

1 **A Comprehensive Phylogenomic Platform for Exploring the Angiosperm Tree of Life**

2

3 William J. Baker^{1,*}, Paul Bailey¹, Vanessa Barber¹, Abigail Barker¹, Sidonie Bellot¹, David
4 Bishop¹, Laura R. Botigué^{1,2}, Grace Brewer¹, Tom Carruthers¹, James J. Clarkson¹, Jeffrey
5 Cook¹, Robyn S. Cowan¹, Steven Dodsworth^{1,3}, Nirosini Epiawalage¹, Elaine Françoso¹,
6 Berta Gallego¹, Matthew G. Johnson⁴, Jan T. Kim^{1,5}, Kevin Leempoel¹, Olivier Maurin¹,
7 Catherine McGinnie¹, Lisa Pokorny^{1,6}, Shyamali Roy¹, Malcolm Stone¹, Eduardo Toledo¹,
8 Norman J. Wickett⁷, Alexandre R. Zuntini¹, Wolf L. Eiserhardt^{1,8,†}, Paul J. Kersey^{1,†}, Ilia J.
9 Leitch^{1,†}, Félix Forest^{1,†}

10

11 ¹Royal Botanic Gardens, Kew, Richmond, Surrey, TW9 3AE, United Kingdom

12 ²Current address: Centre for Research in Agricultural Genomics, Campus UAB, Edifici
13 CRAG, Bellaterra Cerdanyola del Vallès, 08193 Barcelona, Spain

14 ³School of Life Sciences, University of Bedfordshire, University Square, Luton LU1 3JU,
15 United Kingdom

16 ⁴Department of Biological Sciences, Texas Tech University, Lubbock, TX 79409, USA

17 ⁵Current address: Department of Computer Science, School of Physics, Engineering and
18 Computer Science, University of Hertfordshire, Hatfield, Hertfordshire, AL10 9AB, United
19 Kingdom

20 ⁶Current address: Centre for Plant Biotechnology and Genomics (CBGP) UPM-INIA, 28223
21 Pozuelo de Alarcón (Madrid), Spain

22 ⁷Plant Science and Conservation, Chicago Botanic Garden, 1000 Lake Cook Road, Glencoe,
23 IL 60022, USA

24 ⁸Department of Biology, Aarhus University, 8000 Aarhus C, Denmark

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

25 †Joint senior authors

26 *Corresponding author: Royal Botanic Gardens, Kew, Richmond, Surrey, TW9 3AE, United
27 Kingdom, w.baker@kew.org

28

29 *Abstract.*—The tree of life is the fundamental biological roadmap for navigating the evolution
30 and properties of life on Earth, and yet remains largely unknown. Even angiosperms
31 (flowering plants) are fraught with data gaps, despite their critical role in sustaining terrestrial
32 life. Today, high-throughput sequencing promises to significantly deepen our understanding
33 of evolutionary relationships. Here, we describe a comprehensive phylogenomic platform for
34 exploring the angiosperm tree of life, comprising a set of open tools and data based on the
35 353 nuclear genes targeted by the universal Angiosperms353 sequence capture probes. The
36 primary goals of this paper are to (i) document our methods, (ii) describe our first data release
37 and (iii) present a novel open data portal, the Kew Tree of Life Explorer
38 (<https://treeoflife.kew.org>). We aim to generate novel target sequence capture data for all
39 genera of flowering plants, exploiting natural history collections such as herbarium
40 specimens, and augment it with mined public data. Our first data release, described here, is
41 the most extensive nuclear phylogenomic dataset for angiosperms to date, comprising 3,099
42 samples validated by DNA barcode and phylogenetic tests, representing all 64 orders, 404
43 families (96%) and 2,333 genera (17%). A “first pass” angiosperm tree of life was inferred
44 from the data, which totalled 824,878 sequences, 489,086,049 base pairs, and 532,260
45 alignment columns, for interactive presentation in the Kew Tree of Life Explorer. This
46 species tree was generated using methods that were rigorous, yet tractable at our scale of
47 operation. Despite limitations pertaining to taxon and gene sampling, gene recovery, models
48 of sequence evolution and paralogy, the tree strongly supports existing taxonomy, while
49 challenging numerous hypothesized relationships among orders and placing many genera for

Baker et al.

50 the first time. The validated dataset, species tree and all intermediates are openly accessible
51 via the Kew Tree of Life Explorer and will be updated as further data become available. This
52 major milestone towards a complete tree of life for all flowering plant species opens doors to
53 a highly integrated future for angiosperm phylogenomics through the systematic sequencing
54 of standardised nuclear markers. Our approach has the potential to serve as a much-needed
55 bridge between the growing movement to sequence the genomes of all life on Earth and the
56 vast phylogenomic potential of the world's natural history collections.

57 **Keywords:** angiosperms, Angiosperms353, genomics, herbariomics, museomics, nuclear
58 phylogenomics, open access, target sequence capture, tree of life.

59 INTRODUCTION

60

61 Discovering the tree of life is among the most fundamental of the grand challenges in
62 science today (Hinchliff et al. 2015). The tree of life is the biological roadmap that allows us
63 to discover, identify and classify life on Earth, to explore its properties, to understand its
64 origins and evolution, and to predict how it will respond to future environmental change. Of
65 all eukaryotic lineages, the angiosperms (flowering plants) are among the most pressing
66 priorities for tree of life research. Angiosperms sustain the terrestrial living world, including
67 humanity, as primary producers, ecosystem engineers and earth system regulators. They hold
68 potential solutions to global challenges, such as climate change, biodiversity loss, human
69 health, food security and renewable energy (Antonelli et al. 2020). In light of this, a
70 phylogenetic framework with which to navigate and interpret the species, trait and functional
71 diversity of angiosperms has never been more necessary. However, despite substantial
72 progress, the evolutionary connections among Earth's ca. 330,000 flowering plant species
73 (WCVP 2020) remain incompletely known.

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

74 The angiosperm research community were early and organised adopters of the
75 molecular phylogenetic approach, resulting in numerous benchmark tree of life publications
76 (e.g. Chase et al. 1993; Soltis et al. 2008; Soltis et al. 2011), and a community approach to
77 phylogenetic classification (APG 1998; APG II 2003; APG III 2009; APG IV 2016). Through
78 this distributed effort, a wealth of DNA sequence data is now available in public repositories,
79 covering ca. 107,000 (31%) of the ca. 350,000 species of vascular plants (RBG Kew 2016;
80 WCVP 2020), most of which are angiosperms (see also Cornwell et al. 2019). However, the
81 lack of sequence data for the remaining 69% obstructs their accurate placement in the tree of
82 life. In addition, lack of complementarity in gene sampling across public DNA sequence data
83 impedes phylogenetic synthesis (Hinchliff and Smith 2014). For example, data from either
84 one or both of *rbcL* and *matK*, the two most popular plastid genes for phylogenetics, are
85 available for only 54% of the ca. 107,000 sequenced vascular plant species (RBG Kew 2016).
86 Comprehensive phylogenetic trees of flowering plants are in high demand (Hinchliff et al.
87 2015; Eiserhardt et al. 2018), but currently can only be made “complete” using proxies, such
88 as taxonomic classification, to interpolate the unsequenced species (Smith and Brown 2018),
89 which may not accurately reflect relationships. Greater community-wide coordination of both
90 taxon and gene sampling would benefit phylogenetic data integration immensely, creating
91 numerous downstream scientific opportunities.

92 High-throughput sequencing (HTS) now promises to significantly deepen our
93 understanding of evolutionary relationships among Earth’s species, including angiosperms
94 (Li et al. 2019; Yang et al. 2020). For example, the One Thousand Plant Transcriptomes
95 (1KP) initiative has brought an unprecedented scale of data to bear on the plant tree of life
96 (Wickett et al. 2014; Gitzendanner et al. 2018; Leebens-Mack et al. 2019). Nevertheless, with
97 greatly increased data depth come trade-offs in taxon sampling; the pre-eminent HTS studies
98 cited here account for less than 0.01% of angiosperm species. Undeterred by this sampling

Baker et al.

99 gap, the Earth Biogenome Project (EBP) has launched a “moonshot for biology” by
100 proposing to sequence and characterise the genomes of all of Earth’s eukaryotic species over
101 a 10-year period (Lewin et al. 2018). Projects such as the 10,000 Plant Genomes Project
102 (Cheng et al. 2018) and the Darwin Tree of Life Project (<https://www.darwintreeoflife.org/>)
103 aim to contribute to this goal by producing numerous chromosome-level genome assemblies
104 across major lineages and regional biotas. However, taxon sampling remains a significant
105 issue, due to the challenges of obtaining the high molecular weight DNA required by these
106 projects (for long-read HTS) from samples that are both authentically identified and
107 compliant with the spirit and letter of the Nagoya Protocol (Secretariat of the Convention on
108 Biological Diversity 2011). Despite its immense potential, the “whole genome” approach to
109 discovering the tree of life remains a future goal that will not be achieved on a large
110 taxonomic scale in the short term. Methodological compromises are required to accelerate
111 progress.

112 The world’s natural history collections are a goldmine for genomic research (Buerki
113 and Baker 2016), containing tissues of almost all species of life on Earth known to science.
114 However, the condition of these tissues and the DNA therein varies widely, depending on age
115 and preservation techniques, among other factors. In the case of plants, herbarium specimens
116 generally yield degraded DNA, which, though not useful for long-read HTS, is now being
117 intensively exploited for short-read HTS (Bakker et al. 2016; Brewer et al. 2019; Forrest et al.
118 2019; Alsos et al. 2020). In this context, target sequence capture is growing in popularity as
119 the HTS method most widely applied to herbarium DNA (Dodsworth et al. 2019). This
120 approach (also known as target enrichment, target capture, sequence capture, anchored hybrid
121 enrichment) and its variations (e.g. Hyb-Seq, which combines target sequence capture with
122 genome skimming) use RNA or DNA probes to enrich sequencing libraries for specifically
123 targeted loci (Faircloth et al. 2012; Lemmon et al. 2012; Weitemier et al. 2014). It is proving

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

124 to be an increasingly cost-effective means of isolating hundreds of loci for phylogenetic
125 analysis from even centuries-old specimens (Brewer et al. 2019), bringing comprehensive
126 taxon sampling from herbarium collections within the reach of any phylogenomic researcher
127 (Hale et al. 2020).

128 Numerous target sequence probe sets have been developed for specific angiosperm
129 groups (e.g. Annonaceae [Couvreur et al. 2019], Asteraceae [Mandel et al. 2014], *Dioscorea*
130 [Soto Gomez et al. 2019], *Euphorbia* [Villaverde et al. 2018]). The design of these probe sets
131 is informed by available genomic resources, as well as criteria specific to the group of interest
132 and research questions. As a result, locus overlap between probe sets tends to be minimal.
133 Unlike the Sanger sequencing era, in which researchers converged on tractable genes such as
134 *rbcL* and *matK*, the lack of complementarity between probe sets curtails prospects for data
135 integration across broad taxonomic scales. In addition, development of custom probe sets is
136 expensive, requiring considerable genomic resources and bioinformatic expertise. A publicly
137 available, universal probe set for angiosperms targeting a standard set of loci would resolve
138 these issues (Buddenhagen et al. 2016; Chau et al. 2018). In response to this, we designed the
139 Angiosperms353 probe set (Johnson et al. 2019), drawing on 1KP transcriptome data from
140 ca. 650 species across the angiosperms (Leebens-Mack et al. 2019). The probe set targets 353
141 genes from 410 low-copy, protein-coding nuclear orthologs previously selected for
142 phylogenetic analysis across green plants (Leebens-Mack et al. 2019), enriching up to ca. 260
143 kbp from any flowering plant. Angiosperms353 probes are an open data resource that can be
144 used without the expense of design or access to prior genomic data (Baker et al. 2021) and
145 have already been successfully applied across different taxonomic scales (e.g. Larridon et al.
146 2019; Murphy et al. 2020; Pérez-Escobar et al. 2020; Shee et al. 2020), including at the
147 population level (Van Andel et al. 2019; Slimp et al. 2020; Beck et al. 2021).

Baker et al.

148 Here, we describe a large-scale effort to establish a new phylogenomic platform for
149 exploring the angiosperm tree of life, comprising a set of open tools (Angiosperms353
150 probes, laboratory protocols, analysis pipeline, data portal) and data (sequence data,
151 assembled genes, alignments, gene trees, species tree). This platform, which directly
152 addresses the challenges outlined above, is an outcome of the Plant and Fungal Trees of Life
153 project (PAFTOL; www.paftol.org) at the Royal Botanic Gardens, Kew (RBG Kew 2015).
154 As a step towards the ultimate goal of a complete species-level tree, we aim to gather DNA
155 sequence data for the Angiosperms353 genes from one species of all 13,862 angiosperm
156 genera (WCVP 2020). This unprecedented dataset of standard loci draws extensively on
157 herbarium collections for comprehensive sampling, especially of genera that have not been
158 sequenced before (Brewer et al. 2019). Extensive new data have been generated, analysed
159 and released into the public domain, along with corresponding phylogenetic inferences. By
160 providing our data in open and accessible ways, including an interactive tree of life, we aim
161 to foster a transparent and collaborative environment for future data re-use and synthesis.
162 This paper serves as the baseline reference for our platform, (i) documenting our methods, (ii)
163 describing our first data release, comprising 17% of angiosperm genera, including initial
164 insights on phylogenetic performance, and (iii) presenting a novel data portal, the Kew Tree
165 of Life Explorer, through which our data and corresponding tree of life can be interrogated
166 and downloaded. We conclude with reflections on the prospects for our approach, future
167 development requirements and the role of open data for enhancing cross-community
168 collaboration towards a complete tree of life.

169 **MATERIALS AND METHODS**

170

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

171 This section describes the workflow (Fig. 1) used by the PAFTOL project to generate
172 our first data release (i.e. Data Release 1.0), which is publicly accessible through our open
173 data portal, the Kew Tree of Life Explorer (<https://treeoflife.kew.org>), described below. The
174 workflow consists of three main stages: (i) sample processing, encompassing sample
175 selection and laboratory protocols for target sequence capture data generation (Fig. 2), (ii)
176 data analysis, including target gene assembly, data mining, data validation and phylogenetic
177 inference (Figs. 2, 3), and (iii) data publication via the data portal (Fig. 4). The data
178 accessible via the portal comprise raw data (unprocessed sequence reads) and results from
179 “first pass” analyses (gene assemblies, alignments, gene trees, species tree). Though not
180 exhaustive, these first explorations of the data apply methods that are both rigorous and
181 tractable at our scale of operation.

182 Details of the first data release are also given in the data release notes in the portal via
183 our secure FTP (<http://sftp.kew.org/pub/treeoflife/>) and are also archived at the Royal Botanic
184 Gardens, Kew (RBGK) Research Repository (<https://doi.org/10.34885/paftol>). A new release
185 note will be published in the same locations with each future data release and will detail any
186 changes in methods used relative to the first release described here.

187 **Sampling**

188 We aimed to generate novel data from across the angiosperms, using a stratified
189 sampling approach of one species per genus. Our sampling was standardised to the complete
190 list of angiosperms within the World Checklist of Vascular Plants (WCVP 2020), which
191 currently recognises 13,862 accepted genera in 418 families, aligned to the 64 orders of the
192 APG IV classification (APG IV 2016). We prioritised genera that were not represented by
193 published transcriptomic or genomic data in public sequence repositories (e.g. GenBank), and
194 avoided genera that had already been sampled in large genomic initiatives such as the 1KP

Baker et al.

195 project (Leebens-Mack et al. 2019). The selection of species within genera was made
196 pragmatically, although we prioritised the species of the generic type where possible.

197 Plant material was obtained from a variety of sources (Fig. 2), primarily from the
198 collections of RBGK (herbarium, DNA bank, silica gel-dried tissue collection, living
199 collection and the Millennium Seed Bank, [https://www.kew.org/science/collections-and-](https://www.kew.org/science/collections-and-resources/collections)
200 [resources/collections](https://www.kew.org/science/collections-and-resources/collections)). Additional material (tissue samples, extracted DNA) was generously
201 provided by individuals in our collaborative networks (see Acknowledgements). To be
202 selected, the material must have been (i) legally sourced and made available for use in
203 phylogenomic studies, (ii) identified to species level, preferably by an expert in the group,
204 and (iii) ideally collected in the wild. As far as was practically achievable, we ensured that
205 the identity of each sample was substantiated by a voucher specimen deposited in a publicly
206 accessible herbarium.

207 All metadata were captured using a relational database that allowed us to track
208 processing of samples from the selection of material, through the library preparation pipeline
209 to the completion of sequencing. Data were recorded in four main tables (Specimen, Sample,
210 Library, Sequencing). The database architecture allowed us to record multiple sequence
211 datasets (fastq files) from one or several libraries, and one or several DNA extracts from a
212 single specimen. Relevant voucher specimen information was also captured in the database
213 (e.g. collector(s), collector number, herbarium acronym (following Index Herbariorum
214 <http://sweetgum.nybg.org/science/ih/>), country of origin, date of collection, specimen
215 barcodes). Voucher data are available via our data portal (see below). Images of specimens
216 sampled from the RBGK Herbarium are in the process of being captured in RBGK's online
217 herbarium catalogue (<http://apps.kew.org/herbcat/>) and, where available, are linked to the
218 appropriate records in the Kew Tree of Life Explorer.

219

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

220 **DNA extraction**

221 DNA was extracted from 40 mg of herbarium material, 20 mg of silica gel-dried
222 material (Chase and Hills 1991), or 100 mg of fresh material using a modified CTAB
223 extraction method (Doyle and Doyle 1987; Fig. 2). Plant tissue was pulverized using a Mixer
224 Mill MM400 (Retsch GmbH, Germany). DNA extractions were purified by a magnetic bead
225 clean-up using Agencourt AMPure XP beads (Beckman Coulter, Indianapolis, IN, USA),
226 according to the manufacturer's protocols. Samples obtained from the RBGK DNA bank
227 (<http://dnabank.science.kew.org/homepage.html>) had been extracted using a modified CTAB
228 method (Doyle and Doyle 1987) followed by caesium chloride/ethidium bromide density
229 gradient cleaning and dialysis. DNA samples provided by external collaborators had been
230 extracted using a wide variety of extraction methods from living, silica gel-dried and
231 herbarium material.

232 All DNA samples were quality checked for concentration and degree of
233 fragmentation. DNA concentration was measured using a Quantus (Promega, Madison, WI,
234 USA) or Qubit (Thermo Fisher Scientific, Inchinnan, UK) fluorometer. DNA fragment size
235 range was routinely assessed on a 1% agarose gel using ethidium bromide and visualized
236 with a UVP Gel Studio (AnalytikJena, Jena, Germany). For samples with a low DNA
237 concentration (i.e. not visible on a gel), fragment sizes were assessed on a 4200 TapeStation
238 using Genomic DNA ScreenTape (Agilent Technologies, Cheadle, UK).

239 **Library preparation**

240 Genomic DNA samples were diluted to 4 ng/ μ l with 10 mM Tris (pH 8.0). Those with
241 an average fragment size greater than 350 bp were sonicated to an average fragment size ca.
242 400 bp, using a Covaris M220 Focused-ultrasonicator (Covaris, Woburn, MA, USA) by
243 adding 50 μ l of diluted genomic DNA to a 130 μ l Covaris microAFA tube. The sonication
244 time was adjusted for each sample based on its average DNA fragment size (15 to 100 secs,

Baker et al.

245 following the manufacturer's protocols). Additional parameters used were peak incident
246 power to 50W, duty factor to 10% and 200 cycles per burst.

247 Libraries were prepared using the NEBNext Ultra II DNA Library Prep Kit (New
248 England Biolabs, Ipswich, MA, USA; Fig. 2). Size selection was not employed for samples
249 with highly degraded DNA. In the early stages of the project, libraries were prepared
250 following the manufacturer's protocols exactly, but the majority were prepared using half of
251 the recommended volumes throughout to reduce costs. All DNA fragments were indexed
252 using NEBNext Multiplex Oligos for Illumina (Dual Index Primer sets 1 and 2, New England
253 Biolabs, Ipswich, MA, USA).

254 The distribution of fragment sizes in each library was assessed with a 4200
255 TapeStation using standard D1000 tapes. Library concentration was measured using a
256 Quantus fluorometer. If the library concentration was less than 10 nM, up to eight additional
257 PCR cycles were performed, following the NEBNext Ultra II Library Prep Kit protocol with
258 IS5_reamp.P5 and IS6_reamp.P7 primers (Meyer and Kircher 2010). Library quality
259 assessment was then repeated.

260 **Pooling and hybridisation**

261 Prior to hybridisation (Fig. 2), all libraries were normalised to 10 nM, using 10 mM
262 Tris (pH 8.0) and then combined into pools of 20 to 24 libraries, each containing 10 μ l (0.1
263 pmol) of each normalized library (i.e. a total of ca. 600-700 ng DNA in each pool, assuming
264 an average fragment size of ca. 450 bp). To ensure even sequencing across all samples in a
265 pool, species for pooling were selected to minimize the range of DNA fragment sizes and
266 ensure a narrow taxonomic breadth. The latter criterion was needed because samples that are
267 more closely related to the taxa used to construct the probe set tend to preferentially
268 hybridise. This can lead to an over-representation of their sequences in the DNA data if
269 appropriate care is not taken when selecting species for the sequencing pool. In rare cases,

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

270 such as smaller pools (ca. 10 libraries) of short fragment (i.e. <300 bp) libraries, it was
271 necessary to recalculate the standard volume of normalized libraries to be added to ensure
272 that the final pool contained ca. 500 ng of DNA.

273 The pooled libraries were dried in a SpinVac (Eppendorf, Dusseldorf, Germany),
274 resuspended in 8 µl of 10 mM Tris (pH 8.0) and enriched by hybridising with the
275 Angiosperms353 probe kit (Johnson et al. 2019; Arbor Biosciences myBaits Target Sequence
276 Capture Kit, ‘Angiosperms 353 v1’, Catalogue #308196) following the manufacturer’s
277 protocol, version 4.0. Hybridisation was typically performed at 65°C for 24 h, with reactions
278 topped with 30 µl of red Chill-out Liquid Wax (Bio-Rad, Hercules, CA, USA) to prevent
279 evaporation. However, for short libraries (i.e. <350 bp) the temperature was reduced to 60°C,
280 following the recommendations of Arbor Biosciences.

281 The target-enriched pools were amplified using the KAPA HiFi 2X HotStart
282 ReadyMix PCR Kit (Roche, Basel, Switzerland) or NEBNext Q5 HotStart HiFi PCR Master
283 Mix (New England BioLabs, Ipswich, MA, USA) for eight to 14 cycles. Amplified pools
284 were then purified using Agencourt AMPure XP Beads (at 0.9X the sample volume) and
285 eluted in 15 µl of 10 mM Tris (pH 8.0).

286 Products were quantified with a Quantus fluorometer and re-amplified if the
287 concentration was below 6 nM, with three to six PCR cycles (see above). Final products were
288 assessed using the TapeStation to determine the distribution of fragment sizes. The target-
289 enriched pools were normalized to 6 nM (using 10 nM Tris, pH 8.0) and multiplexed for
290 sequencing, with the number of target-enriched pools combined in each sequencing pool
291 varying from two to 20 (comprising a total of 48-384 samples) depending on the sequencing
292 platform and service provider requirements.

293

Baker et al.

294 **DNA sequencing**

295 Initially, DNA sequencing was performed on an Illumina MiSeq at RBGK with
296 version 3 chemistry (Illumina, San Diego, CA, USA) and ran for 600 cycles to generate 2 ×
297 300 bp paired-end reads. Subsequently, DNA sequencing was outsourced (Macrogen, Seoul,
298 South Korea, or Genewiz, Takeley, UK) and performed on an Illumina HiSeq producing 2 ×
299 150 bp paired-end reads. Raw reads were deposited in the European Nucleotide Archive
300 under an umbrella project (accession number PRJEB35285) and can be accessed from the
301 individual sample records in the Kew Tree of Life Explorer.

302

303 **Sequence assembly**

304 Coding sequences were recovered from target-enriched sequence data using our
305 pipeline recoverSeqs (accessible from our GitHub repository
306 <https://github.com/RBGKew/KewTreeOfLife>, pypaftol ‘paftools’ submodule) to retrieve
307 sequences orthologous to the Angiosperms353 target gene set (Johnson et al. 2019;
308 <https://github.com/mossmatters/Angiosperms353>). This target set contained multiple
309 reference sequences per gene, thereby covering a large phylogenetic breadth to facilitate read
310 recovery across angiosperms.

311 The process comprised four main stages (Fig. 2), applied to each sample: (i) sequence
312 reads were trimmed using Trimmomatic (Bolger et al. 2014) with the following parameters:
313 ILLUMINACLIP: <AdapterFastaFile>: 2:30:10:2:true, LEADING: 10, TRAILING: 10,
314 SLIDINGWINDOW: 4:20, MINLEN: 40, with the adaptor fasta file formatted for
315 palindrome trimming, (ii) trimmed read pairs were mapped to the Angiosperms353 target
316 genes with TBLASTN. A representative reference sequence for each gene was then selected
317 by identifying the sequence with the largest number of mapped reads. (iii) This representative
318 gene was used as the reference for assembling the gene-specific reads using an overlap-based

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

319 assembly algorithm (--assembler overlapSerial option). OverlapSerial was developed
320 specifically for this project (see our GitHub repository) with the aim of improving gene
321 recovery, in terms of gene length and number, relative to the widely used HybPiper (Johnson
322 et al. 2016) and was used as follows. First, the reads were aligned to and ordered along the
323 reference sequence based on a minimum alignment size of 50 bases (--windowSizeReference
324 option) with a minimum sequence identity of 70% (--relIdentityThresholdReference option).
325 Consecutive reads ordered along the reference sequence were aligned in a pair-wise manner
326 to find read overlaps. If an overlap of at least 30 bases (--windowSizeReadOverlap option)
327 and 90% sequence identity (--relIdentityThresholdReadOverlap option) was found, the
328 aligned reads were used to construct a consensus contig with ambiguous bases represented by
329 'N'. This last parameter resulted in one or more sets of aligned reads with $\geq 90\%$ sequence
330 identity, each set being merged into a single contig. In the final stage, the exonerate
331 protein2genome program was used to identify the exon-intron structure within each contig.
332 One or more contigs were chosen that best represented the structure of the exon(s) in the
333 reference gene chosen in step (ii). If the exons existed in multiple contigs, those contigs were
334 joined together to form the recovered gene coding sequence.

335 Target gene recovery success was assessed for each sample by calculating the number
336 of genes recovered and the sum of the recovered gene lengths. Samples were removed from
337 downstream analyses if the sum of the recovered gene lengths fell below 20% of the median
338 value across all samples.

339

340 Public data mining

341 In addition to newly generated target sequence capture data, the Angiosperms353
342 genes were mined from publicly available genomic data (Fig. 2). For Data Release 1.0, we
343 focused on mining data from the 1KP Initiative (Carpenter et al. 2019; Leebens-Mack et al.

Baker et al.

344 2019) and published genomes with gene annotations (<https://plants.ensembl.org/>), although
345 other data sources (e.g. the Sequence Read Archive) will be data-mined for future releases.
346 The genes were retrieved from assembled transcript sequences (1KP) or coding sequences
347 (CDS; genomes) using paftools retrievetargets from our pipeline, which uses BLASTN to
348 identify and extract the genomic or transcriptomic sequences corresponding to the 353 genes.
349 BLASTN relies on sequence identity (>70%) and the transcript or CDS with the highest
350 identity is considered to be the ortholog of a given target. Because initial recovery of genes
351 from 1KP transcripts using the standard Angiosperms353 target gene set (Johnson et al.
352 2019) was unsatisfactory, we used an expanded Angiosperms353 target set to improve
353 matching and retrieval of genes. The expanded dataset is a reduced version of the 1KP
354 alignments (Leebens-Mack et al. 2019) produced by Johnson et al. (2019) for the design of
355 the Angiosperms353 probe set from which non-angiosperm sequences had been removed and
356 gap-only sites trimmed. The expanded target set is available from
357 <https://datadryad.org/stash/dataset/doi:10.5061/dryad.s3h9r6j> and a reformatted version from
358 our GitHub. As with the novel target sequence capture assemblies, data were removed from
359 downstream analyses if the sum of the gene lengths fell below 20% of the median value
360 across all samples.

361

362 **Family identification validation**

363 To verify the family identification of our processed samples, we implemented two
364 validation steps, which were run in parallel (Fig. 3). The two steps consisted of (i) DNA
365 barcode validation, which utilised nuclear ribosomal and plastid barcodes for DNA-based
366 identification, and (ii) phylogenetic validation, which checked the placement of each sample
367 in a preliminary tree relative to its expected position based on its initial family assignment.
368 Identification checks below the family level were not conducted due to the incompleteness of

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

369 adequate reference resources for DNA barcode validation and sparseness of sampling for
370 phylogenetic validation at the genus or species level.

371 For barcode validation of target sequence capture data (Fig. 3), plastomes and
372 ribosomal DNA were recovered from raw reads using GetOrganelle (Jin et al. 2020) and
373 subsequently queried against databases of reference plant barcodes using BLASTN
374 (Camacho et al. 2009). For 1KP samples, transcriptome assemblies were directly used as
375 queries in BLASTN. Note that we considered the family identity of annotated genomes to be
376 correct and hence a barcode validation was unnecessary. Six individual barcode reference
377 databases were built from the NCBI nucleotide and BOLD databases
378 (<https://www.ncbi.nlm.nih.gov/nucleotide/>; <https://www.boldsystems.org/>, accessed on
379 29/10/2020), one for the whole plastome, and the remaining five for specific loci (nuclear
380 ribosomal 18S, as well as plastid *rbcL*, *matK*, *trnL*, and *trnH-psbA*). As for samples, the
381 taxonomy of reference sequences was standardized to WCVP (WCVP 2020). BLAST results
382 were further filtered with a minimum identity >95% and a minimum coverage of reference
383 locus $\geq 90\%$ (except for whole plastomes, for which only a filtering based on minimum length
384 was applied).

385 Tests could only be completed if a sample's given family was present in the barcode
386 databases and if at least one BLAST match remained after filtering. Thus, zero to six barcode
387 tests were conducted per sample. A sample passed an individual test if the first ranked
388 BLAST match (ranked by percentage of identity) confirmed its original family identification
389 and failed otherwise. The final result of the barcode validation following the six individual
390 barcode tests were determined as follows: (i) Confirmed, if one or more barcode tests
391 matched the family identification of a sample; (ii) Rejected, if more than half of the barcode
392 tests gave the same incorrect family identification (requires at least two barcode tests); (iii)
393 Inconclusive (otherwise). Further details of the barcode validation methods can be found in

Baker et al.

394 Supplementary Material available on Dryad. The scripts and lists of NCBI and BOLD
395 accessions used in barcode databases are available on our GitHub repository.

396 To conduct phylogenetic validation (Fig. 3), a preliminary phylogenetic tree was built
397 using the complete, unvalidated dataset, following the phylogenetic methods described
398 below. We then assessed which nodes best represented each order and family in the tree. For
399 every node in the tree, two metrics were calculated for all families and orders: (i) the
400 proportion of samples belonging to a given order/family that are descendants of the node, and
401 (ii) the proportion of samples descending from the node that belong to the order/family. The
402 two metrics were then multiplied to produce an overall taxon concordance score. For each
403 family and order, the highest scoring node was subsequently considered to best represent the
404 taxon in the tree (allowing the identification of outlying samples). A node with a score of 1
405 for a given order/family is the crown node (most recent common ancestral node) of that
406 taxon, which is monophyletic in the tree. See Supplementary Figure S1 for an illustration.
407 The family identification of each sample was determined as (i) Confirmed: if identified as
408 belonging to a family whose best scoring node had a taxon concordance score >0.5 and found
409 as a descendant of this node in the tree, (ii) Rejected: if identified as belonging to a family
410 whose best scoring node had a taxon concordance score >0.5 but not found as a descendant of
411 this node, or (iii) Inconclusive: if identified as belonging to a family whose best scoring node
412 had a taxon concordance score ≤ 0.5 . Note that for families represented in the tree by a single
413 sample, the validation was performed with respect to their orders. If the order was
414 represented by a single sample, the validation result was coded as inconclusive.

415 The outputs of the phylogenetic and DNA barcode validation were combined to
416 identify samples for automatic inclusion and exclusion from the final dataset, and samples for
417 which a decision on inclusion/exclusion was subject to expert review (Fig. 3). Exclusions

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

418 after expert review were made based on implausible tree placement (e.g. wrong higher clade)
419 or sample misidentification (e.g. match to another family in the barcode validation).

420 All assembled Angiosperms353 gene data from all samples validated for inclusion
421 form the basis of Data Release 1.0. These were made publicly available via the Kew Tree of
422 Life Explorer.

423

424 Phylogeny estimation

425 We inferred a phylogenetic tree from all validated data (Data Release 1.0) for
426 presentation in an interactive format in the Kew Tree of Life Explorer. This species tree was
427 estimated from gene trees using the multi-species coalescent summary method implemented
428 in ASTRAL-III (Zhang et al. 2018). In addition to the angiosperm samples, ten samples
429 representing seven gymnosperm families from the 1KP initiative were mined for
430 Angiosperms353 orthologs (using retrievetargets, as described above) and included in all
431 analyses as outgroup taxa. Our phylogenomic pipeline, available from our GitHub repository,
432 is summarised below.

433 For each gene, DNA sequences were aligned with UPP 4.3.12 (Nguyen et al. 2015).
434 At the start of the alignment process a set of 1,000 sequences were selected for an initial
435 backbone tree. Option -M was set to '-1' so that sequences could be selected within 25% of
436 the median full-length sequence. Filtering and trimming of the alignment were performed
437 with AMAS (Borowiec 2016) as follows. Sequences with insufficient coverage (<60%)
438 across well occupied columns of each gene alignment were removed. Well occupied columns
439 were defined as those with more than 70% of positions occupied. Then, alignment columns
440 with <0.3% occupancy were removed to remove very rare or unique insertions. Finally,
441 sequences with a total length of less than 80 bases were removed, and genes with <30
442 overlapping bases (at the 70% threshold mentioned above) were excluded.

Baker et al.

443 Gene trees were estimated with IQ-TREE 2.0.5 (Minh et al. 2020) inferring branch
444 support using the ultrafast bootstrap method (option -B; Hoang et al. 2017) with the
445 maximum number of iterations set to 1,000 (option -nm) and using a single model of
446 evolution (option -m GTR+F+R). The use of a single model without testing many models of
447 evolution was a pragmatic choice, following Abadi et al. (2019). TreeShrink 1.3.4 (Mai and
448 Mirarab 2018) was used to remove abnormally long branches from gene trees using default
449 settings, except option -b, which was set to 20. The alignment and gene tree estimation steps
450 were then repeated on the samples retained by TreeShrink. Before reconstructing the species
451 tree using ASTRAL-III, nodes in the gene trees with bootstrap support values less than 30%
452 were collapsed using nw_ed from Newick Utilities 1.6.0 (Junier and Zdobnov 2010). This
453 value was deduced from interpreting Figure 1 in Hoang et al. (2017), adjusting the standard
454 bootstrap threshold of 10% (recommended for ASTRAL-III), to 30 % for the ultrafast
455 bootstrap.

456 All gene alignments, gene trees and the ASTRAL-III species tree are available for
457 download from secure FTP and the Kew Tree of Life Explorer. In addition, the species tree is
458 available to browse through an interactive tree viewer implemented within the Kew Tree of
459 Life Explorer (see also Supplementary Fig. S2).

460

461 **Data portal implementation**

462 To disseminate results, a data portal (the Kew Tree of Life Explorer;
463 <https://treeoflife.kew.org>) was designed and implemented (Fig. 4) with a layered architecture
464 that comprised: (i) a MariaDB running on a Galera multi-master cluster as a database
465 management system; (ii) an API written in Python using the Flask framework and the
466 SQLAlchemy library; (iii) a front-end written using the Vue.js framework and Nuxt.js for the
467 tabular data (used to provide access to gene and specimen data) and content pages; (iv) a tree

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

468 visualisation module developed from the open source application PhyD3 (Kreft et al. 2017)
469 using D3.js (Bostock 2012) for data visualisation; and (v) deployment on a Linux (CentOS 7)
470 server using Nginx as web server and load balancer.

471 The data, with appropriate metadata and documentation, are available for public
472 download over secure FTP (<http://sftp.kew.org/pub/treeoflife/>) and the Kew Tree of Life
473 Explorer under a Creative Commons Attribution 4.0 International (CC BY 4.0) license. When
474 superseded by new releases, archived earlier releases will remain accessible via secure FTP.

475 RESULTS

476 Initial dataset

477 The initial dataset prior to processing and analysis comprised data from 3,272
478 angiosperm samples, representing 413 families of angiosperms (99%) and 2,428 genera
479 (18%; Table 1). We generated novel target sequence capture data for 2,522 of these samples,
480 which included 104 angiosperm genera that have never been sequenced before. Data for the
481 remainder were mined from public sources (689 1KP transcriptomes, 61 annotated genomes).
482 The majority of target sequence capture data were generated from the RBGK collections as
483 follows: DNA Bank (43%), herbarium (28%), silica gel-dried tissue collection (8%), living
484 collection (2%), and Millennium Seed Bank (0.3%). The remaining 19% of samples included
485 in this study were provided by various collaborators of the PAFTOL project, either as DNA
486 samples or as dried tissue (see Acknowledgements).

487 Sequence recovery from all 2,522 target sequence capture samples (prior to any
488 quality controls) is visualised in Figure 5. Eighty-four target sequence capture samples and
489 eleven 1KP transcriptomes were removed from downstream analyses because the sum of
490 gene lengths did not meet the quality threshold of 20% of the median value across all
491 samples.

Baker et al.

492 Family identification validation

493 The remaining 3,177 samples (Table 1) were processed through our sample family
494 identification validation pipeline (Fig. 3, Supplementary Tables S1 and S2). Of these, 3,064
495 (97%) were automatically cleared for inclusion and 67 were automatically excluded
496 (Supplementary Table S1). The remaining 46 samples were held for expert review, after
497 which 35 were cleared for inclusion and 11 were excluded due to implausible tree
498 placements. The majority of excluded samples (64 out of 78) were from the novel target
499 sequence capture data, although 14 were 1KP transcriptomes, highlighting the risk of sample
500 misidentification in even the most highly curated datasets. Further details regarding the
501 results obtained during the family identification validation by DNA barcoding can be found
502 in Supplementary Material available on Dryad.

503 The final validated dataset for Data Release 1.0 consisted of 3,099 angiosperm
504 samples (Table 1), only 5% fewer than were present in the initial dataset. These samples
505 represent all 64 orders, 404 families (96%; 212 represented by >1 sample), 2,333 genera
506 (17%) and 2,956 species (0.01%).

507 Data Release 1.0: sequence quality and gene recovery

508 Nine statistics were used to assess the sequence quality across the 3,099 samples of
509 Data Release 1.0 (Table 2). For the 2,374 target sequence capture samples, the mean
510 percentage of on-target reads was 8%, the mean read depth per sample across all recovered
511 genes was 90x with a median value of 38x and the mean percentage length of recovered
512 genes per sample was 62%. The number of genes and the sum length of gene sequence
513 recovered per sample were tightly associated as expected, varying continuously across the
514 dataset up to the full set of Angiosperms353 genes and a total gene length of 256.9 kbp, close
515 to the maximum expected length of 260 kbp for recovering genes with this target gene set
516 (Fig. 5). The total length of sequence recovered from target sequence capture data was shorter

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

517 than for samples mined for Angiosperms353 genes from 1KP transcriptomes or annotated
518 genomes data (Table 2). The reason for the shorter length of the recovered genes is that some
519 exons were omitted during the process of refining 1KP alignments to select gene instances for
520 the design of the Angiosperms353 probes (Johnson et al. 2019). These missing exons were
521 however present in the expanded Angiosperms353 target set and were therefore retrieved
522 during data mining from 1KP transcriptomes and annotated genomes. The variation in
523 performance of target enrichment across different samples, illustrated by the measures of
524 variability shown in Table 2, likely reflects the variation in structure and metabolite
525 composition of the starting tissue, which is known to impede DNA extraction from various
526 species and its downstream manipulation. This variation is one of the challenges in dealing
527 with samples from a broad taxonomic range such as across the evolutionary diversity of
528 angiosperms. Variation in gene recovery across orders is visualised in Supplementary Figure
529 S3.

530 Phylogenetic results

531 The final phylogenetic tree as inferred from Data Release 1.0 is publicly available in
532 interactive form via the Kew Tree of Life Explorer. In the current release, the tree is
533 annotated with local posterior probabilities (LPP, as given by ASTRAL-III) as indicators of
534 branch support. Other measures of support (e.g. quartet scores) can be found within tree files
535 accessible via the RBGK secure FTP. For completeness, the tree is also available in various
536 formats, including Newick (Supplementary Fig. S2).

537 As a result of filtering and trimming steps during alignment, six genes in Data Release
538 1.0 were excluded from downstream phylogenetic analysis due to insufficient overlap
539 between sequences. All statistics provided below refer to the remaining dataset. Thus, the
540 species tree is based on 347 gene alignments totalling 824,878 sequences, 489,086,049 base
541 pairs and 532,260 alignment columns. Of these, 509,987 columns (96%) are variable and

Baker et al.

542 475,181 columns (89%) are parsimony informative. The proportion of gaps across all
543 alignments is 61.6% and the median number of genes per sample is 284 (mean: 265.3,
544 standard deviation (SD): 64.3, min: 22, max: 347; Supplementary Table S3). The median
545 number of samples per gene alignment is 2,421 (mean: 2,377.2, SD: 359) and median
546 alignment length is 1,259 (mean: 1,533.9, SD: 985.7; Table 3). The resulting gene trees are
547 highly resolved, with a median support across all nodes (ultrafast bootstrap) of 98% (mean:
548 87.8%, standard deviation (SD): 18.560) across all nodes in all gene trees (Supplementary
549 Fig. S4). Only 1.3% of all nodes in all gene trees are very poorly supported (ultrafast
550 bootstrap <30%; Supplementary Fig. S4) and thus collapsed prior to species tree inference.
551 Further statistics for individual gene alignments and gene trees are reported in Table 3 and
552 Supplementary Table S3.

553 The species tree accommodates 82% of the quartet relationships in the gene trees
554 (ASTRAL normalized quartet score of 0.82). The majority (76.8%) of nodes in the species
555 tree were well-supported (LPP \geq 95%, cf. Sayyari and Mirarab 2016), and only seven nodes
556 were informed by too few gene trees (i.e. <20) to evaluate support. Comparing node support
557 in the species tree at different taxonomic levels (Supplementary Fig. S5), median quartet
558 support is progressively higher towards shallower taxonomic levels (Supplementary Fig.
559 S5c), while the effective number of gene trees informing nodes shows the opposite trend
560 (Supplementary Fig. S5e). Local posterior probabilities show a tendency to be lower (1st
561 quartile) at the deepest taxonomic level (Supplementary Fig. S5a). Major groups (i.e.
562 monocots, asterids and rosids) show similar distributions of both local posterior probabilities
563 (Supplementary Fig. S5b) and quartet support values (Supplementary Fig. S5d), despite the
564 fact that the effective number of gene trees supporting nodes is more variable in monocots
565 (Supplementary Fig. S5f), which is the result of the lower recovery rates for some orders in
566 this group such as Alismatales, Commelinales and Liliales (Supplementary Fig. S3).

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

567 Discounting taxa represented by a single sample (193 families, one order), 96% of
568 testable families and 83% of testable orders were resolved as monophyletic in the species
569 tree. Most of the samples of non-monophyletic families and orders could be assigned to a
570 clade that represents the family or order well, despite lacking some samples and/or containing
571 some outlier samples from other taxa (“concordant taxa” where taxon concordance score
572 >0.5 , see Materials and Methods for details). Only five families (Francoaceae,
573 Hernandiaceae, Phyllanthaceae, Pontederiaceae and Schlegeliaceae, represented by 11
574 samples) and two orders (Bruniales and Icaciniales, represented by six samples) were so
575 dispersed that this was not possible (“discordant taxa” where taxon concordance score ≤ 0.5).
576 At the family level, 2,893 samples were resolved in the expected family, two samples were
577 resolved in an unexpected position, and 204 samples were not testable because they belonged
578 to a discordant family or a family represented by a single sample. At the order level, 3,060
579 samples were resolved in the expected order, 32 samples were resolved in an unexpected
580 position, and seven samples were not testable (see Supplementary Tables S4-S6 for lists of
581 specimens from singly represented taxa, poorly resolved taxa, and outliers to well-resolved
582 taxa, respectively). Placements of all but five genera and seven families were consistent with
583 the WCVP/APG IV taxonomic hierarchy of genera, families and orders. Concordance with
584 existing taxonomy was lower at the genus level, with only 74% of testable genera resolving
585 as monophyletic and 47 genera (represented by 130 samples) being discordant; these numbers
586 partly reflect the deliberate inclusion of multiple samples from genera suspected *a priori* to
587 be potentially non-monophyletic.

588 In addition to resolving most genera, families and orders as monophyletic, our tree
589 supports more than half (58%) of the relationships among orders presented by the
590 Angiosperm Phylogeny Group (APG IV 2016; Supplementary Fig. S6). Congruence with
591 APG IV varies among major clades, being notably high in magnoliids (100% of APG IV

Baker et al.

592 relationships supported) and monocots (80%), while being substantially lower in eudicots
593 (47%), especially in rosids (33%). Nodes in our tree that are congruent with APG IV ordinal
594 relationships are slightly better supported on average (mean LPP 0.98, median 1) than nodes
595 that are incongruent with APG IV (mean LPP 0.75, median 0.94).

596 **Tree of Life Explorer**

597 The Kew Tree of Life Explorer (<https://treeoflife.kew.org>) provides open access to
598 taxon, specimen, sequence, alignment and tree data, with associated metadata for the current
599 data release in accordance with the Toronto guidelines on pre-publication data sharing
600 (Toronto International Data Release Workshop Authors 2009). Users can browse by species,
601 gene or interactive phylogenetic tree. The species interface permits searches by order, family,
602 genus or species, and provides voucher specimen metadata (including links to online
603 specimen images, where available), simple sequence metrics, access to assembled genes and
604 raw data. The gene interface documents all Angiosperms 353 genes and associated metrics,
605 links to gene identities in UniProt (<https://www.uniprot.org/>) and provides access to
606 assembled genes across taxa. The tree of life interface enables browsing and taxon searching
607 of the species tree inferred from the current release dataset, as well as tree downloads (as
608 PNG or Newick) and zooming into user-defined subtrees. All processed data (assembled
609 genes, alignments, gene trees, species trees) and archived releases are available from
610 RBGK's secure FTP site (<http://sftp.kew.org/pub/treeoflife/>), whereas raw sequence reads are
611 deposited within the European Nucleotide Archive (project number PRJEB35285) for
612 integration within the Sequence Read Archive.

613 **DISCUSSION**

614

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

615 The new phylogenomic platform described here is a major milestone towards a
616 comprehensive tree of life for all flowering plant species. The sequencing of a standardised
617 nuclear marker set of this scale for so many taxa is unprecedented, opening doors to a highly
618 integrated future for angiosperm phylogenetics in the genomic era. Much like a “next
619 generation” *rbcL*, which underpinned so many Sanger sequencing-based plant phylogenetic
620 studies, the Angiosperms353 genes offer opportunities for continuous synthesis of HTS data
621 across angiosperms. The foundational dataset presented here can be re-used or extended for
622 tree of life research at almost any taxonomic scale (Johnson et al. 2019; Larridon et al. 2019;
623 Van Andel et al. 2019; Murphy et al. 2020; Pérez-Escobar et al. 2020; Shee et al. 2020; Slimp
624 et al. 2020; Beck et al. 2021). This is the first phylogenetic project to gather novel HTS data
625 across angiosperms with a stratified taxon sampling at the genus level. Our sampling strategy
626 systematically and comprehensively represents both the diversity of angiosperms and their
627 deep-time diversification. As genus-level sampling becomes increasingly complete—a target
628 that is well within reach—this backbone will substantially increase our ability to study the
629 dynamics of plant diversity over time and revisit long-standing questions in systematics
630 (Magallón et al. 2018; Sauquet and Magallón 2018; Soltis et al. 2019). Importantly, it will
631 also sharpen the focus on truly intractable phylogenetic problems (Yang et al. 2020; Zhao et
632 al. 2020), encouraging the exploration of the biological drivers of these phenomena.

633 Our approach has already led to a burst of community engagement. More than a
634 dozen studies utilising Angiosperms353 probes are already published (e.g. Larridon et al.
635 2019; Howard et al. 2020; Murphy et al. 2020; Pérez-Escobar et al. 2020; Shee et al. 2020;
636 Slimp et al. 2020; McLay et al. in press), and two journal special issues focused on the probe
637 set are in preparation (Baker et al. in press) arising from a recent symposium (Lagomarsino
638 and Jabaily 2020). The probe set has also been adopted by the Genomics for Australian Plants
639 consortium (<https://www.genomicsforaustralianplants.com/>), which aims to sequence all

Baker et al.

640 Australian angiosperm genera, coordinating with the PAFTOL project to optimise collective
641 taxonomic coverage. A subset of the Angiosperms353 genes is now accessible for non-
642 angiosperm land plants thanks to a probe set developed in parallel (Breinholt et al. 2021),
643 inviting the prospect of data integration across all land plants. Angiosperms353 genes (as
644 distinct from the Angiosperms353 probes) are also being leveraged as components of custom-
645 designed probe sets (e.g. Jantzen et al. 2020; Ogutcen et al. 2021). This approach gives all the
646 integrative benefits of Angiosperms353, while permitting (i) the tailoring of Angiosperms353
647 probes to a taxonomic group by using more specific target data to increase gene recovery, and
648 (ii) the inclusion of additional loci pertinent to the research in question. Angiosperms353
649 probes have also been directly combined with an existing custom probe set (Nikolov et al.
650 2019) as a “probe cocktail” in a single hybridisation, capturing both sets of targets
651 simultaneously with remarkable efficiency (Hendriks et al. in press). These possibilities
652 render the invidious choice between specific and universal probe sets increasingly irrelevant
653 (Kadlec et al. 2017).

654 Although target sequence capture is the most cost-effective way to retrieve the
655 Angiosperms353 genes at the current time, the opportunity to mine the genes from other
656 kinds of HTS data (e.g. shotgun sequence data, RNA sequence data) should not be
657 overlooked. This represents a further opportunity for community engagement, both via
658 mining of public data in the Sequence Read Archive, for example, and by adding value to
659 new data being generated with these methods. A stronger understanding of the sequencing
660 requirements (e.g. coverage) for gene recovery from such data could guide new data
661 generation so that Angiosperms353 genes can be retrieved routinely as a by-product of other
662 research. We took several open data measures to encourage community uptake, in both the
663 design of our tools and the sharing of our data. The Angiosperms353 probe set itself was
664 designed to be a transparent, “off-the-shelf” toolkit that is open, inexpensive and accessible to

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

665 all, especially researchers discouraged by the complexity and cost of custom probe design
666 (Johnson et al. 2019). Our sequence data for Angiosperms353 genes are openly available via
667 the Kew Tree of Life Explorer and the Sequence Read Archive, as a public foundation dataset
668 shared according to pre-publication best practice (Toronto International Data Release
669 Workshop Authors 2009). The Explorer offers enhanced transparency and accessibility by
670 allowing users to navigate the data via a phylogenetic snapshot of the current release, along
671 with metadata (e.g. specimen data) and intermediate data (e.g. gene assemblies, alignments,
672 gene trees). Thanks to these resources, cross-community collaboration via Angiosperms353 is
673 gaining momentum.

674 Our tree, which is based on the most extensive nuclear phylogenomic dataset in
675 flowering plants to date, is strongly supported, credible and highly congruent with existing
676 taxonomy and many hypothesized relationships among orders (APG IV 2016; Supplementary
677 Fig. S6). The data confirm both the effectiveness of Angiosperms353 probes across all major
678 angiosperm clades and the ability of the genes to resolve relationships across taxonomic
679 scales (Supplementary Fig. S5). Variable sequence recovery notwithstanding (Table 2,
680 Supplementary Fig. S3), most nodes in our tree are underpinned by large numbers of gene
681 trees (Supplementary Fig. S5e), allowing the species tree to be inferred with confidence
682 (Supplementary Fig. S5a) despite gene tree conflict (Supplementary Fig. S5c). However,
683 even the most strongly supported phylogenetic hypotheses must be viewed with caution as
684 they may be biased by model misspecification and wrong assumptions. Moreover, our “first
685 pass” analyses based on a set of standard methods may not suit this dataset perfectly (see
686 below). Nevertheless, our findings are rendered credible by their high concordance with
687 taxonomy, an independent point of reference that has been extensively ground-truthed by pre-
688 phylogenomic DNA data, especially plastid loci. Agreement with existing family
689 circumscriptions is particularly striking. In contrast, congruence with previously

Baker et al.

690 hypothesized relationships among orders (APG IV 2016) is much lower (Supplementary Fig.
691 S6). Some of these earlier hypothesized ordinal relationships derive from relatively weak
692 evidence (bootstrap/jackknife >50%; APG IV 2016), which may partly explain this
693 disagreement. However, it may also be due to phylogenetic conflict between nuclear and
694 plastid genomes, as the established ordinal relationships rest primarily on evidence from
695 plastid loci, substantiated more recently by plastid genomes (Li et al. 2019). It is hardly
696 surprising, then, that a large-scale nuclear analysis presents strongly supported, alternative
697 relationships (Supplementary Fig. S6). The conundrum remains that these incongruences are
698 visible at the ordinal backbone, but not the family level. A more comprehensive exploration
699 of these relationships, the underlying phylogenetic signal and their systematic implications is
700 currently underway.

701 The analyses presented here are primarily intended as a window onto the information
702 content of our current data release and are not a complete exploration of the data. Thus,
703 downstream application of the current species tree comes with caveats. We used current,
704 widely accepted methods in a pipeline that can be re-run in a semi-automated fashion
705 whenever we release new data. As a consequence, not all possible analysis options and
706 effects have been explored. We anticipate that users of our data will probe it more rigorously
707 and will tailor both sampling and phylogenomic analyses to their specific questions. For
708 example, users may leverage our data by enriching a subset with denser sampling of their
709 own to address more focused evolutionary questions. A further exemplar use case could be
710 deeper re-analysis of our data from raw sequence reads to deepen understanding of gene
711 history and conflict.

712 Important limitations in our analysis relate to (i) taxon sampling, (ii) gene selection
713 (ii) gene recovery, (iii) models of sequence evolution and (iv) paralogy. Taxon sampling for
714 intermediate data releases is biased by the current state of progress towards our systematic

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

715 sampling strategy. This will be addressed in future data releases and can be adjusted by users
716 of our data. In addition, potential phylogenetic biases attributable to the function or other
717 properties of the Angiosperms353 genes remain poorly understood and require further
718 investigation. Gene recovery relied upon the standard Angiosperms353 target file (Johnson et
719 al. 2019), which, by its universal nature, can yield patchy results. However, it has recently
720 been reported that tailoring target sequences to specific taxonomic groups can improve
721 recovery (McLay et al. in press); this will be tested in future data releases. Moreover, we are
722 yet to exploit intronic data captured in the “splash zone” adjacent to our target exons. By
723 necessity, our “first pass” phylogenetic analysis does not explore the fast-evolving spectrum
724 of methodological options available for phylogenomic analysis. For example, we rely on a
725 simple standard model of sequence evolution, but more sophisticated models accounting for
726 codon positions or amino acids may improve phylogenetic inference. Potential paralogy is
727 not addressed by our current pipeline. The genes underpinning our analysis were carefully
728 chosen to represent single-copy genes across flowering plants (Johnson et al. 2019; Leebens-
729 Mack et al. 2019). The very low proportion of ambiguous bases across all gene alignments
730 (0.01%; Table S3) suggests that gene assembly was not strongly impacted by divergent gene
731 copies, such as paralogs. However, some paralogy may have gone unnoticed due to the
732 pervasiveness of gene and genome duplication in plants (Li and Barker 2020). Overall, we
733 expect that the occasional presence of paralogs in our current analysis would more likely lead
734 to inflated estimates of gene tree incongruence, and thus result in reduced support values,
735 than significant topological biases (Yan et al. 2020). Thus, we consider our tree relatively
736 conservative while acknowledging that we are not yet exploiting the full potential of our data.
737 Although a rigorous analysis of paralogy in Angiosperms353 genes was not tractable for this
738 data release, we look forward to deeper insights emerging as community-wide engagement
739 with Angiosperms353 grows.

Baker et al.

740 **PROSPECTS**

741

742 In the immediate future, we will deliver a further data release through which we
743 expect to reach the milestone of sampling 50% of all angiosperm genera. This target will be
744 achieved through substantial novel data production by PAFTOL and collaborators,
745 augmented by data mined from public sources. In-depth phylogenetic analyses of our data
746 and their evolutionary implications are also underway.

747 Beyond this point, we see three priority areas in which future platform developments
748 might be concentrated, resources permitting. Firstly, taxon sampling to the genus level must
749 be completed. Our original target of sampling all angiosperm genera remains, but the mode of
750 reaching this is likely to evolve. We anticipate an acceleration in production of
751 Angiosperms353 data by the broader community. The completion of generic-level sampling
752 will require both the integration of community data in the broader angiosperm tree of life as
753 well as strategic investment in filling inevitable data gaps for orphan groups. Secondly,
754 numerous opportunities for refinement exist across our methods. For example, insights from
755 our data might permit the optimisation of the Angiosperms353 probes to improve gene
756 capture. Efficiency of gene assembly from sequence data can also be improved
757 bioinformatically (McLay et al. in press). However, as costs of sequencing decline, target
758 sequence capture *in vitro* may no longer be necessary, the target genes simply being mined
759 from sufficiently deeply sequenced genomes. Thirdly, for the full integrative potential of
760 Angiosperms353 genes to be achieved, infrastructure for aggregating and sharing this
761 coherent body of data must be improved. While the Kew Tree of Life Explorer provides a
762 proof-of-concept, it is the public data repositories (e.g. NCBI, ENA) that offer the greatest
763 prospects of a mechanism to achieve this. To fully parallel the earlier success of public
764 repositories for facilitating single-gene phylogenetic trees (e.g. *rbcL*, *matK*), new tools are

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

765 needed to assist with efficient upload and annotation of target capture loci and associated
766 metadata.

767 Even with a completed genus-level angiosperm tree of life well within reach, the
768 monumental task of sampling all species remains. The scale of this challenge is 24-fold
769 greater than the genus-level tree towards which we are currently working. However, with
770 sufficient investment, increased efficiencies and community engagement, such an ambition
771 could potentially be realised. Collections-based institutions are poised to play a critical role in
772 this endeavour through increasingly routine molecular characterisation of their specimens,
773 perhaps as part of digitisation programmes, and are already facilitating the growing trend
774 towards species-complete sampling in phylogenomic studies (e.g. Loiseau et al. 2019;
775 Murphy et al. 2020; Kuhnhäuser et al. 2021). Our platform demonstrates how large-scale
776 phylogenomic projects can capitalise on natural history collections to achieve a much more
777 complete sampling than hitherto possible.

778 The growing movement to sequence the genomes of all life on Earth, inspired by the
779 Earth Biogenome Project (Lewin et al. 2018), significantly boosts the prospects for
780 completing the tree of life for all species, but is hampered by the focus on “gold standard”
781 whole genomes requiring the highest quality input DNA. Our platform offers the opportunity
782 to bridge the gap between the ambition of these projects and the vast phylogenomic potential
783 of natural history collections. However, as life on Earth becomes increasingly imperilled, we
784 cannot afford to wait. To meet the urgent demand for best estimates of the tree of life, we
785 must dynamically integrate phylogenetic information as it is generated, providing synthetic
786 trees of life to the broadest community of potential users (Eiserhardt et al. 2018). Our
787 platform facilitates this crucial synthesis by providing a cross-cutting dataset and directing
788 the community towards universal markers that seem set to play a central role in completing
789 an integrated angiosperm tree of life.

Baker et al.

790

791 **DATA AVAILABILITY AND SUPPLEMENTARY MATERIAL**

792

793 All data generated in this study are publicly released under a Creative Commons
794 Attribution 4.0 International (CC BY 4.0) license and the Toronto guidelines on pre-
795 publication data sharing (Toronto International Data Release Workshop Authors 2009). The
796 data are accessible via the Kew Tree of Life Explorer (<https://treeoflife.kew.org>) and our
797 secure FTP (<http://sftp.kew.org/pub/treeoflife/>). Raw sequence reads are deposited in the
798 European Nucleotide Archive (<https://www.ebi.ac.uk/ena/browser/home>) under umbrella
799 project PRJEB35285. Scripts and other files relating to our phylogenomic pipeline are
800 available at our GitHub (<https://github.com/RBGKew/KewTreeOfLife>). Supplementary
801 materials cited in this paper plus Data Release 1.0 datasets duplicated from our secure FTP
802 (assembled genes, gene alignments, gene trees, species tree, examples of scripts) are available
803 from the Dryad Digital Repository (<https://doi.org/10.5061/dryad.ns1rn8ps7>).

804 **FUNDING**

805

806 This work was supported by grants from the Calleva Foundation and the Sackler Trust
807 to the Plant and Fungal Trees of Life project at the Royal Botanic Gardens, Kew. Additional
808 support was received from the Garfield Weston Foundation, as part of the Global Tree Seed
809 Bank Programme.

810 **ACKNOWLEDGEMENTS**

811 Numerous people have supported this work through collaboration, sharing expertise,
812 providing samples, supporting acquisition of samples from RBGK collections, and assisting

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

813 with laboratory work, data collection, specimen digitisation and computational infrastructure.

814 A full list of acknowledgements is given in Supplementary Material. Particular thanks go to

815 Kathy Willis, former Director of Science at the Royal Botanic Gardens, Kew, for inspiring

816 the establishment of the PAFTOL project.

817 **LITERATURE CITED**

818

819 Abadi S., Azouri D., Pupko T., Mayrose I. 2019. Model selection may not be a mandatory
820 step for phylogeny reconstruction. *Nat. Commun.* 10:934.

821

822 Alsos I.G., Lavergne S., Merkel M.K., Boleda M., Lammers Y., Alberti A., Pouchon C.,
823 Denoeud F., Pitelkova I., Puşcaş M., Roquet C., Hurdu B.-I., Thuiller W., Zimmermann N.E.,
824 Hollingsworth P.M., Coissac E. 2020. The treasure vault can be opened: Large-scale genome
825 skimming works well using herbarium and silica gel dried material. *Plants* 9:432.

826

827 Antonelli A., Fry C., Smith R.J., Simmonds M.S.J., Kersey P.J., Pritchard H.W., Abbo M.S.,
828 Acedo C., Adams J., Ainsworth A.M., Allkin B., Annecke W., Bachman S.P., Bacon K.,
829 Bárrios S., Barstow C., Battison A., Bell E., Bensusan K., Bidartondo M.I., Blackhall-Miles
830 R.J., Borrell J.S., Brearley F.Q., Breman E., Brewer R.F.A., Brodie J., Cámara-Leret R.,
831 Campostrini Forzza R., Cannon P., Carine M., Carretero J., Cavagnaro T.R., Cazar M.E.,
832 Chapman T., Cheek M., Clubbe C., Cockel C., Collemare J., Cooper A., Copeland A.I.,
833 Corcoran M., Couch C., Cowell C., Crous P., da Silva M., Dalle G., Das D., David J.C.,
834 Davies L., Davies N., De Canha M.N., de Lirio E.J., Demissew S., Diazgranados M., Dickie
835 J., Dines T., Douglas B., Dröge G., Dulloo M.E., Fang R., Farlow A., Farrar K., Fay M.F.,
836 Felix J., Forest F., Forrest L.L., Fulcher T., Gafforov Y., Gardiner L.M., Gâteblé G., Gaya E.,
837 Geslin B., Gonçalves S.C., Gore C.J.N., Govaerts R., Gowda B., Grace O.M., Grall A.,

Baker et al.

- 838 Haelewaters D., Halley J.M., Hamilton M.A., Hazra A., Heller T., Hollingsworth P.M.,
 839 Holstein N., Howes M.J.R., Hughes M., Hunter D., Hutchinson N., Hyde K., Iganci J., Jones
 840 M., Kelly L.J., Kirk P., Koch H., Grisai-Greilhuber I., Lall N., Langat M.K., Leaman D.J.,
 841 Leão T.C., Lee M.A., Leitch I.J., Leon C., Lettice E., Lewis G.P., Li L., Lindon H., Liu J.S.,
 842 Liu U., Llewellyn T., Looney B., Lovett J.C., Luczaj L., Lulekal E., Maggassouba S.,
 843 Malécot V., Martin C., Masera O.R., Mattana E., Maxted N., Mba C., McGinn K.J.,
 844 Metheringham C., Miles S., Miller J., Milliken W., Moat J., Moore P.G.P., Morim M.P.,
 845 Mueller G.M., Muminjanov H., Negrão R., Nic Lughadha E., Nicholson N., Niskanen T.,
 846 Nono Womdim R., Noorani A., Obreza M., O'Donnell K., O'Hanlon R., Onana J.M., Ondo I.,
 847 Padulosi S., Paton A., Pearce T., Pérez Escobar O.A., Pieroni A., Pironon S., Prescott T.A.K.,
 848 Qi Y.D., Qin H., Quave C.L., Rajaovelona L., Razanajatovo H., Reich P.B., Rianawati E.,
 849 Rich T.C.G., Richards S.L., Rivers M.C., Ross A., Rumsey F., Ryan M., Ryan P., Sagala S.,
 850 Sanchez M.D., Sharrock S., Shrestha K.K., Sim J., Sirakaya A., Sjöman H., Smidt E.C.,
 851 Smith D., Smith P., Smith S.R., Sofo A., Spence N., Stanworth A., Stara K., Stevenson P.C.,
 852 Stroh P., Suz L.M., Tambam B.B., Tatsis E.C., Taylor I., Thiers B., Thormann I., Vaglica V.,
 853 Vásquez-Londoño C., Victor J., Viruel J., Walker B.E., Walker K., Walsh A., Way M.,
 854 Wilbraham J., Wilkin P., Wilkinson T., Williams C., Winterton D., Wong K.M., Woodfield-
 855 Pascoe N., Woodman J., Wyatt L., Wynberg R., Zhang B.G. 2020. State of the World's Plants
 856 and Fungi 2020. Royal Botanic Gardens, Kew.
- 857
- 858 APG. 1998. An ordinal classification for the families of flowering plants. *Ann. Missouri Bot.*
 859 *Gard.* 85:531-553.
- 860
- 861 APG II. 2003. An update of the Angiosperm Phylogeny Group classification for the orders
 862 and families of flowering plants: Apg II. *Bot. J. Linn. Soc.* 141:399-436.

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

863

864 APG III. 2009. An update of the Angiosperm Phylogeny Group classification for the orders
865 and families of flowering plants: Apg III. Bot. J. Linn. Soc. 161:105-121.

866

867 APG IV. 2016. An update of the Angiosperm Phylogeny Group classification for the orders
868 and families of flowering plants: Apg IV. Bot. J. Linn. Soc. 181:1-20.

869

870 Baker W.J., Dodsworth S., Forest F., Graham S.W., Johnson M.G., McDonnell A., Pokorny
871 L., Tate J.A., Wicke S., Wickett N.J. 2021. Exploring Angiosperms353: An open, community
872 toolkit for collaborative phylogenomic research on flowering plants. Amer. J. Bot. In Press.

873

874 Bakker F.T., Lei D., Yu J., Mohammadin S., Wei Z., van de Kerke S., Gravendeel B.,
875 Nieuwenhuis M., Staats M., Alquezar-Planas D.E., Holmer R. 2016. Herbarium genomics:
876 Plastome sequence assembly from a range of herbarium specimens using an iterative
877 organelle genome assembly pipeline. Biol. J. Linn. Soc. 117:33-43.

878

879 Beck J.B., Markley M.L., Zielke M.G., Thomas J.R., Hale H.J., Williams L.D., Johnson M.G.
880 2021. Is Palmer's elm leaf goldenrod real? The Angiosperms353 kit provides within-species
881 signal in *Solidago ulmifolia* s.L. bioRxiv:2021.2001.2007.425781.

882

883 Bolger A.M., Lohse M., Usadel B. 2014. Trimmomatic: A flexible trimmer for Illumina
884 sequence data. Bioinformatics 30:2114-2120.

885

886 Borowiec M.L. 2016. AMAS: A fast tool for alignment manipulation and computing of
887 summary statistics. PeerJ 4:e1660.

Baker et al.

888

889 Bostock M. 2012. D3.Js - data-driven documents <http://d3js.Org/>.

890

891 Breinholt J.W., Carey S.B., Tiley G.P., Davis E.C., Endara L., McDaniel S.F., Neves L.G.,

892 Sessa E.B., von Konrat M., Chantanaorrapint S., Fawcett S., Ickert-Bond S.M., Labiak P.H.,

893 Larraín J., Lehnert M., Lewis L.R., Nagalingum N.S., Patel N., Rensing S.A., Testo W.,

894 Vasco A., Villarreal J.C., Williams E.W., Burleigh J.G. 2021. A target enrichment probe set

895 for resolving the flagellate land plant tree of life. *Appl. Plant. Sci.* n/a:e11406.

896

897 Brewer G.E., Clarkson J.J., Maurin O., Zuntini A.R., Barber V., Bellot S., Biggs N., Cowan

898 R.S., Davies N.M.J., Dodsworth S., Edwards S.L., Eiserhardt W.L., Epiawalage N., Frisby

899 S., Grall A., Kersey P.J., Pokorny L., Leitch I.J., Forest F., Baker W.J. 2019. Factors

900 affecting targeted sequencing of 353 nuclear genes from herbarium specimens spanning the

901 diversity of angiosperms. *Front. Plant Sci.* 10:1102.

902

903 Buddenhagen C., Lemmon A.R., Lemmon E.M., Bruhl J., Cappa J., Clement W.L.,

904 Donoghue M.J., Edwards E.J., Hipp A.L., Kertyna M. 2016. Anchored phylogenomics of

905 angiosperms I: Assessing the robustness of phylogenetic estimates. *bioRxiv:086298*.

906

907 Buerki S., Baker W.J. 2016. Collections-based research in the genomic era. *Biol. J. Linn.*

908 *Soc.* 117:5-10.

909

910 Camacho C., Coulouris G., Avagyan V., Ma N., Papadopoulos J., Bealer K., Madden T.L.

911 2009. BLAST+: Architecture and applications. *BMC Bioinformatics* 10:421.

912

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 913 Carpenter E.J., Matasci N., Ayyampalayam S., Wu S., Sun J., Yu J., Jimenez Vieira F.R.,
914 Bowler C., Dorrell R.G., Gitzendanner M.A., Li L., Du W., K. Ullrich K., Wickett N.J.,
915 Barkmann T.J., Barker M.S., Leebens-Mack J.H., Wong G.K.-S. 2019. Access to rna-
916 sequencing data from 1,173 plant species: The 1000 Plant Transcriptomes Initiative (1KP).
917 GigaScience 8:giz126.
- 918
- 919 Chase M.W., Hills H.H. 1991. Silica gel: An ideal material for field preservation of leaf
920 samples for DNA studies. *Taxon* 40:215-220.
- 921
- 922 Chase M.W., Soltis D.E., Olmstead R.G., Morgan D., Les D.H., Mishler B.D., Duvall M.R.,
923 Price R.A., Hills H.G., Qiu Y.L., Kron K.A., Rettig J.H., Conti E., Palmer J.D., Manhart J.R.,
924 Sytsma K.J., Michaels H.J., Kress W.J., Karol K.G., Clark W.D., Hedren M., Gaut B.S.,
925 Jansen R.K., Kim K.J., Wimpee C.F., Smith J.F., Furnier G.R., Strauss S.H., Xiang Q.Y.,
926 Plunkett G.M., Soltis P.S., Swensen S.M., Williams S.E., Gadek P.A., Quinn C.J., Eguiarte
927 L.E., Golenberg E., Learn G.H., Graham S.W., Barrett S.C.H., Dayanandan S., Albert V.A.
928 1993. Phylogenetics of seed plants - an analysis of nucleotide sequences from the plastid
929 gene *rbcL*. *Ann. Missouri Bot. Gard.* 80:528-580.
- 930
- 931 Chau J.H., Rahfeldt W.A., Olmstead R.G. 2018. Comparison of taxon-specific versus general
932 locus sets for targeted sequence capture in plant phylogenomics. *Appl. Plant. Sci.* 6:e1032.
- 933
- 934 Cheng S., Melkonian M., Smith S.A., Brockington S., Archibald J.M., Delaux P.-M., Li F.-
935 W., Melkonian B., Mavrodiev E.V., Sun W., Fu Y., Yang H., Soltis D.E., Graham S.W.,
936 Soltis P.S., Liu X., Xu X., Wong G.K.-S. 2018. 10kp: A phylodiverse genome sequencing
937 plan. *GigaScience* 7:giy013.

Baker et al.

938

939 Cornwell W.K., Pearse W.D., Dalrymple R.L., Zanne A.E. 2019. What we (don't) know
940 about global plant diversity. *Ecography* 42:1819-1831.

941

942 Couvreur T.L.P., Helmstetter A.J., Koenen E.J.M., Bethune K., Brandão R.D., Little S.A.,
943 Sauquet H., Erkens R.H.J. 2019. Phylogenomics of the major tropical plant family
944 Annonaceae using targeted enrichment of nuclear genes. *Front. Plant Sci.* 9:1941.

945

946 Dodsworth S., Pokorny L., Johnson M.G., Kim J.T., Maurin O., Wickett N.J., Forest F.,
947 Baker W.J. 2019. Hyb-Seq for flowering plant systematics. *Trends Plant Sci.* 24:887-891.

948

949 Doyle J.J., Doyle J.L. 1987. A rapid DNA isolation procedure from small quantities of fresh
950 leaf tissue. *Phytochem. Bull.* 19:11-15.

951

952 Eiserhardt W.L., Antonelli A., Bennett D.J., Botigué L.R., Burleigh J.G., Dodsworth S.,
953 Enquist B.J., Forest F., Kim J.T., Kozlov A.M., Leitch I.J., Maitner B.S., Mirarab S., Piel
954 W.H., Pérez-Escobar O.A., Pokorny L., Rahbek C., Sandel B., Smith S.A., Stamatakis A.,
955 Vos R.A., Warnow T., Baker W.J. 2018. A roadmap for global synthesis of the plant tree of
956 life. *Amer. J. Bot.* 105:614-622.

957

958 Faircloth B.C., McCormack J.E., Crawford N.G., Harvey M.G., Brumfield R.T., Glenn T.C.
959 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple
960 evolutionary timescales. *Syst. Biol.* 61:717-726.

961

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 962 Forrest L.L., Hart M.L., Hughes M., Wilson H.P., Chung K.-F., Tseng Y.-H., Kidner C.A.
963 2019. The limits of Hyb-Seq for herbarium specimens: Impact of preservation techniques.
964 *Front. Ecol. Evol.* 7:439.
- 965
- 966 Gitzendanner M.A., Soltis P.S., Wong G.K.-S., Ruhfel B.R., Soltis D.E. 2018. Plastid
967 phylogenomic analysis of green plants: A billion years of evolutionary history. *Amer. J. Bot.*
968 105:291-301.
- 969
- 970 Hale H., Gardner E.M., Viruel J., Pokorny L., Johnson M.G. 2020. Strategies for reducing
971 per-sample costs in target capture sequencing for phylogenomics and population genomics in
972 plants. *Appl. Plant. Sci.* 8:e11337.
- 973
- 974 Hendriks K., Mandáková T., Hay N.M., Ly E., Hooft van Huysduynen A., Tamrakar R.,
975 Thomas S.K., Toro-Núñez O., Pires J.C., Nikolov L.A., Koch M.A., Windham M.D., Lysak
976 M.A., Forest F., Mummenhoff K., Baker W.J., Lens F., Bailey C.D. in press. The best of both
977 worlds: Combining lineage specific and universal bait sets in target enrichment hybridization
978 reactions. *Appl. Plant. Sci.*
- 979
- 980 Hinchliff C.E., Smith S.A. 2014. Some limitations of public sequence data for phylogenetic
981 inference (in plants). *PLoS ONE* 9:e98986.
- 982
- 983 Hinchliff C.E., Smith S.A., Allman J.F., Burleigh J.G., Chaudhary R., Coghill L.M., Crandall
984 K.A., Deng J., Drew B.T., Gazis R., Gude K., Hibbett D.S., Katz L.A., Laughinghouse H.D.,
985 McTavish E.J., Midford P.E., Owen C.L., Ree R.H., Rees J.A., Soltis D.E., Williams T.,

Baker et al.

- 986 Cranston K.A. 2015. Synthesis of phylogeny and taxonomy into a comprehensive tree of life.
987 Proc. Natl. Acad. Sci. U.S.A. 112:12764.
988
- 989 Hoang D.T., Chernomor O., von Haeseler A., Minh B.Q., Vinh L.S. 2017. UFBoot2:
990 Improving the ultrafast bootstrap approximation. Mol. Biol. Evol. 35:518-522.
991
- 992 Howard C.C., Crowl A.A., Harvey T.S., Cellinese N. 2020. Peeling back the layers: The
993 complex dynamics shaping the evolution of the Ledebouriinae (Scilloideae, Asparagaceae).
994 bioRxiv:2020.2011.2002.365718.
995
- 996 Jantzen J.R., Amarasinghe P., Folk R.A., Reginato M., Michelangeli F.A., Soltis D.E.,
997 Cellinese N., Soltis P.S. 2020. A two-tier bioinformatic pipeline to develop probes for target
998 capture of nuclear loci with applications in Melastomataceae. Appl. Plant. Sci. 8:e11345.
999
- 1000 Jin J.-J., Yu W.-B., Yang J.-B., Song Y., dePamphilis C.W., Yi T.-S., Li D.-Z. 2020.
1001 GetOrganelle: A fast and versatile toolkit for accurate de novo assembly of organelle
1002 genomes. Genome Biol. 21:241.
1003
- 1004 Johnson M.G., Gardner E.M., Liu Y., Medina R., Goffinet B., Shaw A.J., Zerega N.J.C.,
1005 Wickett N.J. 2016. Hybpiper: Extracting coding sequence and introns for phylogenetics from
1006 high-throughput sequencing reads using target enrichment. Appl. Plant. Sci. 4:1600016.
1007
- 1008 Johnson M.G., Pokorny L., Dodsworth S., Botigue L.R., Cowan R.S., Devault A., Eiserhardt
1009 W.L., Epiawalage N., Forest F., Kim J.T., Leebens-Mack J.H., Leitch I.J., Maurin O., Soltis
1010 D.E., Soltis P.S., Wong G.K., Baker W.J., Wickett N.J. 2019. A universal probe set for

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 1011 targeted sequencing of 353 nuclear genes from any flowering plant designed using k-medoids
1012 clustering. *Syst. Biol.* 68:594-606.
1013
- 1014 Junier T., Zdobnov E.M. 2010. The newick utilities: High-throughput phylogenetic tree
1015 processing in the Unix shell. *Bioinformatics* 26:1669-1670.
1016
- 1017 Kadlec M., Bellstedt D.U., Le Maitre N.C., Pirie M.D. 2017. Targeted NGS for species level
1018 phylogenomics: “Made to measure” or “one size fits all”? *PeerJ* 5:e3569.
1019
- 1020 Kreft L., Botzki A., Coppens F., Vandepoele K., Van Bel M. 2017. Phyd3: A phylogenetic
1021 tree viewer with extended phyloXML support for functional genomics data visualization.
1022 *Bioinformatics* 33:2946-2947.
1023
- 1024 Kuhnhäuser B.G., Bellot S., Couvreur T.L.P., Dransfield J., Henderson A., Schley R.,
1025 Chomicki G., Eiserhardt W.L., Hiscock S.J., Baker W.J. 2021. A robust phylogenomic
1026 framework for the calamoid palms. *Mol. Phylogenet. Evol.*:107067.
1027
- 1028 Lagomarsino L.P., Jabaily R.S. 2020. Virtual Botany Conference 2020 symposium -
1029 Angiosperms353: A new essential tool for plant systematics.
1030 <http://2020.botanyconference.org/engine/search/index.php?func=detail&aid=941>.
1031
- 1032 Larridon I., Villaverde T., Zuntini A.R., Pokorny L., Brewer G.E., Epiawalage N., Fairlie I.,
1033 Hahn M., Kim J., Maguilla E., Maurin O., Xanthos M., Hipp A.L., Forest F., Baker W.J.
1034 2019. Tackling rapid radiations with targeted sequencing. *Front Plant Sci* 10:1655.
1035

Baker et al.

- 1036 Leebens-Mack J.H., Barker M.S., Carpenter E.J., Deyholos M.K., Gitzendanner M.A.,
 1037 Graham S.W., Grosse I., Li Z., Melkonian M., Mirarab S., Porsch M., Quint M., Rensing
 1038 S.A., Soltis D.E., Soltis P.S., Stevenson D.W., Ullrich K.K., Wickett N.J., DeGironimo L.,
 1039 Edger P.P., Jordon-Thaden I.E., Joya S., Liu T., Melkonian B., Miles N.W., Pokorny L.,
 1040 Quigley C., Thomas P., Villarreal J.C., Augustin M.M., Barrett M.D., Baucom R.S., Beerling
 1041 D.J., Benstein R.M., Biffin E., Brockington S.F., Burge D.O., Burris J.N., Burris K.P.,
 1042 Burtet-Sarramegna V., Caicedo A.L., Cannon S.B., Çebi Z., Chang Y., Chater C., Cheeseman
 1043 J.M., Chen T., Clarke N.D., Clayton H., Covshoff S., Crandall-Stotler B.J., Cross H.,
 1044 dePamphilis C.W., Der J.P., Determann R., Dickson R.C., Di Stilio V.S., Ellis S., Fast E.,
 1045 Feja N., Field K.J., Filatov D.A., Finnegan P.M., Floyd S.K., Fogliani B., García N., Gâteblé
 1046 G., Godden G.T., Goh F., Greiner S., Harkess A., Heaney J.M., Helliwell K.E., Heyduk K.,
 1047 Hibberd J.M., Hodel R.G.J., Hollingsworth P.M., Johnson M.T.J., Jost R., Joyce B., Kapralov
 1048 M.V., Kazamia E., Kellogg E.A., Koch M.A., Von Konrat M., Könyves K., Kutchan T.M.,
 1049 Lam V., Larsson A., Leitch A.R., Lentz R., Li F.-W., Lowe A.J., Ludwig M., Manos P.S.,
 1050 Mavrodiev E., McCormick M.K., McKain M., McLellan T., McNeal J.R., Miller R.E.,
 1051 Nelson M.N., Peng Y., Ralph P., Real D., Riggins C.W., Ruhsam M., Sage R.F., Sakai A.K.,
 1052 Scascitella M., Schilling E.E., Schlösser E.-M., Sederoff H., Servick S., Sessa E.B., Shaw
 1053 A.J., Shaw S.W., Sigel E.M., Skema C., Smith A.G., Smithson A., Stewart C.N.,
 1054 Stinchcombe J.R., Szövényi P., Tate J.A., Tiebel H., Trapnell D., Villegente M., Wang C.-N.,
 1055 Weller S.G., Wenzel M., Weststrand S., Westwood J.H., Whigham D.F., Wu S., Wulff A.S.,
 1056 Yang Y., Zhu D., Zhuang C., Zuidof J., Chase M.W., Pires J.C., Rothfels C.J., Yu J., Chen
 1057 C., Chen L., Cheng S., Li J., Li R., Li X., Lu H., Ou Y., Sun X., Tan X., Tang J., Tian Z.,
 1058 Wang F., Wang J., Wei X., Xu X., Yan Z., Yang F., Zhong X., Zhou F., Zhu Y., Zhang Y.,
 1059 Ayyampalayam S., Barkman T.J., Nguyen N.-p., Matasci N., Nelson D.R., Sayyari E.,
 1060 Wafula E.K., Walls R.L., Warnow T., An H., Arrigo N., Baniaga A.E., Galuska S., Jorgensen

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 1061 S.A., Kidder T.I., Kong H., Lu-Irving P., Marx H.E., Qi X., Reardon C.R., Sutherland B.L.,
1062 Tiley G.P., Welles S.R., Yu R., Zhan S., Gramzow L., Theißen G., Wong G.K.-S., One
1063 Thousand Plant Transcriptomes I. 2019. One thousand plant transcriptomes and
1064 the phylogenomics of green plants. *Nature* 574:679-685.
- 1065
- 1066 Lemmon A.R., Emme S.A., Lemmon E.M. 2012. Anchored hybrid enrichment for massively
1067 high-throughput phylogenomics. *Syst. Biol.* 61:727-744.
- 1068
- 1069 Lewin H.A., Robinson G.E., Kress W.J., Baker W.J., Coddington J., Crandall K.A., Durbin
1070 R., Edwards S.V., Forest F., Gilbert M.T.P., Goldstein M.M., Grigoriev I.V., Hackett K.J.,
1071 Haussler D., Jarvis E.D., Johnson W.E., Patrinos A., Richards S., Castilla-Rubio J.C., van
1072 Sluys M.-A., Soltis P.S., Xu X., Yang H., Zhang G. 2018. Earth Biogenome Project:
1073 Sequencing life for the future of life. *Proc. Natl. Acad. Sci. U.S.A.* 115:4325-4333.
- 1074
- 1075 Li H.-T., Yi T.-S., Gao L.-M., Ma P.-F., Zhang T., Yang J.-B., Gitzendanner M.A., Fritsch
1076 P.W., Cai J., Luo Y., Wang H., van der Bank M., Zhang S.-D., Wang Q.-F., Wang J., Zhang
1077 Z.-R., Fu C.-N., Yang J., Hollingsworth P.M., Chase M.W., Soltis D.E., Soltis P.S., Li D.-Z.
1078 2019. Origin of angiosperms and the puzzle of the Jurassic gap. *Nat. Plants* 5:461-470.
- 1079
- 1080 Li Z., Barker M.S. 2020. Inferring putative ancient whole-genome duplications in the 1000
1081 Plants (1KP) Initiative: Access to gene family phylogenies and age distributions. *GigaScience*
1082 9:giaa004.
- 1083
- 1084 Loiseau O., Olivares I., Paris M., de La Harpe M., Weigand A., Koubinova D., Rolland J.,
1085 Bacon C.D., Balslev H., Borchsenius F. 2019. Targeted capture of hundreds of nuclear genes

Baker et al.

1086 unravels phylogenetic relationships of the diverse neotropical palm tribe Geonomateae. *Front.*
1087 *Plant Sci.* 10:864.

1088

1089 Magallón S., Sánchez-Reyes L.L., Gómez-Acevedo S.L. 2018. Thirty clues to the exceptional
1090 diversification of flowering plants. *Ann. Bot.* 123:491-503.

1091

1092 Mai U., Mirarab S. 2018. TreeShrink: Fast and accurate detection of outlier long branches in
1093 collections of phylogenetic trees. *BMC Genomics* 19:272.

1094

1095 Mandel J.R., Dikow R.B., Funk V.A., Masalia R.R., Staton S.E., Kozik A., Michelmore
1096 R.W., Rieseberg L.H., Burke J.M. 2014. A target enrichment method for gathering
1097 phylogenetic information from hundreds of loci: An example from the Compositae. *Appl.*
1098 *Plant. Sci.* 2:1300085.

1099

1100 McLay T.G.B., Gunn B.F., Ning W., Tate J.A., Nauheimer L., Joyce E.M., Simpson L.,
1101 Schmidt-Lebuhn A.N., Baker W.J., Forest F., Jackson C.J. in press. New targets acquired:
1102 Improving locus recovery from the Angiosperms353 probe set. *Appl. Plant. Sci.*

1103

1104 Meyer M., Kircher M. 2010. Illumina sequencing library preparation for highly multiplexed
1105 target capture and sequencing. *Cold Spring Harbor Protocols* 2010:pdb.prot5448.

1106

1107 Minh B.Q., Schmidt H.A., Chernomor O., Schrempf D., Woodhams M.D., Von Haeseler A.,
1108 Lanfear R. 2020. Iq-tree 2: New models and efficient methods for phylogenetic inference in
1109 the genomic era. *Mol. Biol. Evol.* 37:1530-1534.

1110

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 1111 Murphy B., Forest F., Barraclough T., Rosindell J., Bellot S., Cowan R., Golos M., Jebb M.,
1112 Cheek M. 2020. A phylogenomic analysis of *Nepenthes* (Nepenthaceae). *Mol. Phylogenet.*
1113 *Evol.* 144:106668.
- 1114
- 1115 Nguyen N.-P.D., Mirarab S., Kumar K., Warnow T. 2015. Ultra-large alignments using
1116 phylogeny-aware profiles. *Genome Biol.* 16:124.
- 1117
- 1118 Nikolov L.A., Shushkov P., Nevado B., Gan X., Al-Shehbaz I.A., Filatov D., Bailey C.D.,
1119 Tsiantis M. 2019. Resolving the backbone of the Brassicaceae phylogeny for investigating
1120 trait diversity. *New Phytol.* 222:1638-1651.
- 1121
- 1122 Ogutcen E., Christe C., Nishii K., Salamin N., Möller M., Perret M. 2021. Phylogenomics of
1123 Gesneriaceae using targeted capture of nuclear genes. *Mol. Phylogenet. Evol.*:107068.
- 1124
- 1125 Pérez-Escobar O.A., Dodsworth S., Bogarín D., Bellot S., Balbuena J.A., Schley R., Kikuchi
1126 I., Morris S.K., Epiawalage N., Cowan R., Maurin O., Zuntini A., Arias T., Serna A.,
1127 Gravendeel B., Torres M.F., Nargar K., Chomicki G., Chase M.W., Leitch I.J., Forest F.,
1128 Baker W.J. 2020. Hundreds of nuclear and plastid loci yield insights into orchid relationships.
1129 *bioRxiv*:2020.2011.2017.386508.
- 1130
- 1131 RBG Kew. 2015. A global resource for plant and fungal knowledge. Science strategy 2015-
1132 2020. Royal Botanic Gardens, Kew.
- 1133
- 1134 RBG Kew. 2016. The State of the World's Plants report – 2016. Royal Botanic Gardens,
1135 Kew.

Baker et al.

1136

1137 Sauquet H., Magallón S. 2018. Key questions and challenges in angiosperm macroevolution.

1138 New Phytol. 219:1170-1187.

1139

1140 Sayyari E., Mirarab S. 2016. Fast coalescent-based computation of local branch support from

1141 quartet frequencies. Mol. Biol. Evol. 33:1654-1668.

1142

1143 Secretariat of the Convention on Biological Diversity. 2011. Nagoya protocol on access to

1144 genetic resources and the fair and equitable sharing of benefits arising from their utilization to

1145 the convention on biological diversity. Montreal: United Nations Environment Programme.

1146

1147 Shee Z.Q., Frodin D.G., Cámara-Leret R., Pokorny L. 2020. Reconstructing the complex

1148 evolutionary history of the Papuanian *Schefflera* radiation through herbariomics. Front. Plant

1149 Sci. 11:258.

1150

1151 Slimp M., Williams L.D., Hale H., Johnson M.G. 2020. On the potential of Angiosperms353

1152 for population genomics. bioRxiv:2020.2010.2011.335174.

1153

1154 Smith S.A., Brown J.W. 2018. Constructing a broadly inclusive seed plant phylogeny. Amer.

1155 J. Bot. 105:302-314.

1156

1157 Soltis D.E., Smith S.A., Cellinese N., Wurdack K.J., Tank D.C., Brockington S.F., Refulio-

1158 Rodriguez N.F., Walker J.B., Moore M.J., Carlswald B.S., Bell C.D., Latvis M., Crawley S.,

1159 Black C., Diouf D., Xi Z., Rushworth C.A., Gitzendanner M.A., Sytsma K.J., Qiu Y.-L., Hilu

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 1160 K.W., Davis C.C., Sanderson M.J., Beaman R.S., Olmstead R.G., Judd W.S., Donoghue M.J.,
1161 Soltis P.S. 2011. Angiosperm phylogeny: 17 genes, 640 taxa. *Amer. J. Bot.* 98:704-730.
1162
- 1163 Soltis D.E., Soltis P.S., Chase M.W., Mort M.E., Albach D.C., Zanis M., Savolainen V.,
1164 Hahn W.J., Hoot S.B., Fay M.F., Axtell M., Swensen S.M., Prince L.M., Kress W.J., Nixon
1165 K.C., Farris J.S. 2008. Angiosperm phylogeny inferred from 18s rDNA, *rbcL*, and *atpB*
1166 sequences. *Bot. J. Linn. Soc.* 133:381-461.
1167
- 1168 Soltis P.S., Folk R.A., Soltis D.E. 2019. Darwin review: Angiosperm phylogeny and
1169 evolutionary radiations. *Proc. R. Soc. Lond. B Biol. Sci.* 286:20190099.
1170
- 1171 Soto Gomez M., Pokorny L., Kantar M.B., Forest F., Leitch I.J., Gravendeel B., Wilkin P.,
1172 Graham S.W., Viruel J. 2019. A customized nuclear target enrichment approach for
1173 developing a phylogenomic baseline for *Dioscorea* yams (Dioscoreaceae). *Appl. Plant. Sci.*
1174 7:e11254.
1175
- 1176 Toronto International Data Release Workshop Authors. 2009. Prepublication data sharing.
1177 *Nature* 461:168-170.
1178
- 1179 Van Andel T., Veltman M.A., Bertin A., Maat H., Polime T., Hille Ris Lambers D., Tjoe
1180 Awie J., De Boer H., Manzanilla V. 2019. Hidden rice diversity in the Guianas. *Front. Plant*
1181 *Sci.* 10:1161.
1182
- 1183 Villaverde T., Pokorny L., Olsson S., Rincón-Barrado M., Johnson M.G., Gardner E.M.,
1184 Wickett N.J., Molero J., Riina R., Sanmartín I. 2018. Bridging the micro- and

Baker et al.

1185 macroevolutionary levels in phylogenomics: Hyb-Seq solves relationships from populations
1186 to species and above. *New Phytol.* 220:636-650.

1187

1188 WCVP. 2020. World Checklist of Vascular Plants, version 2.0. Facilitated by the Royal
1189 Botanic Gardens, kew. Published on the internet; <http://wcvp.science.kew.org/>, retrieved 18
1190 November 2020.

1191

1192 Weitemier K., Straub S.C.K., Cronn R.C., Fishbein M., Schmickl R., McDonnell A., Liston
1193 A. 2014. Hyb-Seq: Combining target enrichment and genome skimming for plant
1194 phylogenomics. *Appl. Plant. Sci.* 2:1400042.

1195

1196 Wickett N.J., Mirarab S., Nguyen N., Warnow T., Carpenter E., Matasci N., Ayyampalayam
1197 S., Barker M.S., Burleigh J.G., Gitzendanner M.A., Ruhfel B.R., Wafula E., Der J.P., Graham
1198 S.W., Mathews S., Melkonian M., Soltis D.E., Soltis P.S., Miles N.W., Rothfels C.J.,
1199 Pokorny L., Shaw A.J., DeGironimo L., Stevenson D.W., Surek B., Villarreal J.C., Roure B.,
1200 Philippe H., dePamphilis C.W., Chen T., Deyholos M.K., Baucom R.S., Kutchan T.M.,
1201 Augustin M.M., Wang J., Zhang Y., Tian Z., Yan Z., Wu X., Sun X., Wong G.K.-S.,
1202 Leebens-Mack J. 2014. Phylotranscriptomic analysis of the origin and early diversification of
1203 land plants. *Proc. Natl. Acad. Sci. U.S.A.* 111:E4859.

1204

1205 Yan Z., Du P., Hahn M.W., Nakhleh L. 2020. Species tree inference under the multispecies
1206 coalescent on data with paralogs is accurate. *bioRxiv*:498378.

1207

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

- 1208 Yang L., Su D., Chang X., Foster C.S.P., Sun L., Huang C.-H., Zhou X., Zeng L., Ma H.,
1209 Zhong B. 2020. Phylogenomic insights into deep phylogeny of angiosperms based on broad
1210 nuclear gene sampling. *Plant Commun.* 1:100027.
1211
- 1212 Zhang C., Rabiee M., Sayyari E., Mirarab S. 2018. ASTRAL-III: Polynomial time species
1213 tree reconstruction from partially resolved gene trees. *BMC Bioinformatics* 19:153.
1214
- 1215 Zhao T., Xue J., Kao S.-m., Li Z., Zwaenepoel A., Schranz M.E., Van de Peer Y. 2020.
1216 Novel phylogeny of angiosperms inferred from whole-genome microsynteny analysis.
1217 *bioRxiv:2020.2001.2015.908376.*
1218
1219

Baker et al.

1220 **FIGURE LEGENDS**

1221

1222 **Figure 1.** Summary workflow. Overview of steps taken by the PAFTOL project to generate
1223 Data Release 1.0 of the Kew Tree of Life Explorer (<https://treeoflife.kew.org>). The stages of
1224 the workflow are further elaborated in Figs. 2–4.

1225

1226 **Figure 2.** Sample processing and data analysis workflows. Sample processing (left):
1227 processes are indicated by bold headings with reagents and machines used given below;
1228 quality control (QC) checkpoints are indicated in dark grey boxes. Data analysis (right):
1229 pipeline products are shown in blue-green circles (available to download via the Kew Tree of
1230 Life Explorer, <https://treeoflife.kew.org>); processes are indicated by bold headings with
1231 programs used given below.

1232

1233 **Figure 3.** Family identification validation workflow. Processes are indicated by bold
1234 headings. Embedded table (bottom right) indicates decisions made for each sample based on
1235 the two validation steps.

1236

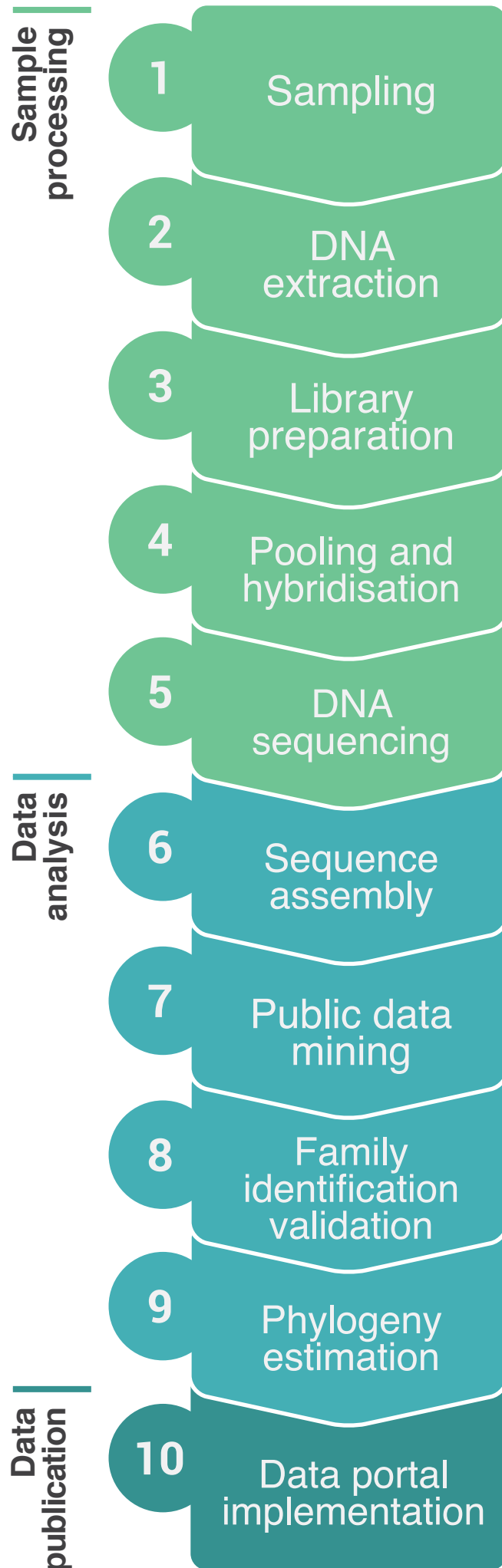
1237 **Figure 4.** Data publication workflow. Implementation of Kew Tree of Life Explorer data
1238 portal is illustrated. Arrows indicate data flow from internal repository to public interface.
1239 Infrastructural components are shown in purple; publicly available information is shown in
1240 green. External links available from the portal are listed in the lower left.

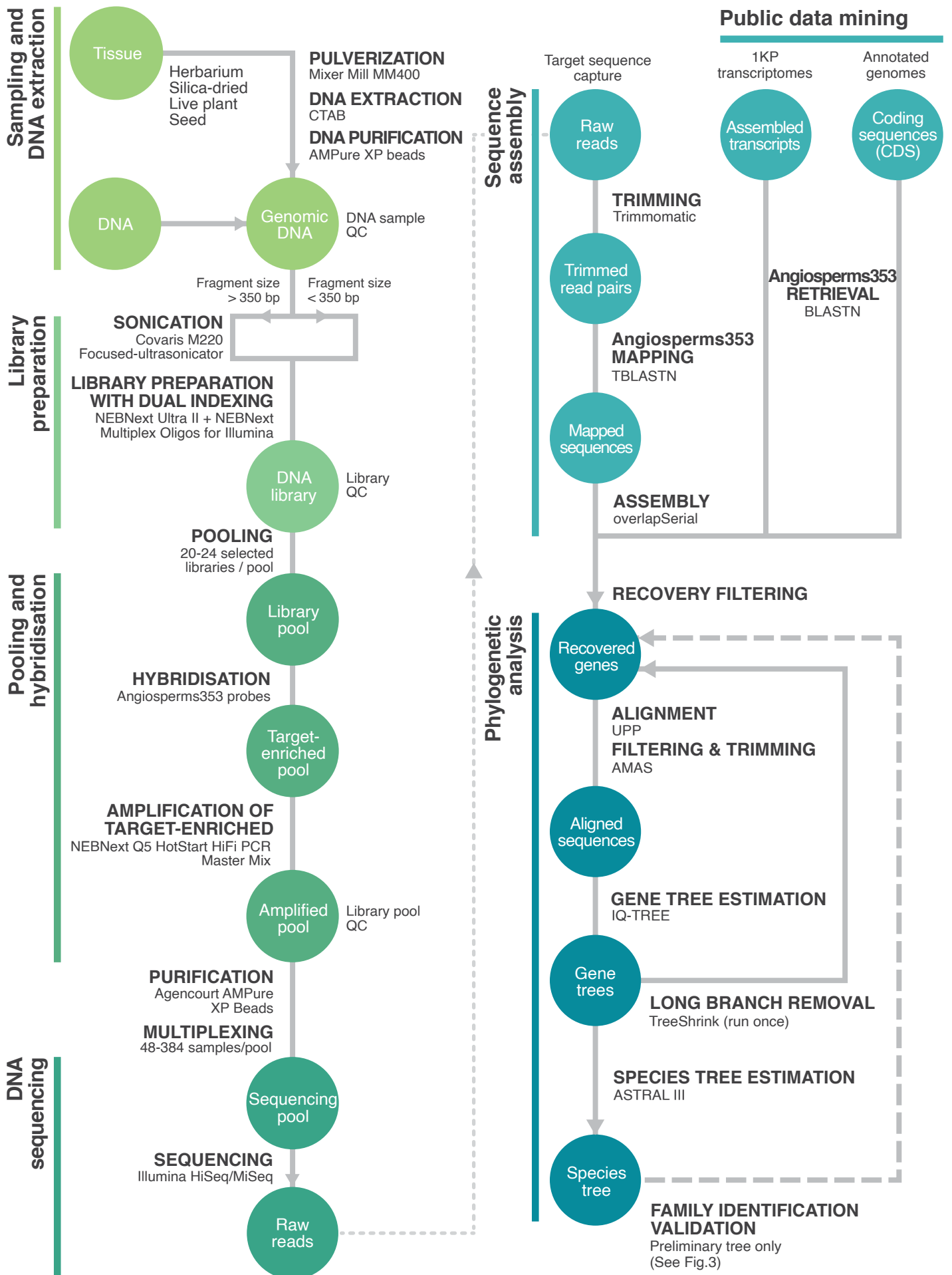
1241

1242 **Figure 5.** Density plots of target sequence recovery from our raw data. Data are presented
1243 prior to any filtering, illustrating relationships of sum of gene lengths (bp) to (a) the number
1244 of mapped reads and (b) the number of recovered genes. Colours indicate density of data

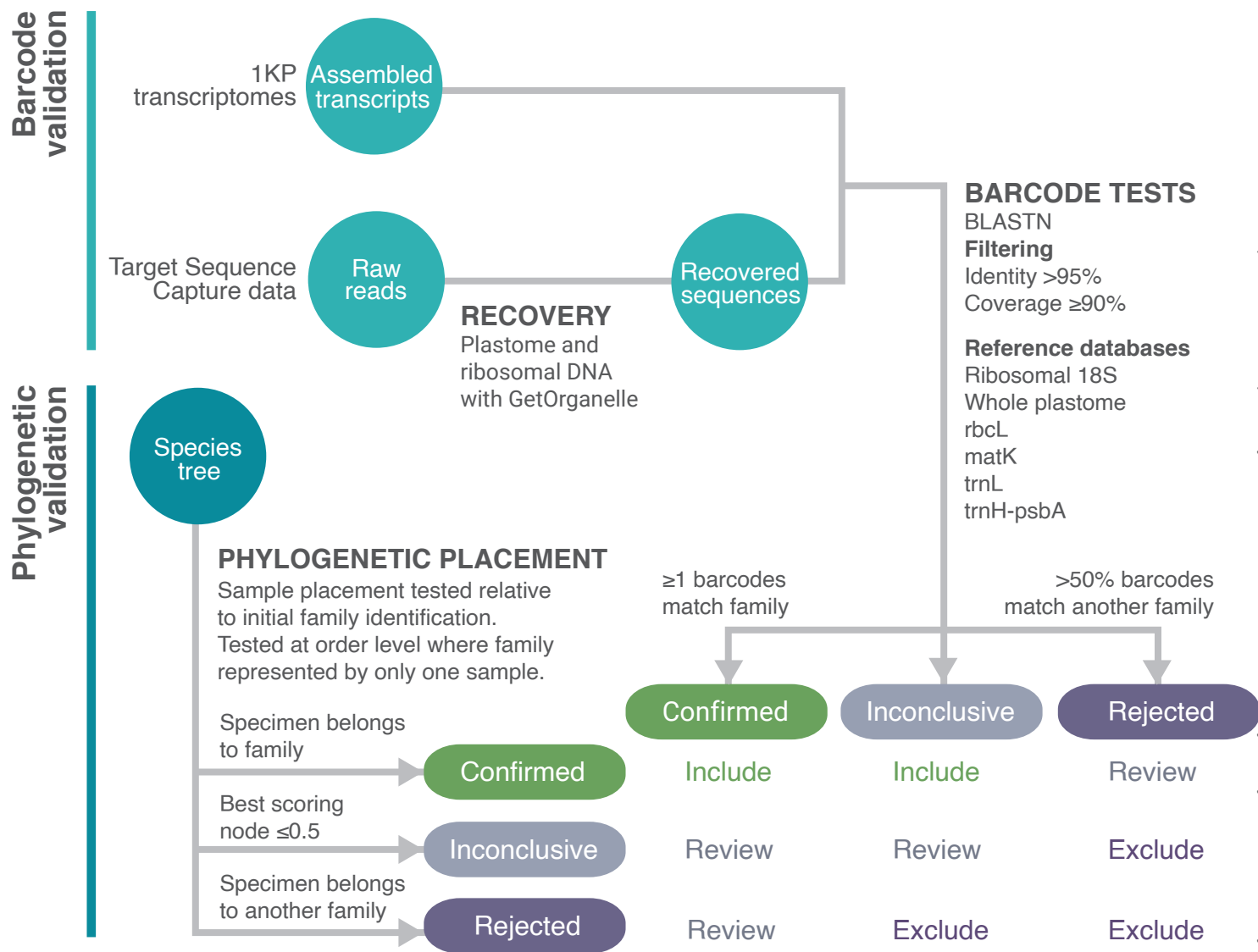
A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

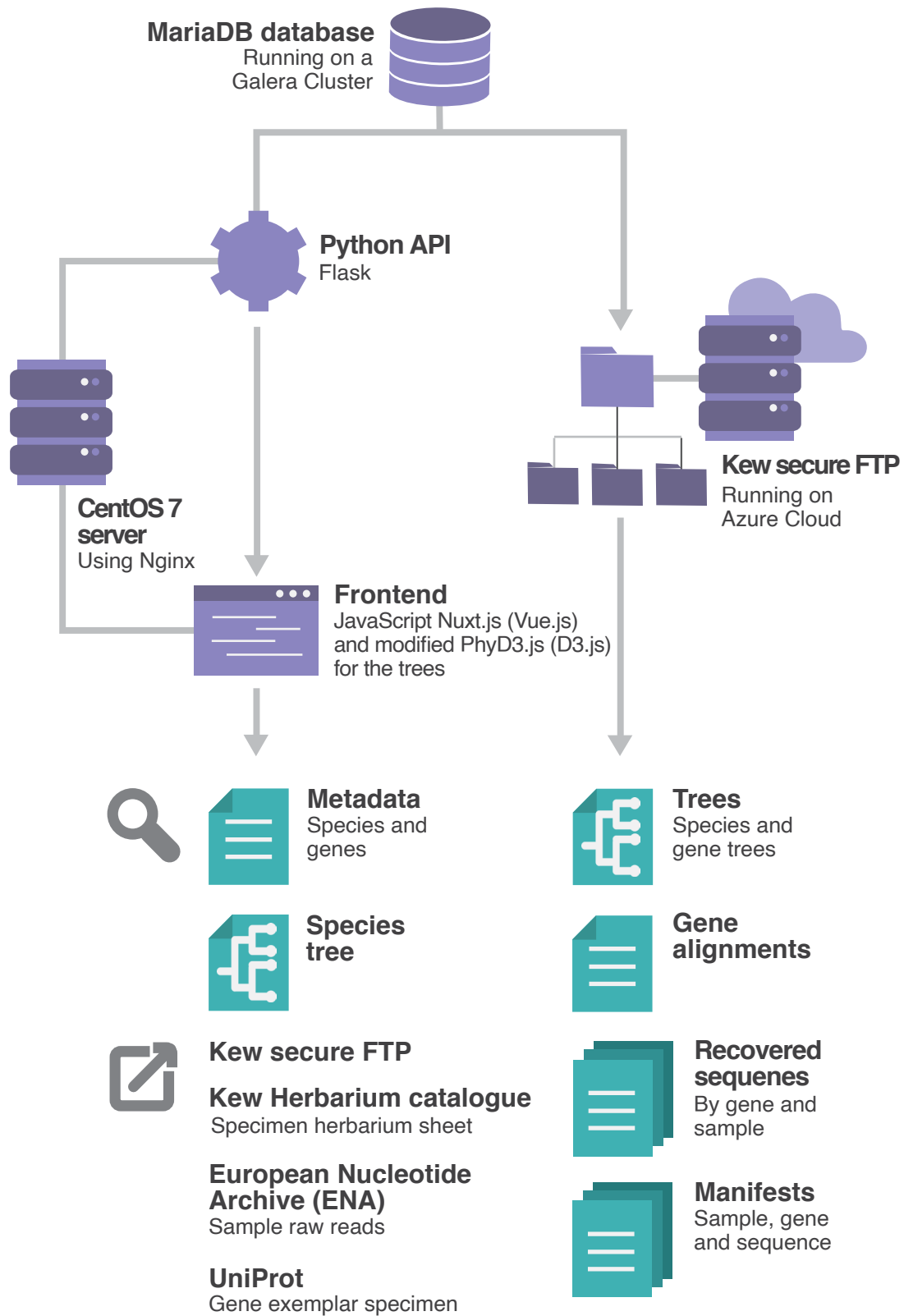
1245 points. Black dotted lines indicate medians of variables and red dotted lines indicate the
1246 threshold used to remove samples from downstream analyses, set as 20% of the median value
1247 across all samples.
1248

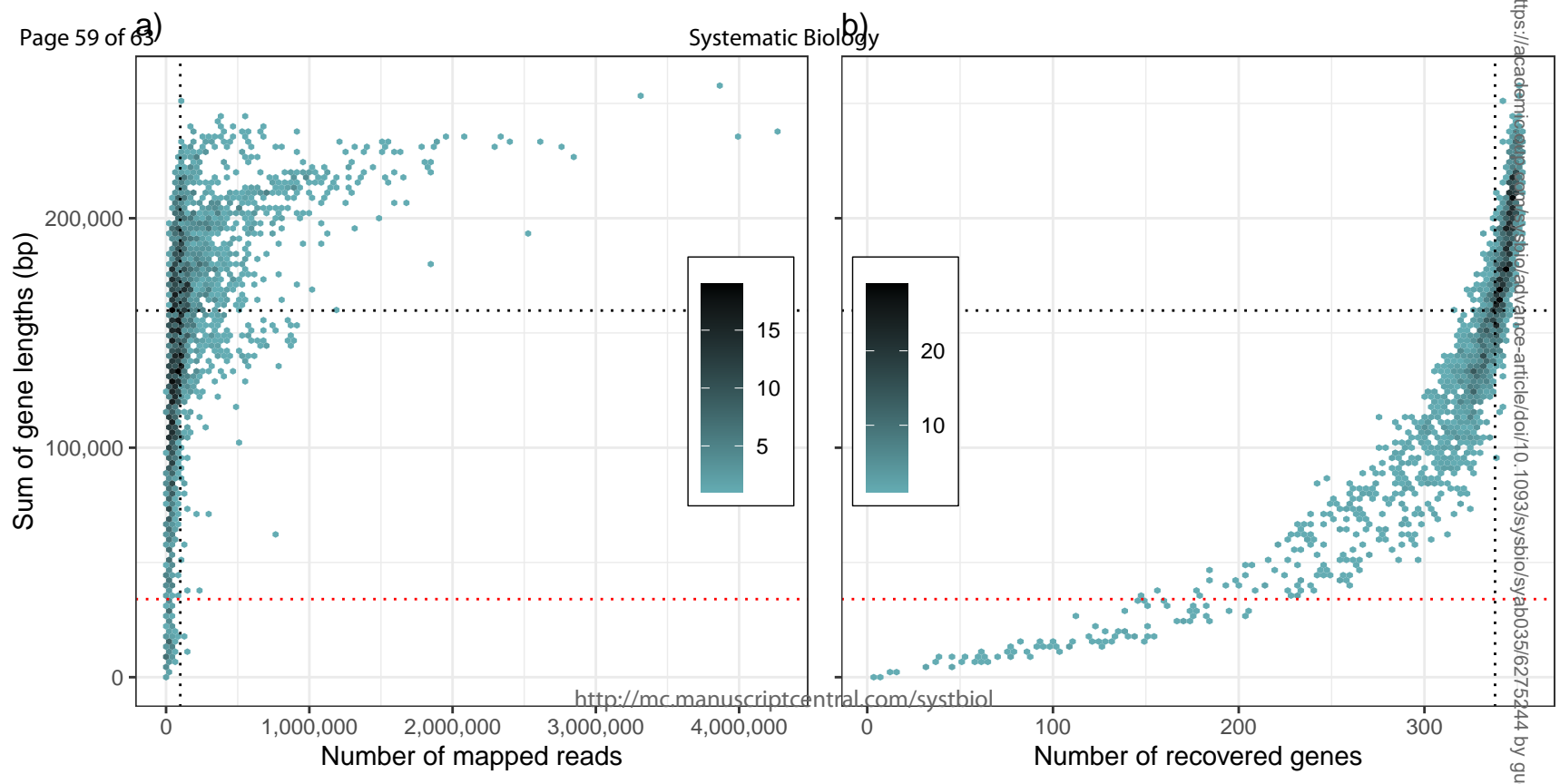




Downloaded from https://academic.oup.com/sysbio/advance-article/doi/10.1093/sysbio/syab035/6275244 by guest on 21 May 2021







<http://mc.manuscriptcentral.com/systbiol>

TABLES

Table 1. Total number of angiosperm samples included at three stages of data release preparation. The first column represents all samples available in the initial dataset. The second column indicates samples included in our preliminary tree, prior to family identification validation, but after removal of samples for which the sum of the gene lengths fell below 20% of the median value across all samples. The third column provides numbers for the samples made public in the Kew Tree of Life Explorer, Data Release 1.0, and included in our final phylogenetic tree. Numbers of angiosperm families, genera and species in each data subset are provided in brackets (as families/genera/species).

Data source	Initial dataset	Preliminary tree pre-validation	Final tree and Data Release 1.0
Target sequence capture data	2,522 (304/1988/2397)	2,438 (297/1947/2340)	2,374 (292/1903/2280)
1KP transcriptomes	689 (254/544/682)	678 (250/530/677)	664 (245/517/663)
Annotated genomes	61 (23/43/59)	61 (23/43/59)	61 (23/43/59)
Total	3,272 (413/2428/3079)	3,177 (410/2388/3028)	3,099 (404/2333/2956)

Baker et al.

Table 2. Target sequence capture and gene recovery statistics by sample or gene for Data Release 1.0, including the results of mining of genes from the 1KP and annotated genome datasets. The upper five rows apply to target sequence capture data only. (SD = standard deviation).

	Median	Mean	SD	Minimum	Maximum
Raw reads per sample	1.757 × 10 ⁶	2.822 × 10 ⁶	3.076 × 10 ⁶	1.676 × 10 ⁴	4.054 × 10 ⁷
Trimmed reads per sample	1.585 × 10 ⁶	2.549 × 10 ⁶	2.791 × 10 ⁶	1.391 × 10 ⁴	3.605 × 10 ⁷
Percentage of reads on-target per sample ^a	5.676	8.020	7.704	0.005	50.953
Read depth per sample ^b	38	90	105	5	2,243
Read depth per gene ^c	38	97	37	27	226
Recovered genes per sample:					
Target sequence capture data	338	330	24	148	353
1KP transcriptomes	341	328	44	30	353
Annotated genomes	346	341	13	287	353
Recovered genes lengths across all samples ^d (bp):					
Target sequence capture data	387	477	347	48	3,564
1KP transcriptomes	717	803	466	50	4,689

A PHYLOGENOMIC PLATFORM FOR ANGIOSPERMS

Annotated genomes	972	1,136	642	45	8,601
Sum of recovered gene lengths per sample (bp):					
Target sequence capture data	1.613 × 10 ⁵	1.576 × 10 ⁵	4.355 × 10 ⁵	3.433 × 10 ⁴	2.569 × 10 ⁵
1KP transcriptomes	2.753 × 10 ⁵	2.627 × 10 ⁵	6.659 × 10 ⁵	6.498 × 10 ⁵	3.674 × 10 ⁵
Annotated genomes	3.901 × 10 ⁵	3.876 × 10 ⁵	1.868 × 10 ⁴	3.217 × 10 ⁵	4.273 × 10 ⁵
Percentage length per recovered gene ^e across all samples:					
Target sequence capture data	63	62	16	27	96
1KP transcriptomes	88	85	10	44	100
Percentage length of recovered genes ^e per sample:					
Target sequence capture data	63	62	14	20	95
1KP transcriptomes	88	84	13	16	100

^aacross all recovered genes

^bat bases with $\geq 4x$ depth across all recovered genes, calculated by Samtools depth program

^cat bases with $\geq 4x$ depth across all samples, calculated by Samtools depth program

^dsee Supplementary Figure S7

Baker et al.

°percentage length calculated against each representative target gene

Baker et al.

Table 3. Properties of the 347 gene alignments and gene trees underpinning the species tree included in the Kew Tree of Life Explorer Data Release 1.0. (SD = standard deviation).

	Median	Mean	SD	Minimum	Maximum
Number of samples	2,421	2,377.2	358.8	491	3,014
% of total samples ^a	77.9	76.5	11.5	15.8	96.9
Alignment length	1,259.0	1,533.9	985.7	250	8,119
% gaps ^b	58.9	57.9	11.3	14.4	85.8
Variable sites	1,224	1,469.7	940.6	240	7,873
% variable sites	96.6	96.0	2.5	81.5	100
Parsimony informative sites	1,137	1,369.4	859.3	233	6,792
% parsimony informative sites	90.7	90.0	4.20	69.1	98.9
% nodes in gene trees above 30% UFBS ^c	98.9	98.5	1.3	90.7	99.9
Mean support ^c of all nodes	88.1	87.8	2.7	78.9	94.3
Median support ^c of all nodes	98.0	97.6	1.8	90.0	100

^apercentage of samples in species tree present in alignment/gene tree^bpercentage of empty cells in each alignment^cUFBS: ultrafast bootstrap