

1 **Forum Article**

2 **Hyb-Seq for flowering plant systematics**

3

4 Steven Dodsworth<sup>1,2†\*</sup>, Lisa Pokorny<sup>1†</sup>, Matthew G. Johnson<sup>3,4†</sup>, Jan T. Kim<sup>1</sup>, Olivier Maurin<sup>1</sup>, Norman  
5 J. Wickett<sup>4,5</sup>, Felix Forest<sup>1</sup>, William J. Baker<sup>1</sup>

6

7 <sup>1</sup>Royal Botanic Gardens, Kew, Richmond TW9 3AE, Surrey, UK.

8 <sup>2</sup>School of Life Sciences, University of Bedfordshire, University Square, Luton LU1 3JU, UK.

9 <sup>3</sup>Department of Biological Sciences, Texas Tech University, Lubbock, TX 79409, USA

10 <sup>4</sup>Chicago Botanic Garden, Glencoe, IL 60022, USA

11 <sup>5</sup>Program in Plant Biology and Conservation, Northwestern University, Evanston, IL 60208, USA

12

13 <sup>†</sup>Authors contributed equally

14 <sup>\*</sup>Correspondence: [steven.dodsworth@beds.ac.uk](mailto:steven.dodsworth@beds.ac.uk)

15

16 **Keywords**

17 High-throughput sequencing – molecular systematics – phylogenetics – Hyb-Seq – sequence  
18 capture – angiosperms – tree of life – genomics

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50 **Abstract**

51 High-throughput DNA sequencing (HTS) presents great opportunities for plant systematics,  
52 yet genomic complexity needs to be reduced for HTS to be effectively applied. We highlight  
53 Hyb-Seq as a promising approach, especially in light of the recent development of probes  
54 enriching 353 low-copy nuclear genes from any flowering plant taxon.

55

56

57 **High-throughput sequencing approaches and plant systematics**

58 Current developments in DNA sequencing, collectively termed high-throughput sequencing  
59 (HTS) technologies, permit many orders of magnitude more DNA data to be routinely  
60 collected compared to standard Sanger sequencing. This has made whole genome  
61 sequencing of diverse plant taxa much more accessible, including both flowering and non-  
62 flowering land plant lineages. However, challenges prevail: plant genome size varies  
63 enormously [1], genome assembly is often non-trivial for even the smallest plant genomes,  
64 and the cost per high-quality genome sequence is still significant. This means that, at least  
65 for the time being, methods are needed to reduce genomic complexity. This is especially the  
66 case for phylogenetics and systematics, in order to find an optimal amount of sequencing  
67 effort per sample whilst reaping the benefits of increased data. In this article, we propose  
68 Hyb-Seq as one of the most promising approaches for plant systematists currently, and  
69 particularly in light of a recent set of probes that target low-copy regions of the nuclear  
70 genome across flowering plants (angiosperms).

71

72 Systematics is primarily concerned with evolutionary relationships and natural classification,  
73 and as such producing reliable phylogenetic frameworks is often of primary concern. This is  
74 not the same as genomic studies, where detailed dissection of phenotypic traits or  
75 speciation processes may be the main goal—though there is a strong overlap between these  
76 fields. Phylogenetic data requires a constant trade-off between the depth (characters as  
77 DNA base pairs) and breadth (number of taxa) of data collected. Different evolutionary  
78 questions may demand different compromises on the depth-breadth spectrum. This is also  
79 a tension between an idealised data source (a complete nuclear genome sequence) and one  
80 that is easier and quicker to produce but far less information-rich (a small DNA barcode of a  
81 few hundred base pairs). Such examples lie at either end of a continuum of DNA sequencing  
82 tactics, making it difficult to find an optimal approach (Table 1).

83

84 Herbarium specimens are the foundation of taxonomic studies in plants. Herbarium DNA is  
85 usually highly fragmented and often contaminated, making PCR-based approaches  
86 challenging [2,3]. HTS can surmount these difficulties as all native DNA fragments present  
87 can potentially be sequenced [3,4], although different approaches have their advantages  
88 and disadvantages (see below).

89

90 *Genome Skimming*

91 Simple approaches such as genome skimming [4] remain popular, although recovery of  
92 orthologous nuclear regions for sequence alignment is limited with these techniques. Whilst  
93 organellar genomes (particularly plastid genomes) are easily reconstructed from such data,  
94 their histories reflect patterns associated with matrilineal genealogy/geography or other  
95 aspects of organelle biology. As such phylogenetic inference based on plastid or organellar  
96 data may not necessarily reflect the evolutionary history of the taxa in question (for a

97 comprehensive view of plastid evolution, see [5]). Ribosomal DNA is easily recovered,  
98 although not always highly variable and concerted evolution can produce incongruent  
99 topologies. Other repetitive elements (e.g. satellite DNA, transposable elements) can be  
100 easily quantified from a genome skim, but sequence divergence of such repeats is low.  
101 Repeat abundance and repeat sequence similarity can be used instead of sequence  
102 alignment for phylogenetic reconstruction [6] although these are very different approaches,  
103 both conceptually and practically.

#### 104 105 *RAD-Seq*

106 Restriction site-associated DNA sequencing (RAD-Seq or similar Genotyping-by-Sequencing  
107 approaches; GBS) is a method to sequence DNA next to restriction sites. The loci are  
108 essentially random, although partial selection for particular genomic contexts (e.g. genic  
109 regions) is possible using methylation-sensitive enzymes [7]. RAD-Seq holds particular  
110 promise at shallow scales, for resolving recent radiations and population-level sampling [8],  
111 where a large number of single nucleotide polymorphisms (SNPs) help. RAD-Seq loci are  
112 often short, however, and not always easy to annotate without a high-quality reference  
113 genome. As genomic DNA is cut with enzymes, high molecular weight DNA is required.  
114 Recent silica-dried collections therefore work well as do very recent herbarium specimens  
115 but degraded DNA from older herbarium specimens will not work. Due to the variability of  
116 restriction sites between taxa, particularly over larger evolutionary distances, securing  
117 enough homologous loci is difficult at deeper (or variable) phylogenetic scales. This also  
118 means that RAD-Seq data in public repositories may not be a very usable resource (e.g. as a  
119 source of outgroup sequences from related taxa).

#### 120 121 *RNA-Seq*

122 Transcriptomics requires high-quality RNA from samples, which usually means flash-frozen  
123 using liquid nitrogen or dry ice or using pricey preservative liquids designed to preserve RNA  
124 in the field and requiring -80 °C storage. Resulting data will include all expressed genes in  
125 that particular sample, which makes RNA-Seq ideal for obtaining large numbers of protein-  
126 coding genes. Due to differences in expression throughout the plant, though, a variety of  
127 tissues should ideally be used (e.g. flower, root, leaf). There are some obvious caveats to  
128 this approach: (i) it requires healthy living plant tissue and access to preservatives/freezers;  
129 and (ii) it may require a variety of tissues; and (iii) it remains relatively expensive per sample  
130 (Table 1).

### 131 132 133 **Sequence capture, target enrichment and Hyb-Seq approaches**

#### 134 *Bait design*

135 Sequence capture approaches are becoming increasingly popular as a method of reducing  
136 genomic complexity, exploiting “baits” (probes) to enrich specific target regions (loci) from  
137 total DNA. This approach has been variously referred to as bait hybridisation, target  
138 enrichment, sequence/target/hybrid capture, Hyb-Seq, or other combinations of such  
139 terms. A common feature is the use of pre-designed RNA or DNA bait sequences, developed  
140 from pre-existing genomic information, such as a closely-related genome sequence or  
141 transcriptome data. Target loci are often nuclear protein-coding sequences or other  
142 conserved genomic regions, such as ultra-conserved elements (UCEs—in animals and fungi).  
143 Typically, low-copy (ideally single-copy) genes are chosen for phylogenetic purposes, thus

144 minimising any orthology issues later on. In many cases, however, multigene families are  
145 also included [e.g. 9], particularly where those genes have known functions of biological  
146 interest to the groups being studied (e.g. photosynthetic transitions, or transcription factors  
147 involved in morphological diversity).

148

149 If protein-coding regions are targeted, phylogenetic inference can employ explicit models  
150 that account for different rates of evolution based on codon position. Such explicit  
151 positional information is often required for reliable inference at deeper phylogenetic scales  
152 [10]. Codon positions are often difficult to infer using RAD-Seq data, protein-coding nuclear  
153 data are lacking in genome skims, and RNA-Seq is expensive. Hyb-Seq can provide protein-  
154 coding data at a fraction of the cost, and a compromise point where these other approaches  
155 fall down.

156

### 157 *Generalised workflow*

158 Genomic DNA extracts are first turned into libraries of genomic fragments. The RNA/DNA  
159 baits are subsequently hybridized to target loci in genomic libraries. Bait-bound DNA is then  
160 separated from the mixture, e.g. by using streptavidin-coated magnetic beads that bind  
161 biotinylated baits (and bait-bound DNA), that can then be separated simply with a magnet  
162 (Figure 1). DNA fragments not bound to baits are discarded through a series of washing  
163 steps, and the result is a pool of fragments enriched for particular target sequences (Figure  
164 1).

165

166 Effective recovery of target loci can be achieved even with surprisingly low levels of  
167 enrichment, as low as 10% of the sequence reads [9]. Consequently, there can be abundant  
168 off-target reads that include high-copy DNA regions, such as repetitive DNA, the ribosomal  
169 operon, and organellar DNA from plastids and mitochondria (Figure 1). This off-target  
170 fraction is similar to a genome skim [4], or low-coverage whole-genome sequencing, and  
171 can also be exploited for systematic analyses [11]. Moreover, regions adjacent to the target  
172 loci (known as the “splash zone”) are also recovered (Figure 1), often including intronic  
173 regions, which may be highly variable and therefore valuable at shallower phylogenetic  
174 levels [12,13].

175

### 176 *Hyb-Seq*

177 The term Hyb-Seq was initially proposed by Weitemier et al. (2014; [12]) to consider the  
178 explicit use of both the on-target reads (i.e. enriched gene sequences) and the off-target  
179 fraction. In recent years, the term Hyb-Seq has had slightly different meanings, such as  
180 mixing the enriched and unenriched (native) libraries [11], or explicitly sequencing both  
181 enriched and unenriched sets of libraries separately. The fundamental meaning remains the  
182 same—utilisation of both low-copy enriched nuclear sequences and high-copy unenriched  
183 ones such as plastid and ribosomal DNA.

184

185 The unenriched category notably and conveniently includes markers that have been  
186 traditionally used for decades in plant systematics, the currently used plant DNA barcodes—  
187 *rbcl*, *matK*, *trnH-psbA* spacer (plastid genome) and nrITS of ribosomal DNA. Sequencing  
188 these loci will facilitate the ongoing global synthesis of plant systematic data for a variety of  
189 use cases. Hyb-Seq has been successfully used in a number of groups at varying levels of  
190 phylogenetic depth [e.g. 11,12]; it has also been used very effectively with herbarium

191 samples, including those over 100 years old and spanning the diversity of angiosperms  
192 [11,14].

193

## 194 **Enriching a core set of genes in flowering plants and future potential**

### 195 *Angiosperms-353 bait set*

196 Probes for sequence capture have traditionally been designed for specific plant groups of  
197 interest. The design of such a kit requires access to (or production of) genomic resources  
198 and at least some bioinformatic expertise. Recent publication of an angiosperm-wide set of  
199 baits makes Hyb-Seq a great deal more accessible for flowering plants and alleviates part of  
200 the financial and bioinformatic burden [4]. Johnson et al. (2018; [15]) have developed a  
201 probe set that targets 353 low-copy orthologous nuclear genes in angiosperms, derived  
202 from an alignment of low-copy genes across all green plants by the 1000 Plant  
203 Transcriptomes Initiative or OneKP project (onekp.com). Their approach includes the use of  
204 up to 15 variants for each of the 353 gene loci (approx. 230 Kbp of nuclear sequence), in  
205 order to capture sequence diversity across angiosperms with one single kit (Angiosperms-  
206 353, available at [arborbiosci.com/products/mybaits-plant-angiosperms](http://arborbiosci.com/products/mybaits-plant-angiosperms), catalog #3081XX).  
207 Including variants means that, on average, DNA from 95% of angiosperm species should  
208 hybridise to one or more gene variants with  $\leq 30\%$  divergence between the sample and the  
209 target sequence. Importantly, hybridisation is reported to be efficient below such a  
210 threshold.

211

### 212 *Future potential*

213 This means that this kit should work for any of the 300,000 currently estimated angiosperm  
214 species, distributed in 416 families, and which dominate terrestrial ecosystems globally.  
215 Johnson et al. [15] show very promising data for 42 samples taken from across the  
216 angiosperms, with no obvious systematic/taxonomic biases, and potential phylogenetic  
217 signal at various levels.

218

219 The Angiosperms-353 kit has enormous potential for studies that combine deep and  
220 shallow-level systematic studies. There is also promise as a powerful new tool in the fields  
221 of molecular and community ecology (e.g. discovering the types of pollen carried by  
222 pollinators, community assembly, or characterising habitats through molecular sampling).  
223 This is potentially possible by building a database of a common set of hundreds of genes per  
224 sample. Such a set of core genes may even be a nuclear solution for the “next generation”  
225 flowering-plant DNA barcode.

226

227

## 228 **References**

229 1. Pellicer, J. *et al.* (2018) Genome size diversity and its impact on the evolution of land  
230 plants. *Genes* 9, 88.

231

232 2. Särkinen, T. *et al.* (2012) How to open the treasure chest? Optimising DNA extraction  
233 from herbarium specimens. *PLoS One*, e43808.

234

235 3. Bakker, F.T. (2017) Herbarium genomics: skimming and plastomics from archival  
236 specimens. *Webbia* 72, 35-45.

- 237  
238 4. Dodsworth, S. (2015) Genome skimming for next-generation biodiversity assessment.  
239 *Trends in Plant Science* 20, 525-527.  
240  
241 5. Gitzendanner, M.A. *et al.* (2018) Plastid phylogenomic analysis of green plants: A billion  
242 years of evolutionary history. *Am. J. Bot.*  
243  
244 6. Dodsworth, S. *et al.* (2015) Genomic repeat abundances contain phylogenetic signal. *Syst.*  
245 *Biol.* 64, 112-126.  
246  
247 7. Elshire, R.J. *et al.* (2011) A robust, simple Genotyping-by-Sequencing (GBS) approach for  
248 high diversity species. *PLoS One* 6, e19379.  
249  
250 8. Paun, O. *et al.* (2015) Processes driving the adaptive radiation of a tropical tree  
251 (*Diospyros*, Ebenaceae) in New Caledonia, a biodiversity hotspot. *Systematic Biology* 65,  
252 212-227.  
253  
254 9. Moore, A.J. *et al.* (2017) Targeted enrichment of large gene families for phylogenetic  
255 inference: Phylogeny and molecular evolution of photosynthesis genes in the Portullugo  
256 clade (Caryophyllales). *Syst. Biol.* 67, 367-383.  
257  
258 10. Wickett, N.J. *et al.* (2014) Phylotranscriptomic analysis of the origin and early  
259 diversification of land plants. *PNAS* 111, E4859-68.  
260  
261 11. Villaverde, T. *et al.* (2018) Bridging the micro and macroevolutionary levels in  
262 phylogenomics: Hyb-Seq solves relationships from populations to species and above. *New*  
263 *Phytologist* 220, 636-650.  
264  
265 12. Weitemier, K. *et al.* (2014) Hyb-Seq: Combining target enrichment and genome  
266 skimming for plant phylogenomics. *Appl. Plant Sci.* 2  
267  
268 13. Johnson, M.G. *et al.* (2016) HybPiper: Extracting coding sequence and introns for  
269 phylogenetics from high-throughput sequencing reads using target enrichment. *Appl. Plant*  
270 *Sci.* 4, 1600016-1600018  
271  
272 14. Hart, M.L. *et al.* (2016) Retrieval of hundreds of nuclear loci from herbarium specimens.  
273 *Taxon* 65, 1081-1092.  
274  
275 15. Johnson, M.G. *et al.* (2018) A universal probe set for sequence capture of 353 nuclear  
276 genes from any flowering plant designed using k-medoids clustering. *Systematic Biology* 68,  
277 594-606.  
278  
279

## 280 **Acknowledgements**

281 This research was supported by grants from the Calleva Foundation, the Sackler Trust and  
282 the Garfield Weston Foundation to the Royal Botanic Gardens, Kew.  
283



285 **Table 1.** Comparison of high-throughput sequencing approaches for plant systematics:  
 286 advantages and disadvantages<sup>a</sup>  
 287

Phylogenomics approach	Genomic resources required	Initial bioinformatic investment	Ultimate bioinformatic investment	Initial laboratory cost	Ultimate cost per sample	Low-copy nuclear genes retrieved
<i>Genome skimming</i>	No	None	Medium	Low	Medium	No/Limited
<i>RAD-Seq</i>	No, but helpful	Medium	High	High	Low	No/SNPs
<i>RNA-Seq</i>	No, but helpful	Low	High	Low	High	Yes-thousands
<i>Hyb-Seq</i>	Varies <sup>b</sup>	High <sup>b</sup>	Medium	Low <sup>b</sup>	Medium	Yes-variable

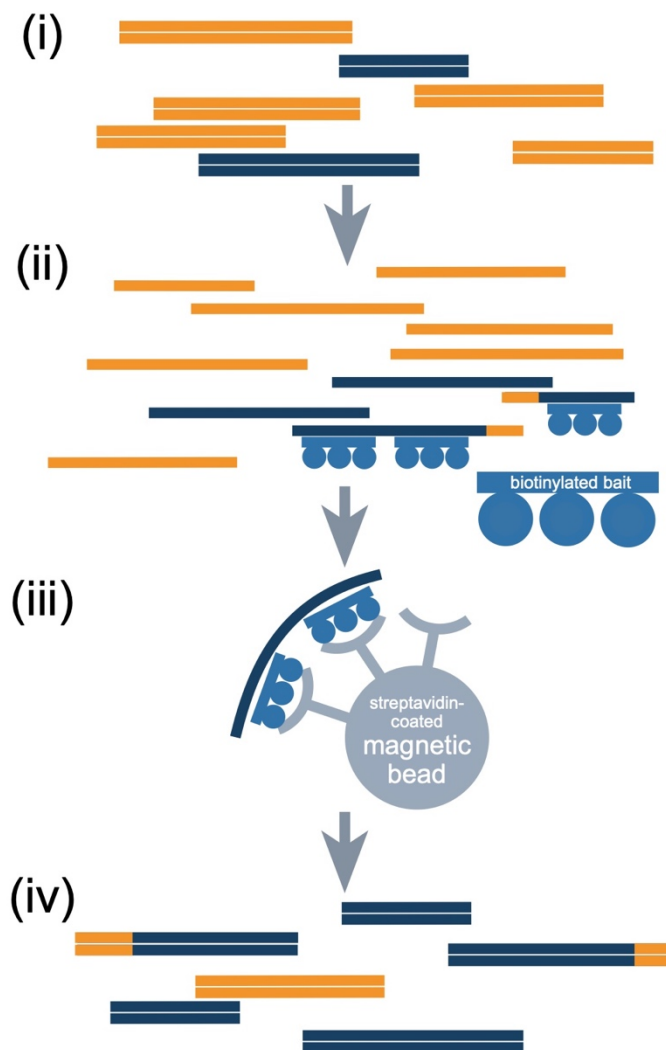
288 <sup>a</sup>Initial costs include the one-time or limited purchase of expensive consumables (e.g.  
 289 biotinylated baits or adapter sequences). Boxes are highlighted from unfavourable (red) to  
 290 favourable (green) under each column.

291 <sup>b</sup>If designing new kit(s) genome or transcriptome resources are required, otherwise readily available kits exist  
 292 for different groups of plants as well as angiosperms as a whole (Angiosperms-353) and are much cheaper  
 293 than designing a new custom bait set.

294  
 295



296 **Figure 1.** Simplified schematic representing the main steps in a typical Hyb-Seq workflow: (i)  
 297 Libraries of double-stranded DNA fragments are prepared from genomic DNA; (ii) Libraries  
 298 are denatured (single-stranded) and bound to biotinylated probes/baits; (iii) streptavidin-  
 299 coated magnetic beads bind to the biotinylated bait-DNA hybrids, these are bound to a  
 300 magnet, and other DNA fragments are washed off; (iv) baited-DNA is PCR-ed and removed  
 301 from the beads for sequencing. Target DNA sequences are in dark blue and non-target  
 302 sequences are in orange. Hyb-Seq has the potential to recover both “splash zone” sequences  
 303 close to targets (edges of dark blue sequences in orange, e.g. introns) as well as some  
 304 completely off-target sequences (orange blocks, e.g. plastid DNA), as indicated in the final  
 305 sequencing library (iv).  
 306



307