted from museum specimens.

62 HybSeq can be used to recover sequences from degraded DNA extracted from old
63 museum specimens because it relies on capturing short fragments of DNA using 120 bp probes,
64 rather than two flanking primer sequences as in amplicon sequencing. Therefore, HybSeq may
65 succeed where direct PCR-based methods might fail (Staats et al. 2013), substantially raising the
66 value of natural history collections for phylogenetic studies and offering the possibility of
67 including taxa that are not possible to collect due to extinction or infeasible fieldwork (Buerki

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

68 and Baker 2016; Brewer et al. 2019). Extreme fragmentation of old DNA and contamination
69 with non-endogenous sequences, such as fungi, still present challenges to assembly of sequences
70 from short-read platforms, leading many studies to focus on high-copy regions of DNA, such as
71 chloroplast sequences (Bakker et al. 2016). Enriching for targeted genes can transform low-copy
72 nuclear genes into high-copy components of a sequencing library, making HybSeq particularly
73 appropriate for museum or herbarium material (Hart et al., 2016; Villaverde et al., 2018).

74 Ensuring that all sequences used for phylogenetic reconstruction are orthologous is a key
75 step in any phylogenetics workflow. HybSeq will recover any portions of the genome that are
76 sufficiently similar to the bait sequences—up to 25–30% divergence in some cases (Johnson et
77 al. 2019; Liu et al. 2019)—including paralogs or genes with enough shared domains (Hart et al.
78 2016; Johnson et al. 2016). For this reason, HybSeq bait development has usually focused on
79 single to low-copy genes (Chamala et al. 2015; Gardner et al. 2016). Nevertheless, whole
80 genome duplications can render single copy genes double copy in entire clades. We previously
81 developed HybSeq baits for phylogenetic reconstruction of the genus *Artocarpus* J.R. Forst. &
82 G. Forst. (Moraceae) from a three-way orthology search between closely related species in the
83 Rosales: *Cannabis sativa* L. (Cannabaceae) (van Bakel et al. 2011), *Morus notabilis* C.K.
84 Schneid. (Moraceae) (He et al. 2013), and *Artocarpus camansi* Blanco (Gardner et al. 2016).
85 Due to an ancient whole-genome duplication in *Artocarpus*, many of the 333 genes were
86 represented as paralogous pairs in that genus, although in almost all cases they were diverged
87 enough to sort and analyze separately (Johnson et al. 2016). The impact of this by-catch on
88 phylogenetic reconstruction remains unclear, but has the potential to greatly increase the number
89 of phylogenetically informative genes, if the paralogs can be easily sorted.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

90 Accurate phylogenetic reconstruction relies on carefully considered parameters for
91 sequence alignment, trimming, and partitioning of the character matrix, among other factors.
92 Because HybSeq recovers not only targeted exons but also flanking non-coding sequences,
93 contigs containing both exons and introns are often aligned and analyzed together (e.g. (Medina
94 et al. 2019). This can make it difficult to ensure that exons are aligned in frame, particularly
95 when frameshifts are present due to the presence of pseudogenes, sequencing errors, or other
96 processes (Ranwez 2011). Aligning coding and noncoding regions together can also hamper
97 partitioning of datasets by codon position, but how these issues impact phylogenetic
98 reconstruction remains unclear (Xi et al. 2012; Lanfear et al. 2014).

99 Concatenating all loci into a supermatrix can result in near-perfect bootstrap support for
100 almost all nodes in a phylogeny (Sayyari and Mirarab 2016). Despite this apparent high support,
101 there may be substantial discordance among gene histories due to incomplete lineage sorting
102 (Kubatko and Degnan 2007; Degnan and Rosenberg 2009). Although there is an increased
103 availability of efficient methods that are consistent under the predictions of the multi-species
104 coalescent model, clear results can be obscured if the underlying gene trees are uninformative
105 (Smith et al. 2015; Sayyari et al. 2017). A major advantage of HybSeq over methods with short,
106 anonymous loci, or large amounts of missing data, is that loci are both long enough to generate
107 single-gene phylogenies and subject to few enough missing taxa per locus for those single-gene
108 phylogenies to be informative. In this paper we explore the possibilities of HybSeq, including the
109 informativeness of paralogs, effects of missing data, and utility of herbarium specimens to
110 reconstruct the most data-rich phylogeny of the breadfruit genus to date.

111

112 *Study system*

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

113 *Artocarpus* J.R. Forst. & G. Forst. (Moraceae, Figure 1) contains approximately 70
114 species of monoecious trees with a center of diversity in Borneo and a native range that extends
115 from India to the Solomon Islands (Williams et al. 2017). The genus is best known for important
116 but underutilized crops such as breadfruit (*A. altilis* (Parkinson) Fosberg) and jackfruit (*A.*
117 *heterophyllus* Lam.), and it also contains crops of regional importance like cempedak (*A. integer*
118 (Thunb.) Merr.) and tarap (*A. odoratissimus* Blanco), and more than a dozen other species with
119 edible fruits whose potential remains largely unexplored (Zerega et al. 2010, 2015; Wang et al.
120 2018; Witherup et al. 2019).

121 *Artocarpus* is characterized by spicate to globose staminate (“male”) inflorescences
122 composed of tiny flowers bearing one stamen each. Carpellate (“female”) inflorescences are
123 composed of tightly packed tiny flowers. In most cases, adjacent carpellate flowers are at least
124 partially fused together. Carpellate inflorescences develop into syncarps, which are tightly
125 packed accessory fruits composed mainly of fleshy floral tissue. Syncarps of different species
126 range in size from a few centimeters in diameter to over half a meter long. *Artocarpus* is the
127 largest genus in the tribe Artocarpeae, which also includes two smaller Neotropical genera,
128 *Batocarpus* H. Karst. (3 spp.) and *Clarisia* Ruiz & Pav. (3 spp.). The neotropical genera always
129 have spicate staminate inflorescences; carpellate flowers may be solitary or condensed into
130 globose heads, but neither tepals nor adjacent flowers are fused.

131 The most recent complete revision of *Artocarpus* (Jarrett 1959, 1960) recognized two
132 subgenera, *Artocarpus* and *Pseudojaca* Tréc., distinguished by phyllotaxy (leaf arrangement),
133 and the degree of fusion between adjacent carpellate flowers. Since then, several new species
134 have been described by Jarrett and others (Jarrett, 1975; Zhengyi and Xiushi, 1989; Kochummen,
135 1998; Berg, 2005; Gardner et al. in prep.; Gardner and Zerega in prep.). Berg et al. (2006)

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

136 revised the Malesian species for the *Flora Malesiana*, in a few cases combining several taxa into
137 a broadly-circumscribed single species. Examples include *A. altilis* (encompassing *A. altilis*, *A.*
138 *camansi*, *A. mariannensis* Tréc., *A. horridus* F.M. Jarrett, *A. blancoi* Merr., *A. pinnatisectus*
139 Merr., and *A. multifidus* F.M. Jarrett) and *A. lacucha* Buch.-Ham. (encompassing *A. lacucha*, *A.*
140 *dadah* Miq., *A. fretessii* Teijsm. & Binnend., *A. ovatus* Blanco, and *A. vrieseanus* var. *refractus*
141 Becc. (F.M. Jarrett)). In these cases, this paper follows Jarrett's nomenclature for clarity. The
142 most recent circumscription of *Artocarpus* recognized four subgenera within *Artocarpus* (Figure
143 1) and was based on just two gene regions and approximately 50% of taxa (Zerega et al. 2010).
144 The subgenera were distinguished by phyllotaxy, the degree of fusion between adjacent
145 carpellate flowers, and the position of inflorescences on the tree (axillary or cauliflorous). A
146 well-sampled phylogenetic framework for *Artocarpus* is necessary to inform future taxonomic
147 revision and to clarify relationships within this important genus, in particular the relationships
148 between crop species and their wild relatives, whose conservation is a priority (Castañeda-
149 Álvarez et al. 2016).

150 In this study, we used near-complete (80/83) taxon sampling (at the subspecies level or
151 above) in *Artocarpus* to explore the impact of paralogs, codon partitions, noncoding sequences,
152 and analysis method (species tree versus concatenated supermatrix) on phylogenetic
153 reconstruction in order to develop best practices for the analysis of HybSeq data and to explore
154 the limits of phylogenomic resolution. We also used this data set to improve the target capture
155 assembly pipeline HybPiper, which is now optimized for accurately scaffolding small
156 disconnected contigs resulting from degraded DNA. The objectives of the study were to (1) Use
157 broad sampling from silica dried material and herbarium specimens over 100 years old to
158 achieve near complete taxon sampling for *Artocarpus*; (2) examine the impact of paralogs,

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

159 partitions, and analysis method on phylogenetic reconstruction; and (3) test the monophyly of the
160 current taxonomic divisions within *Artocarpus* to provide a phylogenetic framework for future
161 studies on the taxonomy, conservation, and ecology of the genus.

162

163 MATERIALS AND METHODS

164

165 *Taxon sampling*

166 We sampled all *Artocarpus* taxa at the subspecies level or above recognized by Jarrett
167 (1959, 1960), Berg et al. (2006), and Kochummen (1998), all three obsolete species that Jarrett
168 (1959) sunk into *A. treculianus* Elmer, and all of the new species described by Wu and Zhang
169 (1989), for a total of 83 named *Artocarpus* taxa. We also sampled nine taxa of questionable
170 affinities. We replicated samples across geographic or morphological ranges when possible, for a
171 total of 167 ingroup samples. As outgroups, we sampled one member of each genus in the
172 Neotropical Artocarpeae (*Batocarpus* and *Clarisia*) and the sister tribe Moreae (*Morus* L.,
173 *Streblus* Lour., *Milicia* Sim., *Trophis* P. Browne, *Bagassa* Aubl., and *Sorocea* A. St.-Hil.). We
174 obtained samples from our own field collections preserved in silica gel (from Malaysia,
175 Thailand, Hong Kong, Bangladesh, and India, and from botanic gardens in Indonesia, Malaysia,
176 and Hawai'i, USA) and from herbarium specimens up to 106 years old (from the following
177 herbaria: BM, BO, CHIC, E, F, HAST, HK, K, KUN, L, MO, NY, KEP, S, SAN, SNP, US). In
178 total we included 179 samples (Table S1).

179

180 *Sample preparation and sequencing*

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

181 We sampled approximately 0.5 cm² of dried leaf from each sample for DNA extraction.
182 For herbarium specimens, we sampled from a fragment packet when feasible and when it was
183 clear that the material in the fragment packet originated from the specimen on the sheet
184 (something that cannot always be assumed with very old specimens). DNA was extracted using
185 one of three methods; (1) the Qiagen DNeasy Plant Mini Kit (Qiagen, Valencia, California,
186 USA) following the manufacturer's protocol; (2) the MoBio PowerPlant Pro DNA Kit, (MoBio
187 Laboratories, Carlsbad, California, USA); or (3) a modified CTAB protocol (Doyle and Doyle
188 1987). For kit extractions, the protocols were modified for herbarium material by extending
189 initial incubation times (Williams et al., 2017) and adding an additional 200 μ L of ethanol to the
190 column-binding step. CTAB extractions of herbarium specimens, which often had high but
191 impure DNA yields, were cleaned using a 1:1.8:5 ratio of sample, SPRI beads, and isopropanol,
192 the latter added to prevent the loss of small fragments (Lee 2014). For herbarium specimens, we
193 sometimes combined two or more separate extractions in order to accumulate enough DNA for
194 library preparation. We assessed degradation of DNA from herbarium specimens using either an
195 agarose gel or a High-Sensitivity DNA Assay on a BioAnalyzer 2100 (Agilent) and did not
196 sonicate samples whose average fragment size was less than 500bp. The remaining DNA
197 samples were sonicated to a mean insert size of 550bp using a Covaris M220 (Covaris, Wobum,
198 Massachusetts, USA). Libraries were prepared with either the Illumina TruSeq Nano HT DNA
199 Library Preparation Kit (Illumina, San Diego, California, USA) or the KAPA Hyper Prep DNA
200 Library Kit following the manufacturer's protocol, except that reactions were performed in one-
201 third volumes to save reagent costs. We used 200ng of input DNA when possible; for some
202 samples, input was as low as 10ng. For herbarium samples with degraded DNA, we usually did
203 not perform size selection, unless there were some fragments that were above 550bp. We also

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

204 diluted the adapters from 15 μ M to 7.5 μ M, and usually performed only a single SPRI bead
205 cleanup between adapter ligation and PCR amplification. Many of these libraries contained
206 substantial amounts of adapter dimer, so we adjusted the post-PCR SPRI bead cleanup ratio to
207 0.8x. Libraries were enriched for 333 phylogenetic markers (Gardner et al., 2016) with a
208 MYbaits kit (MYcroarray, Ann Arbor, Michigan, USA) following the MYbaits manufacturer's
209 protocol (version 3). Hybridization took place in pools of 6–24 libraries; within each pool, we
210 used equal amounts of all libraries (20–100ng, as available), and tried to avoid pooling samples
211 with dramatically different phylogenetic distances to the bait sequences (*Morus* and *Artocarpus*),
212 as closer taxa can out-compete multiplexed distant taxa in hybridization reactions, as we
213 previously found when pooling *Dorstenia* L. and *Parartocarpus* Baill. with *Artocarpus* (Johnson
214 et al. 2016). We reamplified enriched libraries with 14 PCR cycles using the conditions specified
215 in the manufacturer's protocol. In some cases, adapter dimer remained even after hybridization;
216 in those cases, we removed it either using a 0.7x SPRI bead cleanup or, in cases where the
217 library fragments were very short (ca. 200bp, compared to 144bp for the dimer), by size-
218 selecting the final pools to >180bp on a BluePippin size-selector using a 2% agarose gel cassette
219 (Sage Science, Beverley, Massachusetts, USA). Pools of enriched libraries were sequenced on
220 an Illumina MiSeq (600 cycle, version 3 chemistry) alongside samples for other studies in three
221 multiplexed runs each containing 30–99 samples.

222

223 *Sequence quality control and analyses*

224 Demultiplexing and adapter trimming took place automatically through Illumina
225 BaseSpace (basespace.illumina.com). All reads have been deposited in GenBank (BioProject no.
226 PRJNA322184). Raw reads were quality trimmed using Trimmomatic (Bolger et al., 2014), with

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

227 a quality cutoff of 20 in a 4-bp sliding window, discarding any reads trimmed to under 30 bp. In
228 addition to the samples sequenced for this study, reads used for assemblies included all
229 *Artocarpus* samples sequenced in Johnson et al. (2016) (available under the same BioProject
230 number). Common methods for target capture assembly include mapping reads to a reference
231 (Weitemier et al. 2014; Hart et al. 2016) and *de novo* assemblies (Mandel et al. 2014; Faircloth
232 2015), but both have drawbacks. Read mapping can result in lost data, particularly indels and
233 non-coding regions, unless a close reference is available. On the other hand, *de novo* assemblies
234 can also result in lost data if loci cannot be assembled into single scaffolds. A compromise
235 approach, implemented in HybPiper, is to combine local *de novo* assemblies—which may result
236 in many small contigs per locus—with scaffolding based on a reference coding sequence, which
237 need not be closely related; a reference with less than 30% sequence, typically within the same
238 family or order, will usually suffice (Johnson et al. 2016, 2019). The resulting assemblies thus
239 cover the maximum available portion of each locus, notwithstanding the existence of long gaps,
240 and also make use of all available on-target reads, including introns, not simply those that can be
241 aligned to a reference.

242 We assembled sequences using HybPiper 1.2, which represented an update of the original
243 pipeline optimized for short reads from highly-fragmented DNA from museum specimens.
244 HybPiper’s guided assembly method uses the reference to scaffold localized *de novo* assemblies.
245 This is particularly advantageous when dealing with very short reads from degraded DNA,
246 because for those samples, reads covering a single exon may assemble into more than one contig.
247 In those cases, HybPiper uses the reference to scaffold and concatenate multiple contigs into a
248 “supercontig” containing the gene of interest as well as any flanking noncoding sequences
249 (Johnson et al. 2016). The new version of HybPiper is optimized to accurately handle many

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

250 small contigs covering a single gene, deduplicating overlaps and outputting high-confidence
251 predicted coding sequences even in the presence of many gaps caused by fragmentary local
252 assemblies. HybPiper as well as all related scripts used in this study are available at
253 <https://github.com/mossmatters/HybPiper> and <https://github.com/mossmatters/phyloscripts>. We
254 generated a new HybPiper reference for this study, using reads from all four subgenera of
255 *Artocarpus*. Target-enriched reads from *A. camansi* Blanco (the same individual used for whole-
256 genome sequencing in the original marker development (Gardner et al. 2016), *A. limpato* Miq.,
257 *A. heterophyllus*, and *A. lacucha* (the latter three from reads sequenced in Johnson et al. (2016))
258 were assembled de novo using SPAdes (Bankevich et al. 2012), and genes were predicted using
259 Augustus (Keller et al. 2011), with *Arabidopsis* Hehyn. as the reference. Predicted genes were
260 annotated using a BLASTn search seeded with the HybPiper target file of 333 phylogenetic
261 marker genes from Johnson et al. (2016). Paralogs were annotated as follows: genes covering at
262 least 75% of the primary ortholog (labeled “p0”) were labeled as “paralogs” (“p1”, “p2”, etc.).
263 Genes covering less than 75% of the primary ortholog (labeled “e0”) were labeled as “extras”
264 (“e1”, “e2”, etc.), denoting uncertainty as to whether they are paralogs or merely genes with a
265 shared domain. Single copy genes were labeled as “single” in the new reference. We used this
266 new 4-taxon reference to guide all ingroup assemblies, and we used the original set of *Morus*
267 *notabilis* targets (Johnson et al., 2016) to guide all outgroup assemblies.

268 We set the per-gene coverage cutoff to 8x, except for certain low-read samples where
269 gene recovery was improved by lowering the coverage cutoff to 4x (10 samples) or 2x (18
270 samples). HybPiper relies on SPAdes for local *de novo* assemblies. SPAdes creates several
271 assemblies with different k-mer values, with the maximum estimated from the reads (up to
272 127bp), and then merges them into a final assembly. For herbarium samples that initially

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

273 recovered fewer than 400 genes, we reran HybPiper, manually setting the maximum k-mer
274 values for assembly to 55 instead of allowing SPAdes to automatically set it. To extract non-
275 coding sequences and annotate gene features along assembled contigs, we used the HybPiper
276 script “intronerate.py”. We assessed target recovery success using the `get_seq_lengths.py` and
277 `gene_recovery_heatmap.r` scripts from HybPiper.

278 To mask low-coverage regions likely to contain sequencing errors, we mapped each
279 sample’s reads to its HybPiper supercontigs using BWA (Li and Durbin 2009), removed PCR
280 duplicates using Picard (Broad Institute 2016), and calculated the depth at each position with
281 Samtools (Li et al. 2009). Using BedTools (Quinlan and Hall 2010), we then hard-masked all
282 positions covered by less than two unique reads. We then used the masked supercontigs and the
283 HybPiper gene annotation files to generate masked versions of the standard HybPiper outputs
284 (using `intron_exon_extractor.py`): (1) the predicted coding sequence for each target gene
285 (“exon”); (2) the entire contig assembled for each gene (“supercontig”); and (3) the predicted
286 non-coding sequences for each gene (“non-coding”, including introns, UTRs, and intergenic
287 sequences).

288 To the HybPiper output, we added the original orthologs (CDS only) identified in *Morus*
289 *notabilis* (Gardner et al., 2016). Because paralogs were only assembled for ingroup samples (due
290 to an *Artocarpus*-specific whole-genome duplication (Gardner et al. 2016)), we added the
291 corresponding “p0” or “e0” from *Morus* to each paralog alignment to serve as an outgroup.

292 We filtered each set of sequences as follows. For “exon” sequences, we subtracted
293 masked bases (Ns) and removed sequences less than 150 bp and sequences covering less than
294 20% of the average sequence length for that gene. For “supercontig” sequences, we removed

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

295 sequences whose corresponding “exon” sequences had been removed. Samples with less than
296 100 genes remaining after filtering were excluded from the main analyses.

297 Alignment and trimming then proceeded as follows. For “exon” output, after removing
298 the genes and sequences identified during the filtering stage, we created in-frame alignments
299 using MACSE (Ranwez 2011). For “supercontig” output, we used MAFFT for alignment (--
300 maxiter 1000) (Kato and Standley 2013). We trimmed all alignments to remove all columns
301 with >75% gaps using Trimal (Capella-Gutiérrez et al. 2009).

302 To quickly inspect gene trees for artifacts, we built gene trees from the trimmed “exon”
303 alignments using FastTree (Price et al. 2009) and visually inspected the gene trees for outlier
304 long branches within the ingroup to identify alignments containing improperly sorted paralogous
305 sequences. In some cases, we visually inspected alignments using AliView (Larsson 2014). We
306 discarded a small number of genes whose alignments contained paralogous sequences, for a final
307 set of 517 genes, including all of the original 333 genes.

308 We used the trimmed alignments to create three sets of gene alignment datasets:

- 309 1. *CDS*: “exon” alignments, not partitioned by codon position;
- 310 2. *Partitioned CDS*: 333 “exon” alignments, partitioned by codon position; and
- 311 3. *Supercontig*: “supercontig” alignments, not partitioned within genes

312 We also attempted to create a codon-partitioned supercontig alignment by separately
313 aligning “exon” and “intron” sequences and then concatenating them, resulting in three partitions
314 per gene. However, this dataset differed substantially from the *supercontig* dataset, resulting in
315 substantially differing (and nonsensical) topologies even when the partitions were removed;
316 samples with a high proportion of very short or missing non-coding sequences clustered together,
317 perhaps because aligning very short non-coding sequences without longer coding sequences to

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

318 anchor them produced unreliable alignments. We therefore did not include the *partitioned*
319 *exon+intron* dataset in the main analyses (discussed further in Appendix 1).

320 To investigate whether including both copies of a paralogous locus impacted
321 phylogenetic reconstruction, we created versions of each dataset with and without paralogs. We
322 analyzed each of these six datasets using the following two methods, for a total of 12 analyses:
323 (A) *Concatenated supermatrix*: all genes were concatenated into a supermatrix, with each gene
324 partitioned separately (i.e. 1 or 3 partitions per gene, depending on the dataset) and analyzed
325 using RAxML 10 (Stamatakis, 2006) under GTR+CAT model with 200 rapid bootstrap
326 replicates, rooted with the Moreae outgroups; (B) *Species tree*: each gene alignment was
327 analyzed using RAxML 10 under the GTR+CAT model with 200 rapid bootstrap replicates,
328 rooted with the Moreae outgroups. Nodes with <33% support were collapsed into polytomies
329 using SumTrees (Sukumaran and Holder 2010), and the resulting trees were used to estimate a
330 species tree with ASTRAL-III (Mirarab and Warnow, 2015). We estimated node support with
331 multilocus bootstrapping (-r, 160 bootstrap replicates) and by calculating the proportion of
332 quartet trees that support each node (-t 1) (Mirarab and Warnow 2015; Zhang et al. 2017). For
333 the final trees, we also used SumTrees to calculate the proportion of gene trees supporting each
334 split. Quartet support is directly related to the method ASTRAL uses for estimating species
335 trees—decomposing gene trees into quartets (Mirarab and Warnow 2015); it is also less sensitive
336 to occasional out-of-place taxa than raw gene-tree support.

337 Because all RAxML analyses were conducted using the GTRCAT model, we also
338 repeated the analyses of the CDS datasets using the GTRGAMMA model to investigate the
339 robustness of the recovered topologies to slight model differences.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

340 To summarize the overall bootstrap support of each tree with a single statistic, we
341 calculated “percent resolution” as the number of bipartitions with >50% bootstrap support
342 divided by the total number of bipartitions and represents the proportion of nodes that one might
343 consider resolved (Kates et al. 2018). We visualized trees using FigTree (Rambaut 2016) and the
344 APE package in R (Paradis et al. 2004). To compare trees, we used the phytools package in R
345 (Revell 2012) to plot a consensus tree and to calculate a Robison-Foulds (RF) distance matrix for
346 all trees. The RF distance between tree *A* and tree *B* equals the number of bipartitions unique to
347 *A* plus the number of bipartitions unique to *B*. We visualized the first two principal components
348 of the matrix using the Lattice package in R (Sarkar 2008). In addition, we conducted pairwise
349 topology comparisons using the “phylo.diff” function from the Phangorn package in R (Schliep
350 2011) and an updated version of “cophylo” from phytools (github.com/liamrevell/phytools/). All
351 statistical analyses took place in R (R Core Development Team, 2008).

352 Supermatrix analyses took place on the CIPRES Science Gateway (Miller et al. 2010).
353 All other analyses took place on a computing cluster at the Chicago Botanic Garden, and almost
354 all processes were run in parallel using GNU Parallel (Tange 2018). Alignments and trees have
355 been deposited in the Dryad Data Repository (accession no. TBA).

356

357

358 RESULTS

359

360 *Sequencing and assembly*

361 Of the 179 sequenced accessions, 164 resulted in successful HybPiper assemblies (>25
362 genes) (Figure 2, Table S1), including all attempted taxa except for *A. scandens* Miq. sensu

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

363 Jarrett—included by Berg et al. (2006) in *A. scandens*, which was assembled—*A. nigrifolius*
364 C.Y. Wu and *A. nanchuanensis*, C.Y. Wu, two species closely allied to *A. hypargyreus* Hance,
365 which was also assembled. The less successful samples almost all had very few reads and may
366 have been out-competed by other samples during hybridization, reamplification, or both. Fewer
367 reads were also associated with shorter assembled sequences (Figure 2). Five out of the 164
368 HybPiper assemblies included sequences for less than 100 genes after filtering and were
369 excluded from the main analyses. This included the only remaining accession of *A. reticulatus*
370 Miq. Adding the *Morus notabilis* sequences therefore resulted in a main data set containing 160
371 samples representing 80 out of 83 named *Artocarpus* taxa at the subspecies/variety level or
372 above (96%), in addition to nine taxa of uncertain affinity.

373 Overall, samples collected more recently showed improved sequencing results (Figure 2),
374 primarily because the majority of samples collected since 2000 were dried on silica gel. Whether
375 a sample was dried on silica gel was significantly associated with increased gene length as a
376 percentage of average length ($R^2 = 0.33$, $P < 0.0001$) and to a lesser extent with the total number
377 of genes recovered ($R^2 = 0.17$, $P < 0.0001$). All 16 unsuccessful (<25 genes) assemblies were
378 taken from herbarium sheets (with collection years spanning 1917 to 1997), rather than from
379 silica-dried material. Among 67 successfully-assembled samples taken from herbarium sheets,
380 younger age was significantly associated only with increased gene length, although the model
381 was a poor fit ($R^2 = 0.06$, $P = 0.02728$), and not with an increase in the number of genes
382 recovered ($P = 0.2833$) (Figure 3). By the same token, we observed a decrease in average DNA
383 fragment size in older samples (Figure S1). Likewise, lowering the maximum assembly k-mer
384 values for herbarium samples with under 400 genes increased recovery by an average of 20
385 genes.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

386 Gene recovery was high; the average sample (of the 160 passing the final filter) had
387 sequences for 448/517 genes (87%). In the final filtered dataset of 333 genes, the average gene
388 had exon sequences for 151/160 samples (94%, range 57–160, median 154) and non-coding
389 sequences for 130 (81%, range 49–151, median 131). For the 184 paralogs, the average gene in
390 the final filtered dataset had exon sequences for 117/160 samples (73%, range 32–148, median
391 126) and intron sequences for 101 (63%, range 30–132, median 110) (Table S2). The
392 supermatrix of trimmed “exon” alignments for the primary 333 genes contained 407,310
393 characters; and the full set of 517 “exon” alignments, including 184 paralogs, contained 569,796
394 characters. The supermatrix of 333 trimmed “supercontig” alignments contained 813,504
395 characters, and the full set of 517 genes contained 1,181,279 characters. The full set of “exon”
396 alignments had 21% gaps or undetermined characters, while the full set of “supercontig”
397 alignments was 36.87% gaps or undetermined characters.

398

399 *Phylogenetic disagreement*

400 A strict consensus of the 12 species trees under GTR+CAT, with a 20% length cutoff
401 (henceforth, the “main analysis”) had 100/159 (63%) nodes resolved (mean RF distance 53),
402 revealing complete agreement among the various analyses in backbone relationships but
403 substantial disagreement at shallower nodes (Figure 5). The six ASTRAL phylogenies differed
404 little from one another, whereas the supermatrix analyses had somewhat greater divergence
405 (Figure 4, Figure S2).

406 *Partitions and model selection* — In the exon datasets, partitioning by codon position
407 (Figures 4, S2, S3) had little impact on the final topology, resulting in a single rearrangement
408 within *A. lacucha* s.s. in the supermatrix analysis (RF 4), and in the ASTRAL analysis a change

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

409 in the position of an undetermined sample and the change of *A. sepicanus* + *A. altissimus* from a
410 clade to a grade (RF 12). The choice of model (GTRCAT vs GTRGAMMA) also produced only
411 minute changes (Figures S2, S7).

412 *Paralogs*— The addition of paralogs led to slightly more disagreement (Figures 4, S4,
413 Table S3). In the exon dataset, changes to the positions of *A. nitidus* ssp. *lingnanensis* (Merr.)
414 F.M. Jarrett and *A. gomezianus* Wall. ex Tréc. affected the backbone of ser. *Peltati* F.M. Jarrett,
415 subg. *Pseudojaca*, in the supermatrix analysis (RF 58); in the ASTRAL analysis, there were
416 somewhat fewer rearrangements, but the mainly occurred in the same clade (RF 20). However,
417 the disagreement was reduced when noncoding sequences were included (supermatrix RF 22;
418 ASTRAL RF 8).

419 *Introns*— The inclusion of non-coding sequences (Figures 4, S5, Table S3) led to similar
420 amounts of disagreement, with rearrangements at the series level in subg. *Pseudojaca* and subg.
421 *Artocarpus*. Disagreement was greater in the supermatrix analyses (no paralogs) (RF 62) than in
422 the ASTRAL analyses (RF 36). The addition of paralogs reduced disagreement in both cases
423 (supermatrix RF 38; ASTRAL RF 26).

424 *Analysis*— The greatest differences among the 12 trees were between ASTRAL trees and
425 the supermatrix trees (Figures 4, S6, Table S3), with a mean RF distance between the six
426 supermatrix trees and the six ASTRAL trees of 78. Again, the addition of additional sequences in
427 the form of noncoding regions or paralogs reduced the disagreement between supermatrix and
428 ASTRAL analyses; the average RF distance for exons/noparalogs was 85, exons+paralogs 77,
429 supercontig/noparalogs 71, and supercontigs+paralogs 70. Agreement was higher among the
430 ASTRAL trees (mean RF 21, 138/159 nodes in agreement) than among the supermatrix trees
431 (mean RF 48, 116/159 nodes in agreement). The differences (RF 66) between ASTRAL and

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

432 supermatrix analyses for the full dataset (supercontigs+paralogs for all genes) at the species level
433 can be ascribed to mostly minor repositionings of two outgroup taxa (*Bagassa guianensis* Aubl.
434 and *Batocarpus orinoceros* H. Karst.) and 14 ingroup taxa (*A. hypargyreus* Hance, *A.*
435 *thailandicus* C.C. Berg., *A. gomezianus* ssp. *gomezianus*, *A. sp. aff. lacucha*, *A. fulvicortex* F.M.
436 Jarrett, *A. sp. aff nitidus*, *A. sp. aff fretesii*, *A. ovatus*, *A. rubrovenius* Warb., *A. pinnatisectus*, *A.*
437 *cf. horridus*, *A. cf. camansi*, *A. excelsus* F.M. Jarrett, *A. jarrettiae* Kochummen) (Figure 6).

438

439 *Phylogenetic resolution*

440 Bootstrap support was high for all main analysis trees. Percent resolution was 90–95%
441 for all supermatrix trees in the main analysis and did not differ materially between analyses.
442 Among ASTRAL trees, percent resolution was between 84% and 97% for all analyses. By slight
443 margins, the best-resolved trees for both supermatrix and ASTRAL analyses were those based on
444 the largest dataset (Figure 6). Percent resolution based on quartet support for ASTRAL trees was
445 between 57% and 60%. Significant conflict existed at the gene tree level (Table S4, Figure 6).
446 For percent resolution measured by gene tree support (percentage of nodes supported by at least
447 half of the 517 gene trees), scores ranged from 17–24%. In general, analyses including paralogs
448 had reduced gene tree support, and the trees based on supercontigs with no paralogs had the
449 highest scores (24% for both ASTRAL and supermatrix).

450 A more detailed analysis of the differences between the *exon* and *supercontig* datasets
451 revealed that even if the final species trees had similar resolution, the *supercontig* trees were
452 based on more information because the gene trees were significantly more informative. A non-
453 parametric Wilcoxon test indicated that the inclusion of non-coding sequences significantly
454 increased both the mean bootstrap support (+5.47, $P < 0.0001$) and the number of splits with

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

455 over 30% support (+8, $P = 0.0022$). Because nodes under 30% were collapsed for species tree
456 estimation, the species tree in the *supercontig* dataset was based on 9% more splits across the
457 517 gene trees (total of 51,307) than the species tree in the *exon* dataset (total 47,067). These
458 patterns persisted in the no-paralog datasets (difference in mean bootstrap support: +2.06, $P <$
459 0.0001; difference in nodes over 30%: +9, $P < 0.0001$; overall difference in splits for 333
460 collapsed trees: 36,373 vs. 33,404 or 9%). Because the addition of non-coding sequences also
461 increased agreement between supermatrix and ASTRAL analyses (see above), this suggests that
462 at least some disagreement between supermatrix and species-tree analyses arises not only from
463 incomplete lineage sorting but also from lack of resolution at the gene tree level, something that
464 has also been observed at deeper phylogenetic scales (Pease et al. 2018).

465

466 *Phylogenetic relationships*

467 The genus *Artocarpus* was monophyletic in all 12 main analyses, as were subgenera
468 *Cauliflori* F.M. Jarrett and *Prainea* (King) Zerega, Supardi, & Motley (Table 1). Subgenus
469 *Artocarpus* was monophyletic excluding *A. sepicanus* Diels, and subgenus *Pseudojaca* was
470 monophyletic excluding *A. altissimus* J.J. Smith. In all supermatrix analyses and five ASTRAL
471 analyses, *Artocarpus sepicanus* Diels and *A. altissimus* J.J. Smith formed a clade sister to
472 subgenera *Cauliflori* and *Artocarpus*; however, in the codon-partitioned ASTRAL analyses, they
473 formed a grade in the same position, and in the ASTRAL supercontig analysis, only *A. altissimus*
474 was in that position, while *A. sepicanus* was sister to subgenus *Pseudojaca*. The backbone
475 phylogeny was otherwise identical in all twelve trees: subgenus *Prainea* was sister to all other
476 *Artocarpus*, which comprised a grade in this order: subgenus *Pseudojaca*, *A. sepicanus* + *A.*
477 *altissimus* (usually), followed by subgenus *Cauliflori* + subgenus *Artocarpus*. Apart from the

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

478 monophyly of the genus, which was supported by 61% of gene trees in the complete dataset
479 (supercontig, all genes), subgeneric relationships had much less support at the gene tree level.
480 The position of subg. *Prainea* was supported by 28% of gene trees; subg. *Pseudojaca* by 7%,
481 and subgenera *Artocarpus*/*Cauliflori* by only 4%. Quartet support, however, was substantially
482 higher (Figure 6).

483 Within subgenus *Artocarpus*, both of Jarrett's sections were monophyletic (leaving aside
484 *A. sepicanus*, *A. hirsutus*, and *A. nobilis*, which she considered anomalous and did not place in
485 sections), but none of the five series were monophyletic. The closest was series *Rugosi* F.M.
486 Jarrett which is characterized by rugose (sulcate to tuberculate) staminate inflorescences. It was
487 nearly monophyletic in most analyses, requiring only the exclusion of one rugose species (*A.*
488 *obtusus* F.M. Jarrett) and the inclusion of one non-rugose species (*A. teijsmannii* Miq.). Members
489 of series *Incisifolii* F.M. Jarrett, characterized by incised adult leaves (as with breadfruit, *A.*
490 *altilis*), formed two non-sister monophyletic clades, one in the Philippines and one ranging from
491 Indonesia to Oceania, both containing several potential undescribed species. Within subgenus
492 *Pseudojaca*, section *Pseudojaca* was monophyletic (excluding *A. altissimus*), as was series
493 *Clavati* F.M. Jarrett—characterized by clavate interfloral bracts. Series *Peltati* F.M. Jarrett—
494 characterized by peltate interfloral bracts—would be monophyletic if *A. tonkinensis* A. Chev.
495 were excluded, the latter species being sister to series *Clavati* in all main analyses.

496 Most species (for which we included at least two samples) were monophyletic as well,
497 but several were not monophyletic in any analysis, including *A. treculianus* Elmer, *A.*
498 *sarawakensis* F.M. Jarrett, *A. lanceifolius* Roxb., *A. rigidus* Blume, *A. teijsmannii*, and *A. nitidus*
499 Tréc. The type of *A. teijsmannii* ssp. *subglabrus* C.C. Berg was sister to *A. sepicanus* in all

500 analyses, while *ssp. teijsmannii* and other accessions of *ssp. subglabrus* were elsewhere within
501 subgenus *Artocarpus*.

502 The neotropical Artocarpeae formed a clade sister to *Artocarpus* in all twelve trees.

503 While *Batocarpus* was monophyletic in all supermatrix analyses, usually nested within a grade
504 comprising *Clarisia*, neither *Batocarpus* nor *Clarisia* was monophyletic in any ASTRAL tree.

505

506 **DISCUSSION**

507

508 *Taxon sampling*

509 Although other studies have successfully applied target enrichment to recover sequences
510 from herbarium and museum material (Guschanski et al. 2013; Hart et al. 2016), to our
511 knowledge, this is among the first to use herbarium collections to achieve near-complete taxon
512 sampling in a tropical plant genus of this size (ca. 70 spp.). The ability to successfully sequence
513 herbarium material was indispensable for this study. For 34 of 90 (38%) ingroup taxa remaining
514 in the final analyses (including subspecies and the nine individuals of uncertain affinities), we
515 did not have access to any fresh or silica-dried material and relied exclusively on herbarium
516 specimens. In some cases, the only readily available samples were approximately 100 years old,
517 as in the case of *A. treculianus* sensu stricto (coll. 1910–1911: 369–370 genes recovered after
518 filtering), *A. nigrescens* Elmer (coll. 1919: 431 genes), and *A. pinnatisectus* (type coll. 1913: 425
519 genes). Although old samples certainly had a lower success rate than silica-dried material, and
520 sample degradation undoubtedly contributed to shorter assembled contigs, age alone was not
521 significantly associated with recovery of fewer loci (Fig. 3). We hope these results serve as

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

522 encouragement for others to aim for complete taxon sampling with minimally-destructive
523 sampling from natural history collections when newly-collected material is not available.

524 While fieldwork remains among the most important aspects for systematic biology
525 studies, phylogenetic reconstruction can benefit dramatically from the incorporation of DNA
526 from museum specimens. In this study, we were able to sequence several DNA extractions that
527 had been prepared several years ago for previous studies but had been unusable, because PCR
528 amplification for Sanger sequencing failed (presumably due to small fragment size) (Zerega et al.
529 2010; Williams et al., 2017). The ability to achieve near-complete taxon sampling from museum
530 material will open up new opportunities for phylogeny-based analyses of clades with species that
531 are difficult to collect, rare, or extinct, but present in herbarium collections. Our results suggest
532 that near-complete taxon sampling can improve consistency between analyses, resulting in more
533 reliable phylogenies. A previous study (Kates et al. 2018) using a smaller dataset of 22
534 *Artocarpus* species, found substantial disagreement between analyses in the backbone phylogeny
535 of *Artocarpus*, in particular the positions of *Prainea* and *A. sepicanus*. Here, all 12 main analyses
536 recovered almost the same backbone, disagreeing occasionally as to the positions of *A. altissimus*
537 and *A. sepicanus*. Others have likewise found that missing taxa can substantially impact
538 phylogenetic reconstructions (de la Torre-Bárcena et al. 2009). Robust taxon sampling also has
539 serious implications for biodiversity conservation. *Artocarpus treculianus* is listed as Vulnerable
540 by the IUCN (World Conservation Monitoring Centre 1998). Due to the availability of sequences
541 from century-old herbarium sheets, we now know that this species is not monophyletic and that
542 the two obsolete taxa (*A. ovatifolius* Merr. and *A. nigrescens*, an unusual taxon with black fruits)
543 that Jarrett sunk into *A. treculianus* (Jarrett 1959) should probably be reinstated. Splitting a
544 Vulnerable species into three will result, at the very least, in three Vulnerable species. The

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

545 availability of material from collections has also revealed new species including *A. bergii* E.M.
546 Gardner, Zerega, and Arifiani (Gardner et al., in review), a close ally of breadfruit from the
547 Maluku Islands and *A. vietnamicus* E.M. Gardner and N.J.C. Zerega (Gardner and Zerega, in
548 review), a montane species endemic to Vietnam that resembles *A. excelsus* F.M. Jarrett and *A.*
549 *lowii* King (*A. aff. excelsus*).

550

551 *Impact of various analysis methods*

552 Of all the analytical variants we tested, partitioning the “exon” analyses by codon
553 position had the least impact, resulting in no major topological changes in any analysis. This is
554 not surprising, given that RAxML’s GTRCAT model provides for rate heterogeneity even absent
555 explicit partitioning (Stamatakis 2006). Comparisons of analyses with and without paralogs, with
556 and without non-coding sequences, and ASTRAL versus supermatrix revealed moderate
557 disagreement, mostly at shallow phylogenetic depths. However, in all cases, disagreement
558 decreased if additional sequences were added to a dataset (i.e., adding paralogs to the non-coding
559 present/absent comparison, adding non-coding sequences to the paralogs/no-paralogs
560 comparison, or adding either non-coding sequences or paralogs to the supermatrix/ASTRAL
561 comparison). This suggests that more data can lead to a certain amount of convergence in
562 analyses, even though simply adding more data to a supermatrix may not improve the accuracy
563 of the resulting species tree (Degnan and Rosenberg 2009).

564 Although adding or extending loci may reduce disagreement between analyses, it may
565 not always increase phylogenetic resolution. Certainly, with small numbers of genes, there may
566 not be enough informative characters to resolve a phylogeny, and resolution may increase as loci
567 are added. With hundreds of genes, however, lack of informative characters is not the problem.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

568 Here, although adding paralogs slightly reduced disagreement between analyses, overall percent
569 resolution remained static, and gene tree concordance with the species tree decreased. Other
570 phylogenomic studies have also found high rates of gene tree discordance (Degnan and
571 Rosenberg 2009; Wickett et al. 2014; Copetti et al. 2017; Pease et al. 2018; Liu et al. 2019).
572 Gene tree discordance results from, inter alia, biological processes such as incomplete lineage
573 sorting or ancient hybridization, and therefore may reflect not a lack of phylogenetic resolution,
574 but the non-existence of a fixed, absolute species tree. Nonetheless, just as bootstrap support can
575 convey a misleading sense of certainty, support measured by the rate of gene-tree support can
576 exaggerate uncertainty. For example, if a gene tree generally supports a clade, but has one out-
577 of-place taxon, perhaps due to an incomplete or erroneous sequence, that gene tree will not be
578 counted as supporting the clade in question. Support measure as the proportion of gene tree
579 quartets supporting each node, not the frequency of the exact clade being tested, may provide a
580 more realistic measure of support (Sayyari and Mirarab 2016); in our analyses, they were
581 generally lower than bootstrap values but substantially higher than gene-tree support.

582 Based on these results, we conclude that partitioning by codon position is not necessary
583 for our data, and analyses at similar phylogenetic scales may also not benefit from such a
584 partitioning scheme, at least for analyses that provide for built-in rate heterogeneity such as
585 RAxML. We do however recommend that when possible, flanking non-coding sequences be
586 included in analyses. The benefits of more informative gene trees, and thus more reliable final
587 analyses, outweigh any minimal advantage gained in partitioning by codon position, at least for a
588 data set like ours. Finally, in light of the increased congruence between analyses as our data set
589 was enlarged, we suggest using as many loci and as much flanking noncoding sequence as is
590 available, with the caveat to exercise caution with regard to taxa with excessive missing data.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

591 The cutoffs we used, >20% of the average sequence length and >~20% of loci, might be made
592 more stringent, as some inter-analysis disagreement appeared to center around samples with
593 more missing data.

594

595 *Taxonomic considerations*

596 Our results provide a phylogenetic framework for a taxonomic revision of *Artocarpus*,
597 currently in progress. The subgeneric divisions made by Jarrett (1959, 1960) and Zerega et al.
598 (2010) can be maintained with minor modifications to account for the anomalous *A. sepicanus*
599 and *A. altissimus*, which in 9/12 main analyses formed a clade. It is extremely curious that these
600 two species should be closely allied, and the occasional disagreement as to their affinity in the
601 ASTRAL analyses warrants further investigation. *Artocarpus sepicanus* is a fairly ordinary
602 member of subgenus *Artocarpus*, remarkable only in that its characters seem intermediate
603 between those defining section *Artocarpus* and section *Duricarpus* F.M. Jarrett (Jarrett 1959).
604 *Artocarpus altissimus*, on the other hand, is a most unusual member of subgenus *Pseudojaca*,
605 placed in that division solely on the basis of distichous leaves, and the fusion of perianth tissue of
606 adjacent carpellate flowers. Its leaves, with their trinerved bases and glandular-crenate margins,
607 are unique in the genus *Artocarpus* and are reminiscent of *Morus*, reflected in its basionym
608 *Morus altissima* Miq. Anatomical studies may reveal more, but the only apparent morphological
609 affinity between *A. sepicanus* and *A. altissimus* are bifid styles, a moraceous plesiomorphy
610 (Clement and Weiblen 2009) present occasionally in subgenus *Artocarpus* but unique to *A.*
611 *altissimus* in subgenus *Pseudojaca*. The leaf margins of *A. sepicanus*, which are entire in mature
612 trees but toothed in juveniles, are perhaps reminiscent of the glandular-crenate leaves of *A.*
613 *altissimus*.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

614 In addition, the phylogeny supports the broad outlines of Jarrett's (1959, 1960) sections,
615 validating her careful morphological and anatomical studies, which built on those of Renner
616 (1907). The sections within subgenera *Artocarpus* might be maintained with the minor
617 adjustment of including *A. hirsutus* Lam. and *A. nobilis* Thwaites in section *Duricarpus*. Jarrett
618 noted that those species had characters intermediate between sections *Artoarpus* and *Duriarpus*,
619 and indeed, their positions in all main analyses was that of sister to most of the rest of section
620 *Duricarpus*, possibly indicating the preservation of plesiomorphic characters. The two
621 subspecies of *A. lanceifolius* Roxb. were not sister taxa and may be distinguished by varying
622 stamen lengths and leaf sizes. In *A. sarawakensis*, the Sumatran individual was not sister to the
623 sample from Sarawak and the former might more appropriately be considered a variant of *A.*
624 *lanceifolius ssp. lanceifolius*, to which it bears similarity in syncarp characters, leaf shape,
625 differing only in the presence of a dense indumentum on the stipules.

626 At the series level within subgenus *Artocarpus*, a wholesale reconsideration is probably
627 necessary, although in a rough sense, the relevant characters match the clades. The exception is
628 series *Angusticarpus*, which did not form a consistent clade or grade. *Artoarpus teijsmannii*
629 belongs with the species of series *Rugosi*, which clade, in the broad sense, is characterized by
630 species with either rugose staminate inflorescences and/or dimorphic perianths (long+short) on
631 carpellate inflorescences, the latter of which applies to *A. teijsmannii* as well. The remaining
632 species, *A. lowii*, has morphological affinities to *A. excelsus* and *A. vietnamicus*. All three species
633 have smallish syncarps (to 6.5–7 cm), with shallow surface protrusions (except in *A.*
634 *vietnamicus*), and smallish elliptical leaves (6–36 cm long). These characters are closer to those
635 of *A. sepicanus* than most other members of section *Artocarpus*, which often have broadly ovate
636 to obovate leaves (up to ca. 70 cm long), suggesting that the three species, scattered throughout

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

637 section *Artocarpus*, may preserve some plesiomorphic characters. The two non-sister clades
638 comprising series *Incisifolii* do indeed contain all species with incised adult leaves, but many
639 others as well; they are defined more by geography than by morphology. The clade containing *A.*
640 *altilis* contains perhaps three new species, of interest as previously-unknown wild relatives of
641 breadfruit. These include *A. bergii*, known only from the Maluku Islands, and two accessions of
642 uncertain affinity from the living collections of Bogor Botanical Gardens (*cf. camansi* and *cf.*
643 *horridus*) that formed a clade yet have starkly divergent vegetative morphology, requiring further
644 study to determine whether they are distinct taxa. The status of *A. horridus*— morphologically
645 similar to *A. camansi*—is unclear; one accession fell in its expected place (sister to *A. camansi*),
646 but the position of the other, sister to the entire clade, must be treated with caution, as that
647 sample had among the highest proportions of missing data.

648 Within subgenus *Pseudojaca*, to the extent we included multiple accessions per species,
649 our results mostly supported Jarrett's (1960) revision. The series were largely monophyletic,
650 with the exception of the position of *A. tonkinensis* (with peltate interfloral bracts) nested within
651 the clade distinguished by clavate interfloral bracts. The ancestral state for interfloral bracts is
652 likely peltate (Clement and Weiblen 2009), so *A. tonkinensis* may simply represent a
653 plesiomorphic taxon sister to a derived clade. As Williams et al. (2017) found, the four species
654 sunk into *A. lacucha* by Berg et al. (2006) (*A. dadah* Miq., *A. ovatus* Blanco, *A. fretesii*, and *A.*
655 *vrieseanus* var. *refractus*) do not belong together. The varieties of *A. vrieseanus* do in fact form a
656 clade, and *A. longifolius* ssp. *adpressus* C.C. Berg, described by Berg (2005) ahead of the *Flora*
657 *Malesiana*, does indeed belong with the type subspecies. However, *A. gomezianus* ssp.
658 *zeylanicus* Jarrett was not sister to the type subspecies, but instead apparently belonged with *A.*
659 *lacucha* sensu Jarrett. *Artocarpus nitidus* ssp. *humilis* and ssp. *griffithii* usually formed a clade,

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

660 but otherwise the other two *A. nitidus* subspecies were unrelated. Among the Chinese species
661 described since Jarrett's (1960) revision, the rather distinctive *A. pithecogallus* C.Y. Wu and *A.*
662 *gongshanensis* S. K. Wu ex C. Y. Wu & S. S. Chang fell in the expected clavate-bracted clade.
663 Our sampling did not include *A. nanchuanensis*, but this species is morphologically similar to *A.*
664 *hypargyreus*, and subsequent sequencing after the main analyses were complete confirmed the
665 affinity. We were unable to successfully sequence *A. nigrifolius*, but an examination of type
666 specimen confirmed a close affinity to *A. hypargyreus*. Although the original description
667 mentions possible affinities with *A. styracifolius*, the clavate rather than ovoid staminate
668 inflorescences are much closer to *A. hypargyreus*. Moreover, the primary distinguishing
669 characters—leaves drying black—is frequently found on *A. hypargyreus* herbarium sheets as
670 well, including the type of the latter species (*Hance 4484*, P) and might therefore not be a proper
671 diagnostic.

672

673 CONCLUSION

674 The increasing availability of phylogenomic datasets has dramatically changed the
675 practice of revisionary systematics. Data sets containing hundreds or thousands of loci produce
676 trees with extremely high statistical support, apparently providing ironclad frameworks for
677 making taxonomic decisions. However, apparent high support for relationships may often be an
678 artifact of the massive number of characters available for phylogenetic inference, masking real
679 uncertainties that are revealed only by employing a variety of analytical methods. By the same
680 token, focusing on exclusively conserved coding regions—an inherent feature of some reference-
681 based assembly methods—can result in unnecessarily uninformative gene trees, leading to poor
682 support at the species tree level. Using a data set with near-complete taxon sampling, we

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

683 demonstrated that decisions made in how to conduct analyses can substantially affect
684 phylogenetic reconstruction, resulting in discordant phylogenies, each with high statistical
685 support. Employing multiple analytical methods can help separate truly robust phylogenetic
686 relationships from those that only appear to be well-supported but are not consistent across
687 analyses. While codon partitioning and model choice did not substantially alter our phylogeny,
688 the inclusion of flanking non-coding sequences in analyses significantly increased the number of
689 informative splits at the gene tree level, resulting ultimately in more robust species trees. In
690 general, increasing the size of datasets through the inclusion of paralogous genes increased
691 convergence between analysis methods but did not reduce gene tree conflict, which likely
692 resulted from biological and not analytical processes; for this reason, we prefer quartet-based
693 scoring methods as the most informative ways of determining support for species trees.

694

695 We provide a robust phylogenetic framework for *Artocarpus*, making use of herbarium
696 specimens up to 106 years old to supplement our own collections and achieve near-complete
697 taxon sampling, demonstrating the value of even very old natural history collections in
698 improving phylogenetic studies. Our results will inform future evolutionary and systematic
699 studies of this important group of plants. More generally, the results may guide future analyses
700 of HybSeq datasets, particularly those combining fresh with museum material, by counseling
701 careful attention to dataset construction and analysis method to produce the most informative
702 phylogenetic hypotheses.

703

704 FUNDING

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

705 This work was supported by the United States National Science Foundation (DEB award
706 numbers 0919119 to NJCZ, 1501373 to NJCZ and EMG, 1342873 and 1239992 to NJW, and
707 DBI award number 1711391 to EMG); the Northwestern University Plant Biology and
708 Conservation Program; The Initiative for Sustainability and Energy at Northwestern University;
709 the Garden Club of America; the American Society of Plant Taxonomists; a Systematics
710 Research Fund grant from the Linnaean Society and the Systematics Association; and the
711 Botanical Society of America.

712

713 ACKNOWLEDGEMENTS

714 We thank Postar Miun, Jeisin Jumian, Markus Gubilil, Aloysius Laim, S. Brono, Jegong anak
715 Suka, Salang anak Nyegang, Jugah anak Tagi, Wan Nuur Fatiha Wan Zakaria, and Harto for
716 assistance in the field; the Pritzker Laboratory for Molecular Systematics at the Field Museum of
717 Natural History (K. Feldheim) for the use of sequencing facilities; J. Fant, E. Williams, H.
718 Noble, R. Overson, and B. Cooper for assistance in the lab; the Sabah Agriculture Department
719 for access to collections; and the following herbaria for access to collections for examination
720 and/or sampling: BM, BO, F, IBSC, K, FRIM, FTBG, L, MO, NY, PNH, SAN, SING, SAR,
721 SNP.

722

723 LITERATURE CITED

724

725 van Bakel H., Stout J.M., Cote A.G., Tallon C.M., Sharpe A.G., Hughes T.R., Page J.E. 2011.

726 The draft genome and transcriptome of *Cannabis sativa*. *Genome Biol.* 12:R102.

727 Bakker F.T., Lei D., Yu J., Mohammadin S., Wei Z., van de Kerke S., Gravendeel B.,

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 728 Nieuwenhuis M., Staats M., Alquezar-Planas D.E., Holmer R. 2016. Herbarium genomics:
729 Plastome sequence assembly from a range of herbarium specimens using an Iterative
730 Organelle Genome Assembly pipeline. *Biol. J. Linn. Soc.* 117:33–43.
- 731 Bankevich A., Nurk S., Antipov D., Gurevich A.A., Dvorkin M., Kulikov A.S., Lesin V.M.,
732 Nikolenko S.I., Pham S., Prjibelski A.D., Pyshkin A. V., Sirotkin A. V., Vyahhi N., Tesler
733 G., Alekseyev M.A., Pevzner P.A. 2012. SPAdes: A New Genome Assembly Algorithm
734 and Its Applications to Single-Cell Sequencing. *J. Comput. Biol.* 19:455–477.
- 735 Berg C. 2005. Flora Malesiana precursor for the treatment of Moraceae 8: other genera than
736 *Ficus*. *Blumea*. 50:535–550.
- 737 Brewer G.E., Clarkson J.J., Maurin O., Zuntini A.R., Barber V., Bellot S., Biggs N., Cowan R.S.,
738 Davies N.M.J., Dodsworth S., others. 2019. Factors affecting targeted sequencing of 353
739 nuclear genes from herbarium specimens spanning the diversity of angiosperms. *Front.*
740 *Plant Sci.* 10:1102.
- 741 Broad Institute. 2016. Picard tools. Available from
742 <https://broadinstitute.github.io/picard/>
<http://broadinstitute.github.io/picard/>.
- 743 Buerki S., Baker W.J. 2016. Collections-based research in the genomic era. *Biol. J. Linn. Soc.*
744 117:5–10.
- 745 Capella-Gutiérrez S., Silla-Martínez J.M., Gabaldón T. 2009. trimAl: A tool for automated
746 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*. 25:1972–1973.
- 747 Castañeda-Álvarez N.P., Khoury C.K., Achicanoy H.A., Bernau V., Dempewolf H., Eastwood
748 R.J., Guarino L., Harker R.H., Jarvis A., Maxted N., Müller J. V., Ramirez-Villegas J., Sosa
749 C.C., Struik P.C., Vincent H., Toll J. 2016. Global conservation priorities for crop wild
750 relatives. *Nat. Plants*.:1–6.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 751 Chamala S., García N., Godden G.T., Krishnakumar V., Jordon-Thaden I.E., De Smet R.,
752 Barbazuk W.B., Soltis D.E., Soltis P.S. 2015. MarkerMiner 1.0: A new application for
753 phylogenetic marker development using angiosperm transcriptomes. *Appl. Plant Sci.*
754 3:1400115.
- 755 Clement W.L., Weiblen G.D. 2009. Morphological Evolution in the Mulberry Family
756 (Moraceae). *Syst. Bot.* 34:530–552.
- 757 Copetti D., Búrquez A., Bustamante E., Charboneau J.L.M., Childs K.L., Eguiarte L.E., Lee S.,
758 Liu T.L., McMahon M.M., Whiteman N.K., Wing R.A., Wojciechowski M.F., Sanderson
759 M.J. 2017. Extensive gene tree discordance and hemiplasy shaped the genomes of North
760 American columnar cacti. *Proc. Natl. Acad. Sci. U. S. A.*
- 761 Degnan J.H., Rosenberg N.A. 2009. Gene tree discordance, phylogenetic inference and the
762 multispecies coalescent. *Trends Ecol. Evol.*
- 763 Doyle J., Doyle J. 1987. Genomic plant DNA preparation from fresh tissue-CTAB method.
764 *Phytochem. Bull.* 19:11–15.
- 765 Faircloth B.C. 2015. PHYLUCE is a software package for the analysis of conserved genomic
766 loci. *Bioinformatics.* 32:786–788.
- 767 Faircloth B.C., McCormack J.E., Crawford N.G., Harvey M.G., Brumfield R.T., Glenn T.C.
768 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple
769 evolutionary timescales. *Syst. Biol.* 61:717–26.
- 770 Gardner E.M., Arifiani D., Zerega N.J.C. In review. *Artocarpus bergii* (Moraceae): a new
771 species in the breadfruit clade from the Moluccas. *Syst. Bot.*
- 772 Gardner E.M., Johnson M.G., Ragone D., Wickett N.J., Zerega N.J.C. 2016. Low-Coverage,
773 Whole-Genome Sequencing of *Artocarpus camansi* (Moraceae) for Phylogenetic Marker

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 774 Development and Gene Discovery. Appl. Plant Sci. 4:1600017.
- 775 Gardner E.M., Zerega N.J.C. In review. *Artocarpus vietnamicus* (Moraceae): a new species from
776 the mountain forests of Vietnam. Phytotaxa.
- 777 Guschanski K., Krause J., Sawyer S., Valente L.M., Bailey S., Finstermeier K., Sabin R.,
778 Gilissen E., Sonet G., Nagy Z.T., Lenglet G., Mayer F., Savolainen V. 2013. Next-
779 generation museomics disentangles one of the largest primate radiations. Syst. Biol.
780 62:539–554.
- 781 Hart M.L., Forrest L.L., Nicholls J.A., Kidner C.A. 2016. Retrieval of hundreds of nuclear loci
782 from herbarium specimens. Taxon. 65:1081–1092.
- 783 He N., Zhang C., Qi X., Zhao S., Tao Y. 2013. Draft genome sequence of the mulberry tree
784 *Morus notabilis*. Nat. Commun. 4:2445.
- 785 Jarrett F.M. 1975. Four new *Artocarpus* species from Indo-Malesia (Moraceae). Blumea.
786 22:409–410.
- 787 Johnson M.G., Gardner E.M., Liu Y., Medina R., Goffinet B., Shaw A.J., Zerega N.J.C., Wickett
788 N.J. 2016. HybPiper: Extracting Coding Sequence and Introns for Phylogenetics from High-
789 Throughput Sequencing Reads Using Target Enrichment. Appl. Plant Sci. 4:1600016.
- 790 Johnson M.G., Pokorny L., Dodsworth S., Botigué L.R., Cowan R.S., Devault A., Eiserhardt
791 W.L., Epiawalage N., Forest F., Kim J.T., Leebens-Mack J.H., Leitch I.J., Maurin O.,
792 Soltis D.E., Soltis P.S., Wong G.K., Baker W.J., Wickett N.J. 2019. A Universal Probe Set
793 for Targeted Sequencing of 353 Nuclear Genes from Any Flowering Plant Designed Using
794 k-Medoids Clustering. Syst. Biol. 68:594–606.
- 795 Kates H.R., Johnson M.G., Gardner E.M., Zerega N.J.C., Wickett N.J. 2018. Allele phasing has
796 minimal impact on phylogenetic reconstruction from targeted nuclear gene sequences in a

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 797 case study of *Artocarpus*. *Am. J. Bot.*
- 798 Katoh K., Standley D.M. 2013. MAFFT multiple sequence alignment software version 7:
799 improvements in performance and usability. *Mol. Biol. Evol.* 30:772–80.
- 800 Keller O., Kollmar M., Stanke M., Waack S. 2011. A novel hybrid gene prediction method
801 employing protein multiple sequence alignments. *Bioinformatics.* 27:757–63.
- 802 Kochummen K.M. 1998. New species and varieties of Moraceae from Malaysia. *Gard. Bull.*
803 Singapore. 50:197–219.
- 804 Kubatko L.S., Degnan J.H. 2007. Inconsistency of phylogenetic estimates from concatenated
805 data under coalescence. *Syst. Biol.* 56:17–24.
- 806 de la Torre-Bárcena J.E., Kolokotronis S.-O., Lee E.K., Stevenson D.W., Brenner E.D., Katari
807 M.S., Coruzzi G.M., DeSalle R. 2009. The impact of outgroup choice and missing data on
808 major seed plant phylogenetics using genome-wide EST data. *PLoS One.* 4:e5764.
- 809 Lanfear R., Calcott B., Kainer D., Mayer C., Stamatakis A. 2014. Selecting optimal partitioning
810 schemes for phylogenomic datasets. *BMC Evol. Biol.* 14:82.
- 811 Larsson A. 2014. AliView: A fast and lightweight alignment viewer and editor for large datasets.
812 *Bioinformatics.* 30:3276–3278.
- 813 Lee B.N. 2014. Solid Phase Reverse Immobilization (SPRI) Bead Technology for Micro RNA
814 Clean Up using the . .
- 815 Li H., Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
816 *Bioinformatics.* 25:1754–1760.
- 817 Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin
818 R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics.* 25:2078–
819 2079.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 820 Liu Y., Johnson M.G., Cox C.J., Medina R., Devos N., Vanderpoorten A., Hedenäs L., Bell N.E.,
821 Shevock J.R., Aguero B., Quandt D., Wickett N.J., Shaw A.J., Goffinet B. 2019. Resolution
822 of the ordinal phylogeny of mosses using targeted exons from organellar and nuclear
823 genomes. *Nat. Commun.*
- 824 Mandel J.R., Dikow R.B., Funk V. a., Masalia R.R., Staton S.E., Kozik A., Michelmore R.W.,
825 Rieseberg L.H., Burke J.M. 2014. A Target Enrichment Method for Gathering Phylogenetic
826 Information from Hundreds of Loci: An Example from the Compositae. *Appl. Plant Sci.*
827 2:1300085.
- 828 Medina R., Johnson M.G., Liu Y., Wickett N.J., Shaw A.J., Goffinet B. 2019. Phylogenomic
829 delineation of *Physcomitrium* (Bryophyta: Funariaceae) based on targeted sequencing of
830 nuclear exons and their flanking regions rejects the retention of *Physcomitrella* ,
831 *Physcomitridium* and *Aphanorrhagma* . *J. Syst. Evol.*
- 832 Miller M.A., Pfeiffer W., Schwartz T. 2010. Creating the CIPRES Science Gateway for
833 inference of large phylogenetic trees. 2010 Gatew. Comput. Environ. Work. GCE 2010.
- 834 Mirarab S., Warnow T. 2015. ASTRAL-II: coalescent-based species tree estimation with many
835 hundreds of taxa and thousands of genes. *Bioinformatics.* 31:i44–i52.
- 836 Paradis E., Claude J., Strimmer K. 2004. APE: Analyses of phylogenetics and evolution in R
837 language. *Bioinformatics.* 20:289–290.
- 838 Pease J.B., Brown J.W., Walker J.F., Hinchliff C.E., Smith S.A. 2018. Quartet Sampling
839 distinguishes lack of support from conflicting support in the green plant tree of life. *Am. J.*
840 *Bot.*
- 841 Price M.N., Dehal P.S., Arkin A.P. 2009. Fasttree: Computing large minimum evolution trees
842 with profiles instead of a distance matrix. *Mol. Biol. Evol.* 26:1641–1650.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 843 Prum R.O., Berv J.S., Dornburg A., Field D.J., Townsend J.P., Lemmon E.M., Lemmon A.R.
844 2015. A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA
845 sequencing. *Nature*. 526:569–573.
- 846 Quinlan A.R., Hall I.M. 2010. BEDTools: A flexible suite of utilities for comparing genomic
847 features. *Bioinformatics*. 26:841–842.
- 848 Rambaut A. 2016. FigTree v1.4.3. *Mol. Evol. phylogenetics Epidemiol.*
- 849 Ranwez V. 2011. MACSE : Multiple Alignment of Coding SEquences Accounting for
850 Frameshifts and Stop Codons. 6.
- 851 Revell L.J. 2012. phytools: An R package for phylogenetic comparative biology (and other
852 things). *Methods Ecol. Evol.* 3:217–223.
- 853 Sarkar D. 2008. *Lattice: Multivariate Data Visualization with R*. New York: Springer.
- 854 Sayyari E., Mirarab S. 2016. Fast Coalescent-Based Computation of Local Branch Support from
855 Quartet Frequencies. *Mol. Biol. Evol.* 33:1654–1668.
- 856 Sayyari E., Whitfield J.B., Mirarab S. 2017. Fragmentary Gene Sequences Negatively Impact
857 Gene Tree and Species Tree Reconstruction. *Mol. Biol. Evol.* 34:3279–3291.
- 858 Schliep K.P. 2011. phangorn: Phylogenetic analysis in R. *Bioinformatics*. 27:592–593.
- 859 Smith S.A., Moore M.J., Brown J.W., Yang Y. 2015. Analysis of phylogenomic datasets reveals
860 conflict, concordance, and gene duplications with examples from animals and plants. *BMC*
861 *Evol. Biol.*
- 862 Staats M., Erkens R.H.J., van de Vossen B., Wieringa J.J., Kraaijeveld K., Stielow B., Geml
863 J., Richardson J.E., Bakker F.T. 2013. Genomic treasure troves: complete genome
864 sequencing of herbarium and insect museum specimens. *PLoS One*. 8:e69189.
- 865 Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 866 thousands of taxa and mixed models. *Bioinformatics*. 22:2688–90.
- 867 Sukumaran J., Holder M.T. 2010. DendroPy: a Python library for phylogenetic computing.
868 *Bioinformatics*. 26:1569–71.
- 869 Tange O. 2018. GNU Parallel 2018. <https://doi.org/10.5281/zenodo.1146014>.
- 870 Team R.D.C. 2008. R: A language and environment for statistical computing. Vienna, Austria: R
871 Foundation for Statistical Computing.
- 872 Villaverde T., Pokorny L., Olsson S., Rincón-Barrado M., Johnson M.G., Gardner E.M., Wickett
873 N.J., Molero J., Riina R., Sanmartín I. 2018. Bridging the micro- and macroevolutionary
874 levels in phylogenomics: Hyb-Seq solves relationships from populations to species and
875 above. *New Phytol*.
- 876 Wang M.M.H., Gardner E.M., Chung R.C.K., Chew M.Y., Milan A.R., Pereira J.T., Zerega
877 N.J.C. 2018. Origin and diversity of an underutilized fruit tree crop, cempedak (*Artocarpus*
878 *integer*, Moraceae). *Am. J. Bot.*
- 879 Weitemier K., Straub S.C.K., Cronn R.C., Fishbein M., Schmickl R., McDonnell A., Liston A.
880 2014. Hyb-Seq: Combining Target Enrichment and Genome Skimming for Plant
881 Phylogenomics. *Appl. Plant Sci.* 2:1400042.
- 882 Wickett N.J., Mirarab S., Nguyen N., Warnow T., Carpenter E., Matasci N., Ayyampalayam S.,
883 Barker M.S., Burleigh J.G., Gitzendanner M.A., Ruhfel B.R., Wafula E., Der J.P., Graham
884 S.W., Mathews S., Melkonian M., Soltis D.E., Soltis P.S., Miles N.W., Rothfels C.J.,
885 Pokorny L., Shaw A.J., DeGironimo L., Stevenson D.W., Surek B., Villarreal J.C., Roure
886 B., Philippe H., DePamphilis C.W., Chen T., Deyholos M.K., Baucom R.S., Kutchan T.M.,
887 Augustin M.M., Wang J., Zhang Y., Tian Z., Yan Z., Wu X., Sun X., Wong G.K.-S.,
888 Leebens-Mack J. 2014. Phylotranscriptomic analysis of the origin and early diversification

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

- 889 of land plants. Proc. Natl. Acad. Sci. 111:E4859–E4868.
- 890 Williams E.W., Gardner E.M., Harris R., Chaveerach A., Pereira J.T., Zerega N.J.C. 2017. Out
891 of Borneo: Biogeography, phylogeny, and divergence date estimates of *Artocarpus*
892 (Moraceae). Ann. Bot. 119:611–627.
- 893 Witherup C., Zuberi M.I., Hossain S., Zerega N.J.C. 2019. Genetic Diversity of Bangladeshi
894 Jackfruit (*Artocarpus heterophyllus*) over Time and Across Seedling Sources. Econ. Bot.
895 73:233–248.
- 896 World Conservation Monitoring Centre. 1998. *Artocarpus treculianus*. Available from
897 <http://dx.doi.org/10.2305/IUCN.UK.1998.RLTS.T33246A97711111.en>.
- 898 Xi Z., Ruhfel B.R., Schaefer H., Amorim A.M., Sugumaran M., Wurdack K.J., Endress P.K.,
899 Matthews M.L., Stevens P.F., Mathews S., Davis C.C. 2012. Phylogenomics and a
900 posteriori data partitioning resolve the Cretaceous angiosperm radiation Malpighiales. Proc.
901 Natl. Acad. Sci. U. S. A. 109:17519–24.
- 902 Zerega N.J.C., Nur Supardi M.N., Motley T.J. 2010. Phylogeny and Recircumscription of
903 *Artocarpeae* (Moraceae) with a Focus on *Artocarpus*. Syst. Bot. 35:766–782.
- 904 Zerega N.J.C., Wiesner-Hanks T., Ragone D., Irish B., Scheffler B., Simpson S., Zee F. 2015.
905 Diversity in the breadfruit complex (*Artocarpus*, Moraceae): genetic characterization of
906 critical germplasm. Tree Genet. Genomes. 11:1–26.
- 907 Zhang C., Sayyari E., Mirarab S. 2017. ASTRAL-III: Increased Scalability and Impacts of
908 Contracting Low Support Branches BT - Comparative Genomics: 15th International
909 Workshop, RECOMB CG 2017, Barcelona, Spain, October 4-6, 2017, Proceedings. In:
910 Meidanis J., Nakhleh L., editors. Cham: Springer International Publishing. p. 53–75.
- 911 Zhengyi W., Xiushi Z. 1989. TAXA NOVA NONNULLA MORACEARUM SINENSIIUM.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

912 Acta Bot. Yunnanica. 11:24–34.

913

914

915

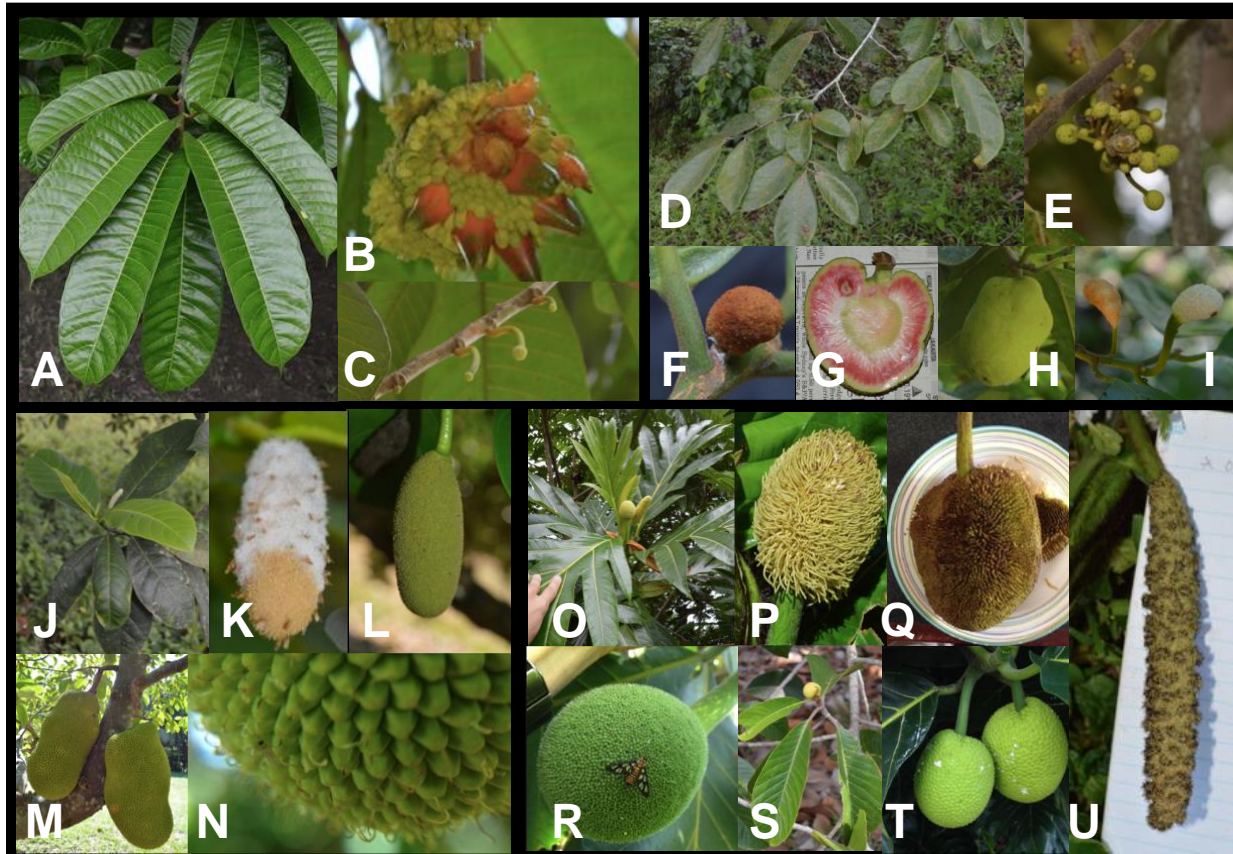
GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

Table 1. A summary of *Artocarpus* taxonomy following Zerega et al. (2010) at the subgeneric level, and Jarrett (1959–1960) at the section and series level. Species marked with an asterisk (*) were described after Jarrett's revision; we have listed them with the taxonomic divisions in which they place based on the phylogeny presented in this study. Species in gray text were not included in the phylogeny.

Subgenus	Section	Series	Species	Monophyletic
<i>Artocarpus</i>				Yes, if <i>A. sepicanus</i> is excluded
	<i>Artocarpus</i>			Yes, if <i>A. sepicanus</i> is excluded
		<i>Angusticarpa</i>	<i>A. lowii</i> , <i>A. teijsmannii</i>	Yes, if <i>A. teijsmannii</i> ssp. <i>subglabrous</i> is excluded, but see also series <i>Rugosi</i> for comments
		<i>Incisifolii</i>	[<i>A. altitilis</i> , <i>A. camansi</i> , <i>A. mariannensis</i> , <i>A. horridus</i> , <i>A. bergii</i> *], [<i>A. blancoi</i> , <i>A. treculianus</i> , <i>A. pinnatisectus</i> , <i>A. multifidus</i>]	No, but it consists of two monophyletic clades (separated by brackets to the left) defined by geography
		<i>Rugosi</i>	<i>A. scortechinii</i> , <i>A. elasticus</i> , <i>A. sericicarpus</i> , <i>A. tamaran</i> , <i>A. sumatranus</i> , <i>A. kemando</i> , <i>A. maingayi</i> , <i>A. corneri</i> *, <i>A. jarrettiae</i> *, <i>A. excelsus</i> *, <i>A. obtusus</i> *	In most analyses, yes if <i>A. lowii</i> , and <i>A. teijsmannii</i> ssp. <i>teijsmannii</i> are included
		Unplaced	<i>A. hirsutus</i> , <i>A. nobilis</i> , <i>A. sepicanus</i> , <i>A. vietnamicus</i> *	
	<i>Duricarpus</i>			Yes, if <i>A. hirsutus</i> and <i>A. nobilis</i> are included
		<i>Asperifolii</i>	<i>A. melinocylus</i> , <i>A. odoratissimus</i> , <i>A. hispidus</i> , <i>A. rigidus</i> , <i>A. chama</i> , <i>A. brevipedunculatus</i> , <i>A. sarawakensis</i> *	Yes, if <i>A. hirsutus</i> and <i>A. nobilis</i> are included, and <i>A. sarawakensis</i> and <i>A. brevipedunculatus</i> are excluded
		<i>Laevifolii</i>	<i>A. anisophyllus</i> , <i>A. lanceifolius</i>	Yes, if <i>A. sarawakensis</i> and <i>A. brevipedunculatus</i> are included
<i>Canliflori</i>			<i>A. heterophyllus</i> , <i>A. integer</i> , <i>A. annulatus</i> *	Yes
<i>Pseudojaca</i>				Yes, if <i>A. altissimus</i> is excluded.
	<i>Glandulifolium</i>		<i>A. altissimus</i>	Yes
	<i>Pseudojaca</i>	<i>Clavati</i>	<i>A. hypargyraeus</i> , <i>A. styracifolius</i> , <i>A. petalotii</i> , <i>A. pitheogallus</i> *, <i>A. gongshanensis</i> *, <i>A. nigrifolius</i> *, <i>A. nanchuanensis</i> *	Yes, if <i>A. tonkinensis</i> is included.
		<i>Peltati</i>	<i>A. glaucus</i> , <i>A. vrieseanus</i> , <i>A. xanthocarpus</i> , <i>A. longifolius</i> , <i>A. subtundifolius</i> , <i>A. reticulatus</i> , <i>A. lacucha</i> , <i>A. gomezianus</i> , <i>A. tomentosulus</i> , <i>A. ovatus</i> , <i>A. tonkinensis</i> , <i>A. fretessii</i> , <i>A. dadah</i> , <i>A. rubrovenius</i> , <i>A. nitidus</i> , <i>A. fulvicortex</i> , <i>A. lacucha</i> , <i>A. albobrunneus</i> *, <i>A. thailandicus</i> *, <i>A. primackii</i> *	Yes, if <i>A. tonkinensis</i> is excluded
<i>Prainea</i>			<i>A. limpato</i> , <i>A. papuanus</i> , <i>A. scandens</i> , <i>A. frutescens</i>	Yes

916
917

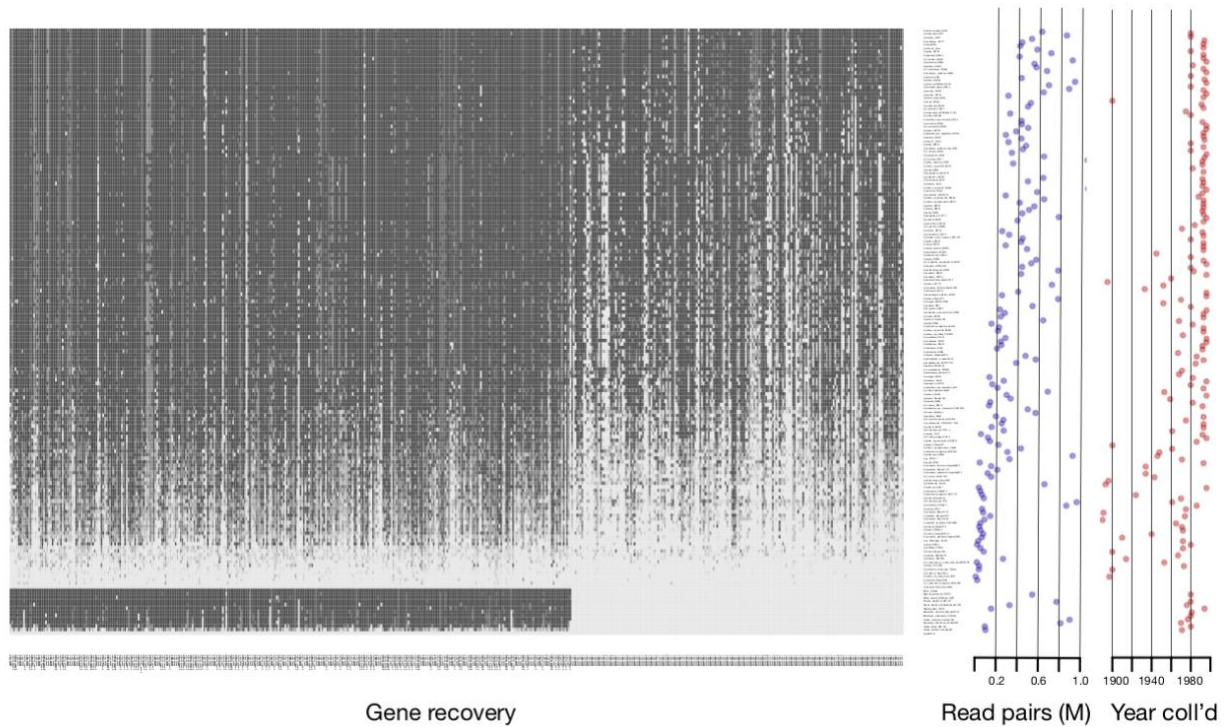
918 FIGURES



919

920 Figure 1. Diversity of *Artocarpus*. **Subg. Prainea**— (A) leaves, (B) syncarp, and (C) immature
921 inflorescences of *A. limpatum*. **Subg. Pseudojaca** — (D) leaves and (E) staminate
922 inflorescences of *A. fretessii*; (F) carpellate inflorescence of *A. nitidus ssp. borneensis*; (G)
923 syncarp of *A. primackii*; (H) syncarp of *A. nitidus ssp. lingnanensis*; and (I) staminate (left)
924 and carpellate (right) inflorescences of *A. hypargyreus*. **Subg. Cauliflori**— (J) leaves of *A.*
925 *heterophyllus*; (K–L) staminate inflorescences, (M) syncarp, and (N) carpellate inflorescence of
926 *A. heterophyllus*. **Subg. Artocarpus**— (O) leaves and inflorescences of *A. altilis*; (P)
927 carpellate inflorescence of *A. tamaran*; (Q) syncarp and (R) carpellate inflorescence of *A.*
928 *odoratissimus*; (S) leaves and staminate inflorescence of *A. rigidus*; (T) syncarp of *A.*
929 *altilis*; and (U) staminate inflorescence of *A. tamaran*.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*



930

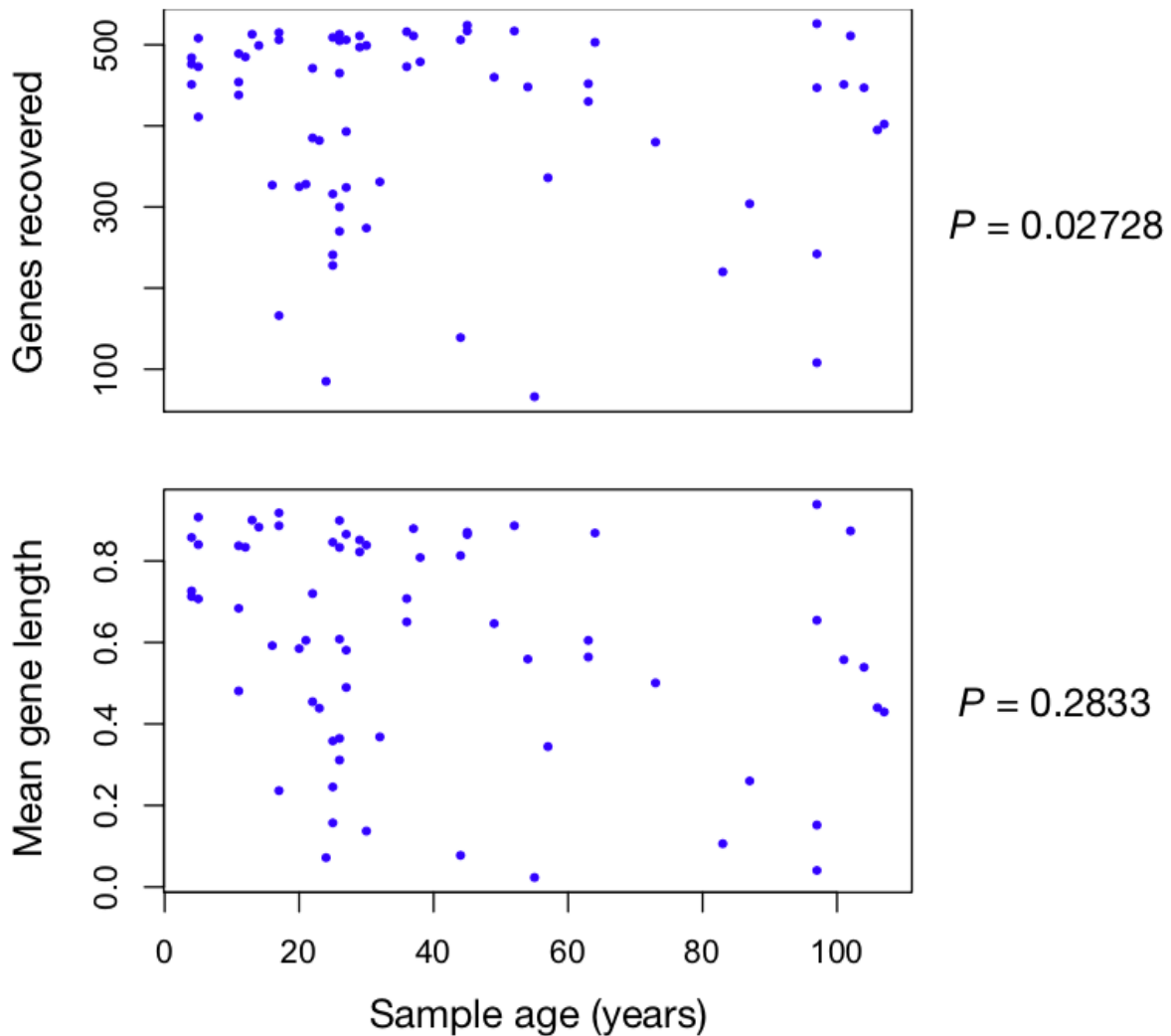
931 Figure 2. (A) Heatmap of gene recovery as a percentage of the average recovered sequence

932 length for each gene. Rows represent samples, and columns represent genes. Darker

933 colors indicate more complete recovery; white indicates no recovery. (B) Age (x-axis) for each

934 sample. (C) Number of reads on target (x-axis) for each sample.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

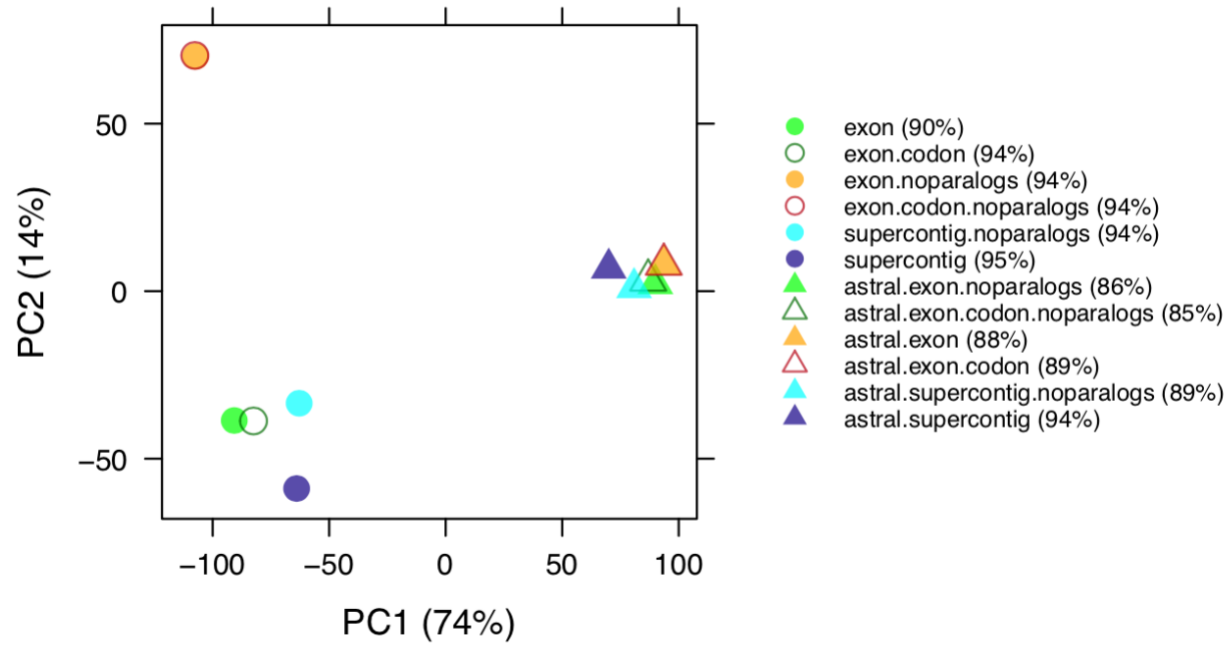


935

936 Figure 3. Comparison between herbarium sample age and total genes recovered (top) and mean

937 gene length as a proportion of average gene length for each gene (bottom)

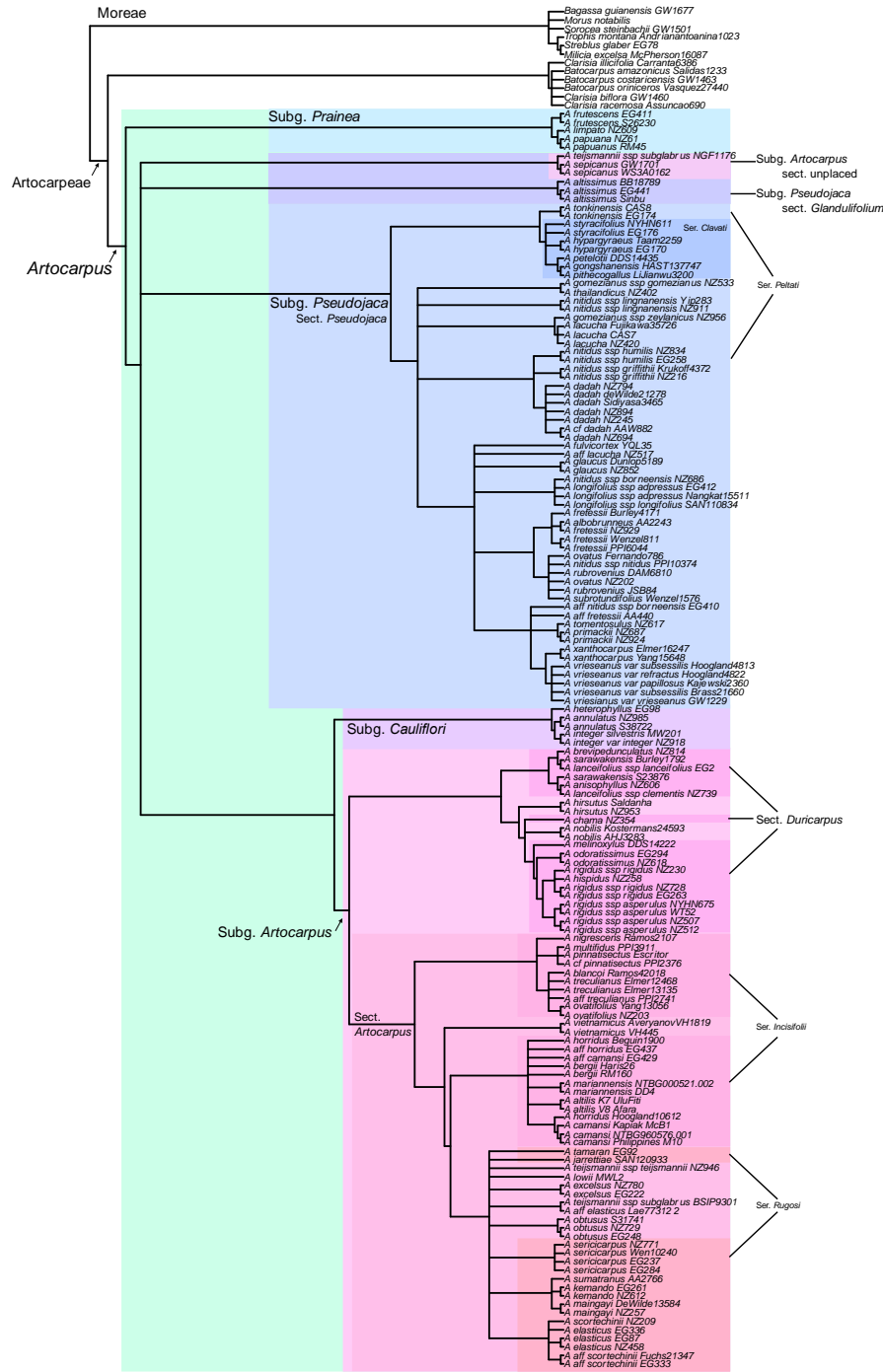
GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*



938

939 Figure 4. PCA of Robinson-Foulds (RF) distances between all 12 main analyses.

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*



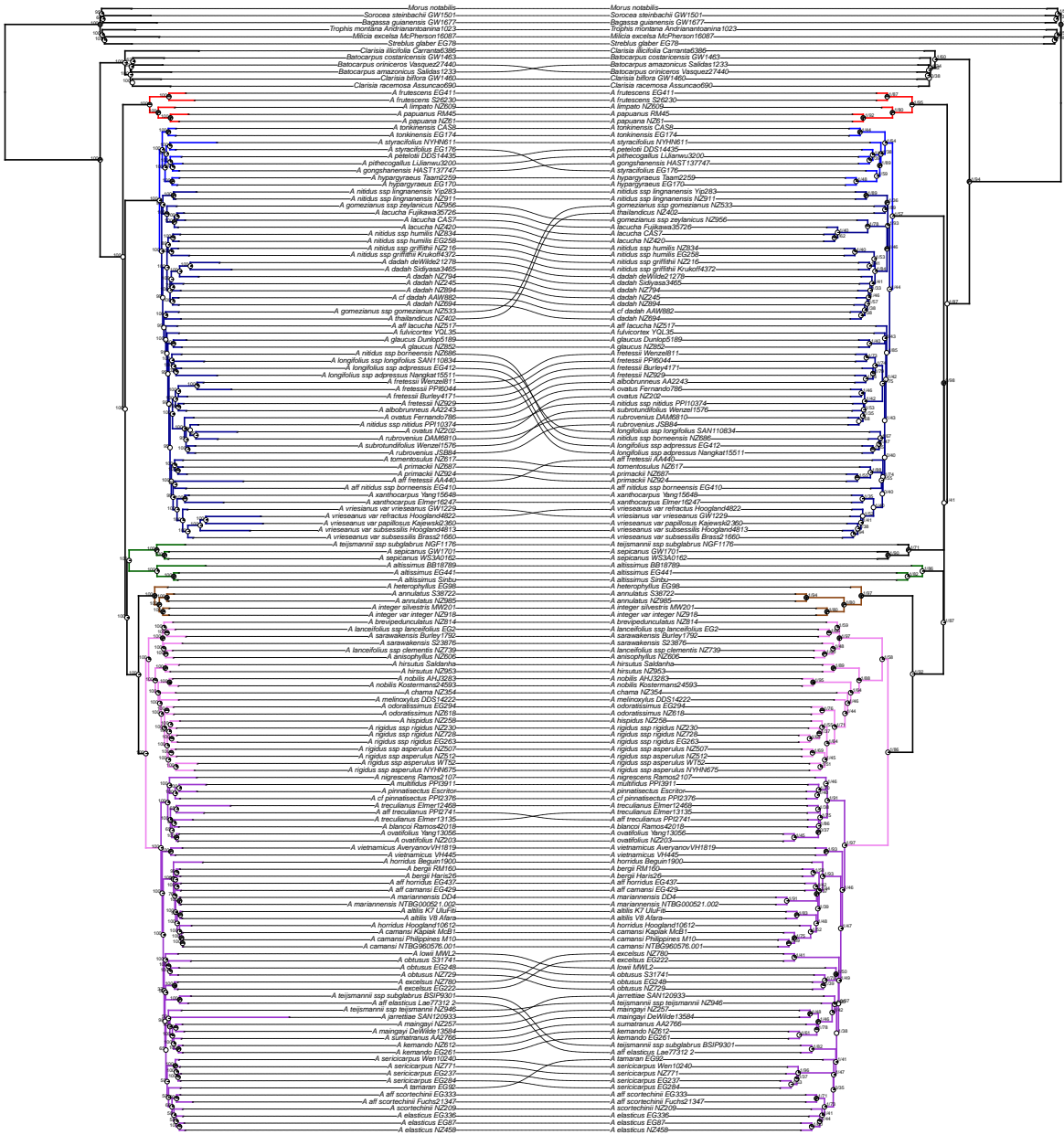
940

941 Figure 5. Strict consensus of all 12 main-analysis trees (excluding only those analyses in which

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

942 exons and introns were aligned separately, for reasons discussed in the text). Colored boxes
943 reflect Jarrett's (1959, 1960) taxonomic divisions, as modified by Zerega et al. (2010).
944 Recently-described taxa that were split from older taxa recognized by Jarrett are classified
945 according to Jarrett's species concepts. Labels to the right of the tree denote major non-
946 monophyletic taxonomic divisions
947

GARDNER ET AL., PHYLOGENOMICS OF ARTOCARPUS



948

949 Figure 6. Comparison between the full-dataset (supercontigs for all genes) supermatrix and

950 ASTRAL trees, showing moderate disagreement at shallow phylogenetic depths but

951 complete agreement at deeper nodes. Left: maximum-likelihood tree based on all

952 supercontigs, partitioned by gene, including all paralogs; all branch lengths are proportional

953 to mean substitutions per site. Right: ASTRAL tree based on all supercontigs; internal

GARDNER ET AL., PHYLOGENOMICS OF *ARTOCARPUS*

954 branch lengths are proportional to coalescent units; terminal branch lengths were arbitrarily
955 assigned to improve visualization. Pie charts at nodes represent the proportion of gene trees
956 supporting each split, and numbers represent bootstrap support
957