

1 Adapterama IV: Sequence Capture of Dual-digest RADseq Libraries with Identifiable Duplicates
2 (RADcap)
3

4 Sandra L. Hoffberg^{1,*}, Troy J. Kieran², Julian M. Catchen³, Alison Devault⁴, Brant C. Faircloth⁵,
5 Rodney Mauricio¹, Travis C. Glenn^{1,2,*}
6

7 ¹Department of Genetics, University of Georgia, Athens, GA 30602, USA

8 ²Department of Environmental Health Science, University of Georgia, Athens, GA 30602, USA

9 ³Department of Animal Biology, University of Illinois, Urbana, IL 61801, USA

10 ⁴MycroArray, 5692 Plymouth Rd., Ann Arbor, MI 48105, USA

11 ⁵Department of Biological Sciences and Museum of Natural Science, Louisiana State University,
12 Baton Rouge, LA 70803, USA
13

14 **Keywords:** ddRAD, Illumina, reduced representation libraries, target enrichment, *Wisteria*,
15 3RAD

16 ***Correspondence:**

17 Sandra Hoffberg Department of Genetics, University of Georgia, Athens, GA 30602, USA; Tel:
18 706-542-1417; Fax: 706-542-3910; E-mail: sandra@hoffberg.org

19 Travis C. Glenn, Dept. of EHS, Environmental Health Science Bldg., University of Georgia,
20 Athens, GA 30602, USA; Tel: 706-583-0662; fax: 706-542-7472; E-mail: travisg@uga.edu
21

22 **Running Title:** RADcap: Sequence Capture of RADseq Libraries

23

24 **Abstract:**

25 Molecular ecologists seek to genotype hundreds to thousands of loci from hundreds to thousands
26 of individuals at minimal cost per sample. Current methods such as restriction site associated
27 DNA sequencing (RADseq) and sequence capture are constrained by costs associated with
28 inefficient use of sequencing data and sample preparation, respectively. Here, we demonstrate
29 RADcap, an approach that combines the major benefits of RADseq (low cost with specific start
30 positions) with those of sequence capture (repeatable sequencing of specific loci) to significantly
31 increase efficiency and reduce costs relative to current approaches. The RADcap approach uses
32 a new version of dual-digest RADseq (3RAD) to identify candidate SNP loci for capture bait
33 design, and subsequently uses custom sequence capture baits to consistently enrich candidate
34 SNP loci across many individuals. We combined this approach with a new library preparation
35 method for identifying and removing PCR duplicates from 3RAD libraries, which allows
36 researchers to process RADseq data using traditional pipelines, and we tested the RADcap
37 method by genotyping sets of 96 to 384 *Wisteria* plants. Our results demonstrate that our
38 RADcap method: 1) can methodologically reduce (to <5%) and computationally remove PCR
39 duplicate reads from data; (2) achieves 80-90% reads-on-target in 11 of 12 enrichments; (3)
40 returns consistent coverage ($\geq 4x$) across >90% of individuals at up to 99.9% of the targeted loci;
41 (4) produces consistently high occupancy matrices of genotypes across hundreds of individuals;
42 and (5) is inexpensive, with reagent and sequencing costs totaling <\$6/sample and adapter and
43 primer costs of only a few hundred dollars.

44

45 **Introduction**

46 Massively parallel sequencing is changing molecular ecology and other life science
47 disciplines (Rogers & Venter 2005; Tautz *et al.* 2010). While the costs of whole genome
48 sequencing and genome resequencing have declined, the time investment, cost, and
49 computational complexity of genome assembly and genome resequencing remain significant
50 drawbacks. Fortunately, many biological hypotheses can be tested with smaller samples of the
51 genome that collect data from several hundred to several thousand variable loci (Cariou *et al.*
52 2013; Pante *et al.* 2015) rather than requiring the millions of variable sites identified during
53 genome sequencing/re-sequencing workflows. Although genome reduction techniques that
54 collect data from hundreds or thousands of loci are an appealing and inexpensive proxy for
55 whole genome resequencing, the matter of how best to collect genotypes from hundreds or
56 thousands of loci across hundreds or thousands of individuals remains. Thus, researchers still
57 face the decades-old dilemma of choosing among methods that make trade-offs in the number
58 and kind of data (loci) collected versus the number of individuals surveyed.

59 Genome reduction techniques fall into a broad class of so-called “reduced representation”
60 approaches, and these methods are meant to collect data from a small and repeatable fraction of
61 the genome across a population of individuals - enabling the population under study to be
62 compared at identical loci without sequencing the entire genome (Altshuler *et al.* 2000; Novaes
63 *et al.* 2008; Wiedmann *et al.* 2008). Two general types of reduced representation library
64 approaches for massively parallel sequencing are widely used – sequence capture (Gnirke *et al.*
65 2009; Okou *et al.* 2007) and restriction-site associated DNA sequencing (RADseq; Baird *et al.*
66 2008; Davey & Blaxter 2010; Davey *et al.* 2011; Miller *et al.* 2007; Peterson *et al.* 2012).

67 Although both methods have advantages and disadvantages (Harvey *et al.* 2013), neither is
68 entirely capable of achieving a primary goal of many population genetic studies: consistently
69 obtaining a set of hundreds or thousands of putatively unlinked single nucleotide polymorphisms
70 (SNPs) from hundreds to thousands of individuals at low cost (e.g., <\$10/sample).

71 Sequence capture is a powerful technique that combines a custom set of long,
72 biotinylated, oligonucleotide baits with in-solution hybridization to target and enrich any number
73 of genomic regions of varying size (Gnirke *et al.* 2009), from entire chromosomal segments (Cao
74 *et al.* 2013; Pröll *et al.* 2011; Rabenstein *et al.* 2015) to sets of smaller loci containing SNPs
75 (Kenny *et al.* 2010; Saintenac *et al.* 2011). Sequence capture requires some form of prior
76 sequence information to design capture baits (Gnirke *et al.* 2009), which can be a challenge for
77 certain non-model species and marker types. Several groups have designed bait sets that target
78 conserved sequences including ultraconserved elements (UCEs; Faircloth *et al.* 2012), anchor
79 regions (Lemmon *et al.* 2012), and exons (Bi *et al.* 2012), which allow sets of baits to be used
80 across many species (e.g., at the level of taxonomic class; Li *et al.* 2013). Although useful for
81 many situations, sequence capture is constrained by high library preparation costs, expensive
82 baits, randomness of where the collected sequences start and stop, and off-target sequence reads
83 (Harvey *et al.* 2013). Although off-target reads can be beneficial (Raposo do Amaral *et al.* 2015),
84 researchers must account for such reads to achieve the desired sequencing coverage (e.g., if 50%
85 of reads are off-target, researchers must obtain twice as many reads to yield the same coverage as
86 when all reads are on-target). Conversely, the randomness inherent in the beginning and ending
87 positions of the collected sequences enables the removal of PCR duplicates and probabilistic
88 variant calling methods, which are widely used in genome resequencing research.

89 RADseq methods reduce the genome by sequencing thousands to hundreds of thousands
90 of DNA fragments that are located near restriction enzyme cut sites (Baird *et al.* 2008; Davey *et al.*
91 *al.* 2011; Miller *et al.* 2007). Many RADseq derivatives have been developed (Andrews *et al.*
92 2016), and throughout this manuscript, we will use the term “traditional RADseq” to mean
93 methods where one end of the sequenced DNA insert derives from a restriction-site and the other
94 end is randomly sheared (Baird *et al.* 2008; Davey *et al.* 2011; Miller *et al.* 2007), whereas we
95 will use “RADseq” to generically refer to any of the derivative forms of RADseq. Dual-digest
96 RADseq (ddRAD; Peterson *et al.* 2012) methods, including our 3RAD variant (Glenn *et al.*
97 2016b; Graham *et al.* 2015), sequence DNA inserts that fall precisely between two restriction
98 enzyme cut sites, giving both ends of the sequenced DNA precise start and stop positions.
99 Sequencing these ddRAD-type libraries is particularly efficient when compared to libraries
100 derived from sequence capture or traditional RADseq approaches because sequencing reads pile-
101 up on the ends of the loci, boosting coverage and efficiently increasing the accuracy of
102 downstream SNP calling (Fountain *et al.* 2016). Compared to sequence capture, RADseq
103 methods generally have lower library preparation costs and do not explicitly require existing
104 genomic information from the taxa of interest.

105 The primary disadvantages of RADseq (see Mastretta-Yanes *et al.* (2015) for review)
106 include sequencing many monomorphic loci and stochastic variation (mutation and methylation)
107 at the restriction enzyme cut sites that produce sparse genotype matrices. Overall quality and
108 utility of RADseq data sets can also be affected by abiotic factors such as the molecular and
109 bioinformatic protocols used to generate the RADseq data. For example, ddRAD loci may be
110 lost due to imprecise size selection methods; errors within loci may be introduced by low-quality

111 reagents; and/or PCR bias may preferentially amplify smaller fragments, GC rich regions (Puritz
112 *et al.* 2014), or one allele over another (Casbon *et al.* 2011).

113 Particularly problematic are errors introduced to RADseq libraries during PCR.
114 Incorporation errors early during the PCR reaction can be amplified to high coverage as PCR
115 proceeds (Tin *et al.* 2015), and PCR duplication of loci can give false confidence in the accuracy
116 of downstream variant calls. For example, many RADseq processing pipelines use coverage to
117 validate the accuracy of SNP calls and PCR duplicates can comprise 20-90% of reads in RADseq
118 libraries (Ali *et al.* 2015; Andrews *et al.* 2014; Schweyen *et al.* 2014; Tin *et al.* 2015); therefore,
119 duplicate reads can produce inflated coverage across loci, resulting in falsely high confidence
120 assigned to the genotypes obtained (Casbon *et al.* 2011; Schweyen *et al.* 2014; Tin *et al.* 2015).

121 The traditional approach for distinguishing duplicates in standard genomic libraries,
122 which are randomly sheared on both ends, and traditional RADseq libraries, which are randomly
123 sheared on one end, is to identify duplicate reads as those having identical start and stop
124 positions when aligned to a reference sequence. However, this technique cannot be applied to
125 ddRAD-type approaches, where all sequence reads from each RAD locus are identical (Andrews
126 *et al.* 2014), and it is thought that using less template DNA can exacerbate the problem of read
127 duplication (Casbon *et al.* 2011). Single molecule tagging has been employed to identify and
128 remove PCR duplicates in a variety of approaches (Hiatt *et al.* 2013; Jabara *et al.* 2011; Kivioja
129 *et al.* 2012; Miner *et al.* 2004; Shiroguchi *et al.* 2012; Smith *et al.* 2014), including deep
130 sequencing of limited input DNA. Recently, this approach has been employed in RADseq and
131 ddRAD experiments by incorporating degenerate bases in adapters (Casbon *et al.* 2011;
132 Schweyen *et al.* 2014; Tin *et al.* 2015), but all of the methods developed, to date, have some

133 limitations in their general implementation (see Discussion for details). Single molecule tagging
134 approaches that are easy to implement and have high power to distinguish PCR duplicates are
135 still needed.

136 Here, we introduce RADcap, a novel method that combines the benefits of single-
137 molecule tagging with 3RAD and sequence capture to collect a consistent and repeatable sample
138 of hundreds of loci across hundreds of individuals, remove PCR duplicates from the resulting
139 data, and call SNPs using a probabilistic base-calling pipeline (GATK; DePristo *et al.* 2011;
140 McKenna *et al.* 2010). The RADcap workflow begins with a pilot experiment using 3RAD to
141 collect genetic information from a small sample of individuals. After processing the resulting
142 sequence reads using Stacks (Catchen *et al.* 2013; Catchen *et al.* 2011) to identify variable RAD
143 loci, the workflow proceeds by designing a set of biotinylated ssRNA baits targeting a subset of
144 the variable RAD loci, and enriching the targeted loci from a pool of DNA libraries prepared
145 using our inexpensive 3RAD library preparation process. To ameliorate the problem of false
146 confidence in genotype calls bolstered by PCR duplicates, the RADcap approach incorporates a
147 random 8nt sequence tag in place of the iTru5 primer index (Figures 1 and 2) to each library
148 molecule, which allows researchers to distinguish PCR duplicates from unique template
149 molecules during post-processing of the sequence data. Finally, following a GATK workflow,
150 we created a RADcap data processing package, which calls SNPs in the duplicate-free reads
151 using a “radnome” (those RAD loci we targeted with capture baits) as a reference sequence. We
152 empirically tested the RADcap method by measuring genetic diversity of 96 samples of *Wisteria*
153 across an urban center.

154

155 **Methods**

156 Study Species and Experimental Design

157 Previous research on *Wisteria* in the southeastern United States (Trusty et al. 2007;
158 Trusty et al. 2008) and within Athens, GA (Glenn *et al.* 2016c) has shown that most *Wisteria*
159 plants are hybrids of *W. sinensis* and *W. floribunda*. Both species were introduced to the United
160 States in the early 1800s (Wilson 1916; Wyman 1949) as ornamental plants and both species
161 reproduce sexually and vegetatively (Valder 1995). While currently available genetic markers
162 can distinguish species, there are no markers available with enough resolution to distinguish
163 among individuals from the same population. Understanding the population genetics of this
164 invasive plant requires many more markers, and is crucial to understanding how it is spreading.
165 We use SNPs to estimate the genetic diversity of 96 samples across Athens, Georgia, USA. We
166 compare these estimates to estimates of genetic diversity obtained from SNPs from the same loci
167 in the same samples, that were prepared via the more traditional 3RAD instead of RADcap.

168

169 3RAD SNP Discovery for Bait Design

170 Because sequence capture uses baits designed from pre-existing sequence information,
171 we collected these data using a pilot 3RAD study of four individual *Wisteria* plants: three
172 samples collected around Athens, GA (wist69-3, wist124-1, and wist276-4) and one sample
173 collected from greenhouse-grown seedlings (Wmat9-7-P5-S1). We prepared samples using
174 3RAD (Glenn *et al.* 2016b; Graham *et al.* 2015), which we summarize below and explain, in
175 detail, in the Supplemental Methods. We added short forward and reverse adapters containing a
176 unique tag combination to extracted DNA from each of the four samples, and we performed a

177 restriction digest of this solution using XbaI, EcoRI-HF, and NheI-HF. Following initial
178 digestion, we added T4 DNA ligase to the digested DNA without disabling the restriction
179 enzymes, and we cycled temperatures to sequentially promote ligation of adapters followed by
180 digestion of chimeras and dimers. We cleaned the resulting reactions with NaCl-PEG diluted
181 Speedbeads (Rohland & Reich 2012), and we completed the adapter sequences using PCR with
182 iTru5 and iTru7 primers (see Glenn *et al.* (2016a) for details about primers). We pooled the
183 resulting libraries, size-selected fragments of 550bp (+/- 10%) with a PippinPrep (Sage Science,
184 Inc.), and performed a final round of low-cycle PCR-recovery using the P5 and P7 primers to
185 increase the concentration of fragments in the desired size range. We sequenced samples on an
186 Illumina NextSeq v2 300 cycle kit to obtain paired-end 150nt (PE150) reads (Figure 1).

187 We used the *process_radtags* program in Stacks v1.29 (Catchen *et al.* 2013; Catchen *et*
188 *al.* 2011) to clean and demultiplex the resulting sequence data. We “rescued” sequence tags and
189 RAD-tags within 2bp of their expected sequence; otherwise, we removed reads with an uncalled
190 base or containing the wrong adapter or wrong cut site. Because our 3RAD adapter sequences
191 vary in length, and because Stacks requires all reads to be the same length, we truncated reads to
192 140bp, removing 0-3 bases of sequence per read, in *process_radtags*. We ran the Stacks pipeline
193 with the following modifications: in the *ustacks* program, we removed highly repetitive stacks,
194 we used the deleveraging algorithm, and we set the maximum distance between stacks (M) to 3;
195 in the *cstacks* program, we set the number of mismatches allowed between sample tags when
196 generating the catalog (n) to 4; in the program *populations*, we required at least 3 individuals to
197 have reads to retain a given locus (r), and we set the minimum stack depth required for
198 individuals at a locus (m) to 3. We output the full sequence from each allele identified across

199 our pilot samples in FASTA format. We selected loci that were polymorphic, but had less than
200 five SNPs across both paired end reads, and that were present in three or four of the samples,
201 which resulted in 1740 paired reads (candidate loci) for bait design.

202

203 Bait Design and Synthesis

204 We selected bait sequences to minimize target redundancy and bait-to-bait hybridization,
205 which can compromise the synthesis of ssRNA baits as well as the target capture hybridization
206 reaction. To perform these steps, we subjected sequences to self-analysis using BLAST 2.2.19
207 (*filter query sequence = false, word size = 11, e-value = 1e-13, number of sequences to show*
208 *alignments for = 2,000*; Boratyn *et al.* 2013). We discarded any locus with one or both
209 sequences having a BLAST hit of at least 140bp to another sequence (682 loci). Next, we
210 subjected sequences to a same-strand self-analysis in BLAST (as above, *query strand = bottom*).
211 We discarded 94 additional loci in which one or both paired sequences had a BLAST hit to
212 another sequence, leaving 964 loci. Then, we designed two sets of 90mer baits targeting the
213 remaining 964 loci. In the first set, we chose a single bait from both paired sequences for every
214 locus, and we positioned baits to start at the 20th base of their parent sequences (creating 1,928
215 *Wisteria* baits; Supplementary File: Wist-Probes-Set1.fasta). In the second set, we added
216 additional baits from both sequences corresponding to a random subset of 200 loci (creating 400
217 additional baits; Supplementary File: Wist-Probes-Set2-SUBSET-400.fasta), and we positioned
218 these baits to start at the 40th base of their parent sequences. The two sets produced a total of
219 2,328 baits targeting *Wisteria* library molecules. To reduce synthesis costs, we combined this
220 bait set design with a similar number of baits designed in the same way for another species
221 (*Pueraria montana* var. *lobata*, kudzu). We subjected the bait sequences for both species to a

222 final same-strand self-analysis using BLAST (same process as above), and we did not find
223 evidence of additional bait-to-bait hybridization. Before bait synthesis and because MYbaits[®]
224 cannot be synthesized with a mixture of bases, we replaced any variable positions in any bait
225 sequence with a random candidate base, and we replaced all unknown (“N”) positions with a
226 thymine. We created a custom set of biotinylated RNA baits by having them synthesized as a
227 MYbaits-1 kit (MYcroarray, Ann Arbor, MI, USA).

228

229 Sample Library Preparation and Description of Treatments

230 We provide detailed sample collection and sample preparation methods in the
231 Supplemental Methods. Briefly, we randomly arranged 192 of the 203 greenhouse-grown
232 *Wisteria* samples in 2 plates (RADcap_Plate1 and RADcap_Plate2; Supplementary Tables 1 and
233 2). We placed the remaining 9 greenhouse-grown samples into a third plate, to which we added
234 randomly selected columns of the DNA in plates 1 and 2 (RADcap_Plate3; Supplementary Table
235 3). This arrangement allowed us to re-process 133 libraries independently prepared from the
236 same samples, and we used these replicates to compute the amount of missing data between
237 replicate samples that was not caused by genetic variation. Five samples were included twice
238 across plates 1-3, and 41 samples were included three times across plates 1-3. We arranged the
239 192 samples from wild-collected individuals by DNA concentration in a fourth and fifth plate
240 (RADcap_Plate4 and RADcap_Plate5; Supplementary Table 4 and 5). We normalized DNA
241 concentration, then we digested the plated DNA with XbaI, NheI-HF, and EcoRI-HF in a
242 reaction that included forward and reverse adapters. As before, we added T4 DNA ligase to the
243 digested DNA without disabling the restriction enzymes, and we cycled temperatures to

244 sequentially promote ligation of adapters followed by digestion of chimeras and dimers (Figure
245 2; see Supplemental Methods for full details). Following adapter ligation, we combined
246 approximately 66% of the ligation volume from each sample in each plate into plate-specific
247 pools, we cleaned each pool with speedbeads, and we re-suspended cleaned pools in 33ul of
248 TLE. For plates 1-4, we split each pool into three aliquots of 20μl, 10μl, and 3μl, and we used
249 these aliquots to test the effect of different PCR conditions on the efficiency of RADcap
250 (described below; Table 1).

251

252 *Single molecule tagging*

253 To tag and track duplicate reads that resulted from the PCR amplification process, we
254 designed a new iTru5 i5 primer (Glenn *et al.* 2016a) that incorporated a random 8 nucleotide
255 sequence tag (i.e., the i5 index sequence was specified as NNNNNNNN when ordering the
256 iTru5-8N primer). This resulted in the synthesis of a mixture of 65,536 iTru5 i5 primers with
257 unique, 8 nucleotide index sequences. In the experimental treatments, below, we incorporated
258 these uniquely tagged iTru5-8N i5 primers to our DNA library constructs using different PCR
259 conditions to determine what methods produce the fewest PCR duplicates.

260

261 *One-primer amplifications*

262 Following adapter ligation and cleaning, we split each 20μl aliquot into two tubes to
263 increase the total PCR volume possible, and we performed a single-cycle, one-primer PCR.
264 Each reaction contained 10μl template DNA and the iTru5-8N primer. Because we amplified
265 each reaction using only one cycle, the primers did not denature from the library molecules and

266 re-anneal to different library molecules. We pooled the two resulting reactions and cleaned them
267 with speedbeads, and we split them into 2 tubes for a 6-cycle PCR where we included the P5
268 primer and the plate-specific iTru7 primer (Supplementary Table 13). This second reaction
269 completed the library construct, added the plate-specific i7 index sequence to each library
270 construct, and increased the total amount of library available for capture. We called the plates in
271 this treatment RADcap_1cycle_Plate1-4.

272

273 *Two-primer amplifications*

274 For the second aliquots of 10 μ l, we performed four PCRs for each pooled plate with 2 μ l
275 template DNA in each. We included both the iTru5-8N primer and the iTru7 primer in each
276 PCR, and we ran PCR for 5 cycles. Because we included the iTru5-8N primer in the PCR
277 reaction for multiple cycles, newly synthesized molecules could receive new i5 tags (Casbon *et*
278 *al.* 2011) and thus a single template DNA molecule could generate multiple library constructs
279 with unique i5 sequence tags (i.e., this method produced ≤ 10 undetectable PCR duplicates per
280 template molecule). Because we used libraries in these treatments that were identical to those
281 used above (i.e., one-primer amplification with a single-cycle PCR), this experiment allowed us
282 to determine the effect of low-efficiency first-strand replication and test how additional PCR
283 cycles affect the identification of PCR duplicates and subsequent variant calling. We called the
284 plates in this treatment RADcap_5cycle_Plate1-4.

285

286 *Low-template, one-primer amplifications*

287 We used the 3 μ l aliquot from plate 1 to determine the effect of low DNA concentrations
288 on PCR duplication and subsequent variant calling. We added the iTru5-8N primer to 3 μ l of
289 template from plate 1, and we performed a single-cycle PCR. As we described for the one-
290 primer amplifications above, we cleaned the resulting PCR product and performed another 6-
291 cycle PCR to add the iTru7 primer. We called this treatment RADcap_Low_Template_Plate1.

292

293 *RAD locus capture*

294 Following all final PCRs described above, we pooled the replicate PCRs, cleaned each
295 plate pool with speedbeads, and performed a separate capture hybridization reaction on each pool
296 from each treatment (9 total) according to the MYcroarray MYbaits v3.0 protocol, with a
297 hybridization temperature of 65°C for 21 hours. Following capture, we split each capture into
298 three tubes and amplified loci with the P5 and P7 primers in an 18 cycle PCR. PCRs from each
299 capture were pooled and cleaned with speedbeads. Following PCR recovery and cleanup, we
300 quantified the 9 experimental treatments, and we pooled these libraries with unrelated libraries
301 from other experiments at a ratio that would return 20% of the reads from an Illumina
302 sequencing run. We sequenced the pooled libraries using an Illumina NextSeq High Output v2
303 150 cycle kit to achieve PE75 reads.

304

305 *Blocking over-enrichment and testing RADcap versus size selection*

306 After sequencing, the data from plates 1-4 included 8 loci with an average coverage 20x
307 higher than other loci in the one-primer treatment, 10x higher than other loci in the two-primer
308 treatment, and 28x higher than other loci in the low-template treatment. To block the over-

309 enrichment of these loci, we designed and ordered 29 custom oligonucleotides (Supplementary
310 Table 6) between 26 and 60 bp long that were complementary to the baits targeting these 8 loci
311 and which had a DNA to RNA T_m greater than 70°C. We also optimized the PCR for plate 5,
312 based on the one-primer treatment above, by increasing reaction volumes 3-fold for the PCR
313 reaction to add the iTru5-8N primer and 1.3-fold for the PCR to add the iTru7 primer, and we
314 included the locus-specific bait blockers during the hybridization reaction. To compare the
315 results of capturing RAD loci to those of size selection normally done in 3RAD (and other
316 RADseq protocols), we split plate 5 in half following the one-primer PCR, and we captured loci
317 from one-half of the plate, as described above. We called this treatment
318 RADcap_optimized_Plate5. We size-selected the remaining half of plate 5 samples as described
319 above, for SNP discovery. We called this treatment 3RAD_SizeSelect_Plate5. We pooled these
320 two libraries with unrelated libraries to obtain 7% of the reads on a second Illumina NextSeq run
321 using conditions described above for the other RADcap libraries.

322

323 Data Analysis

324 *Modification of Stacks software*

325 Stacks (Catchen *et al.* 2013; Catchen *et al.* 2011) had previously been modified to
326 identify the variable-length internal tags that distinguish individual samples in 3RAD data.
327 However, no software program existed to properly identify and remove the PCR duplicates from
328 RADseq data. We developed and implemented new code as part of the *clone_filter* module
329 within Stacks v1.35 to remove PCR duplicates. *Clone_filter* can be used before or after
330 *process_radtags* and can use any combination of inline or index sequence tags, in addition to

331 using read sequences, to reduce duplicated reads to a single representative in the output.

332 Importantly, *clone_filter* does not modify FASTQ headers, allowing repeated use of

333 *process_radtags* and *clone_filter* for read demultiplexing and duplicate removal.

334

335 *Sample demultiplexing, alignment, and SNP calling*

336 After sequencing, we converted BCL files to FASTQ format using *bcl2fastq2* v2.16.0.10

337 (Illumina, Inc.), and we modified the default parameters to create a separate FASTQ file for

338 index reads (Figure 1). We demultiplexed and removed PCR duplicates from the FASTQ data

339 using *Stacks* v 1.35. First, we demultiplexed reads by *i7* tag (Supplementary Table 13) using

340 *process_radtags*. We discarded reads with an uncalled base, reads having low quality (using

341 default settings), or reads having a sequence tag or RAD tag more than 2 bases distant from the

342 expected sequence, and we rescued reads having sequence tags or RAD tags and within 2 bases

343 of the expected sequence. This initial demultiplexing produced paired end files corresponding to

344 each plate in each treatment. We ran *process_radtags* again on each plate of samples, with the

345 same parameters, to demultiplex reads by inner adapter, which produced paired end files for each

346 individual in each plate. Finally, we used the *clone_filter* program to remove any read having

347 the same combination of random *i5* tag and insert sequence, which likely represent duplicates

348 created during PCR amplification.

349 We created a FASTA-formatted “radnome” file that contained the 964 paired sequences

350 from which we designed baits, and we used this file as a reference sequence for read alignment

351 and SNP calling (Supplementary File: *wisteria_reference.fasta*). Within this FASTA file, paired

352 reads were separate entries given arbitrary locus names, and we inserted 20 Ns between the

353 sequences for Read 1 and Read 2. We aligned RADcap reads to the reference using BWA v
354 0.7.7 (Li & Durbin 2009) with the mem algorithm and shorter split hits marked as secondary
355 (M), and we called SNPs using an automated pipeline (<https://github.com/faircloth-lab/radcap>)
356 that incorporates BWA, PICARD, and the open-source GATK-lite package (DePristo *et al.*
357 2011; McKenna *et al.* 2010). Following automated BWA alignment, the pipeline merged
358 individual alignments, re-aligned BAM files around indels, called SNPs and indels, and filtered
359 problematic or low-quality SNP calls from the total set of raw SNP calls to create a passing file
360 of SNPs.

361 Variant calling is an inherently population-based process in that errors can be
362 distinguished from variants at a specific position by considering that position in all individuals in
363 the population (Bansal *et al.* 2010; Catchen *et al.* 2011; Craig *et al.* 2008). Therefore, the
364 detection and statistical properties of variant genotypes are dependent on how the population
365 under study is sampled, with fewer variant sites recovered with lower statistical support from
366 smaller populations. To mimic this effect of population sampling and to facilitate comparisons
367 among our experimental treatments, we called SNPs in two ways. First, we treated all 384
368 individuals from plates 1-4 as a single population, and we called SNPs separately in each of the
369 one-primer (n=384 individuals) and two-primer experimental treatments (n=384 individuals).
370 Second, we treated the 96 samples in plate 1 as a single population, and we called SNPs for the
371 plate 1 population in the one-primer, two-primer, and low-template treatments, as well as the
372 plate 5 optimized and size selected treatments (Table 1). After SNP calling, we filtered the
373 resulting VCF files using vcftools v0.1.12b (Danecek *et al.* 2011) to exclude sites with more than
374 50%, 20%, or 10% missing data (i.e., 50%, 80%, or 90% complete data), and we computed

375 summary statistics across captured loci and variant sites using a program

376 (`radcap_summarize_snp_calls.py`) from the *radcap* software package.

377

378 *Efficiency of Random Tagging at the i5 Index Position*

379 Because PCR can be biased by the composition of certain primers, we wanted to estimate
380 how well our iTru5-8N primers were incorporated into our library constructs. Using the FASTQ
381 file of index reads as input, we determined the count of each iTru5-8N sequence tag using FastX
382 v0.0.14 (Gordon & Hannon 2010; Supplementary file: `i5_and_coverage_code.md`). We plotted
383 the cumulative count of iTru5-8N sequence tags incorporated to DNA libraries for all possible
384 sequence tag combinations, except for iTru5-8N tags that do not return a signal on the NextSeq
385 (GGGGGGGG), those DNA inserts that have no apparent i5 sequence tag (AGATCTCG), and
386 those iTru5 sequence tags in the adapters ligated to other libraries on the sequencing run.

387

388 *RADcap Efficiency and Coverage*

389 We expected sequence capture to be more efficient than size selection and that the
390 resulting data from captured RAD loci would include fewer off-target reads, have higher
391 coverage at target loci, and consistently recover a larger number of target loci from reads. To
392 investigate these parameters, we computed the coverage of each position in each sample from
393 BAM files using SAMtools v1.2 (Li *et al.* 2009; Supplementary file:
394 `i5_and_coverage_code.md`). For this analysis, we used the BAM files produced directly from
395 BWA to avoid effects of the BAM re-alignment on our coverage computations and because we
396 wanted to assess which loci were present in the dataset (where coverage of loci in the radnome

397 reference was greater than 0), despite being monomorphic or having errors. We report the
398 average coverage for all loci and samples within each plate, normalized by million reads per
399 sample. In order to determine whether the variation in coverage between loci in a treatment
400 decreased in the optimized treatment, we plotted the log transformed coverage of each locus and
401 tested whether the optimized treatment had less variation in log transformed coverage using a
402 one-sided Siegel-Tukey Test for equality in variability with adjusted medians in DescTools
403 (Signorell 2015) in R. We then calculated the average coverage per locus per million reads per
404 sample for loci with at least 4x coverage in plates RADcap_1cycle_Plate1,
405 RADcap_5cycle_Plate1, RADcap_Low_Template_Plate1, RADcap_optimized_Plate5, and
406 3RAD_SizeSelect_Plate5. As a measure of consistency and to see if the same loci were
407 recovered in each treatment, we identified the loci with at least 4x coverage in 90% of samples
408 from each treatment and determined the loci in common between treatments using VennDiagram
409 (Chen & Boutros 2011) in R. In addition, we plotted the density kernel of the coverage for Read
410 1 and Read 2 for each of the five treatments, and compared the distributions of coverage between
411 treatments in a one-sided two-sample Kolmogorov–Smirnov test in R.

412 To determine how many reads were necessary to recover all of the targeted loci at
413 reasonable coverage, we plotted the number of loci at or above 4x coverage and the number of
414 reads for each sample. To get more resolution at lower read numbers, we took the median
415 coverage for all samples at each locus and divided that to get corresponding coverage between
416 1,000 and 100,000 reads per sample. We plotted the number of loci at or above 4x, 10x, and 20x
417 coverage as a function of the reads per sample.

418

419 *Error Rate and Genetic Diversity of Wisteria*

420 We calculated the frequency of missing data between replicate samples within the one-primer
421 and two-primer treatments by converting VCF output files to Genepop format in PGDSpider v.
422 2.0.9.1 (Lischer & Excoffier 2012) and counting the number of SNPs at which one sample had a
423 base called while another did not. Since there were no replicate samples within plate 5, we could
424 not assess the amount of missing data. Instead, we compared estimates of genetic diversity of
425 *Wisteria* in plates RADcap_optimized_Plate5 and 3RAD_SizeSelect_Plate5 from 80% filled
426 matrices in GenAIEx v6.502 (Peakall & Smouse 2006; Smouse & Peakall 2012). For each plate,
427 we report the average number of samples genotyped (out of 96) across all loci, the number of
428 alleles identified, the effective number of alleles, Shannon's Information Index, the observed and
429 expected heterozygosity, and the fixation index, along with standard error estimates for each
430 parameter.

431

432 **Results:**

433 Initial 3RAD SNP Discovery for Bait Design

434 Following SNP discovery using four *Wisteria* samples, we obtained 1.4 to 2.5 million
435 PE150 reads per sample, and we retained an average of 83.7% of reads after quality filtering.
436 We identified 31,686 loci placed in the Stacks catalog, 1350 of which were sequenced in all four
437 samples, and 3483 of which were sequenced in three samples. Of the loci recovered in at least
438 three samples, 2573 loci were polymorphic and contained a total of 6531 variant sites. After
439 filtering these loci, there were 1428 putative variants in the 964 loci we used to design our
440 capture baits.

441

442 Random Tagging at the i5 Index Position Allows Removal of PCR Duplicates

443 For the RADcap samples in plates 1 to 5, we obtained 3-40 million reads per plate
444 (average 17 million), and we retained >94% of reads after quality filtering (Table 2). We
445 incorporated and sequenced all 65,536 of the expected i5 random sequence tags in both of the
446 Illumina NextSeq runs performed to generate our data, and after removing tags from other
447 libraries on the run, and tags indicating no i5 or no signal, we did not observe major sequence
448 dependent biases that affected PCR efficiency of primers containing different random i5
449 sequence tags (Supplementary Figure 1). All plates from which we collected data using RADcap
450 had a similar percent of reads retained after quality filtering by *process_radtags* in Stacks (Table
451 2). We retained an average of 68.9% of reads after decloning (range 20.4 to 95.7%; Table 2,
452 Supplementary Tables 7-11), with the most reads retained from the optimized PCR protocol
453 (which we performed on RADcap_optimized_Plate5 and 3RAD_SizeSelect_Plate5) and
454 RADcap_5cycle_Plate3.

455

456 Optimizing RADcap Efficiency and Coverage

457 All but one of the capture treatments yielded $\geq 80\%$ of reads on target (Table 2), while the
458 optimized treatment (RADcap_optimized_Plate5) yielded the highest proportion of reads on
459 target (90%). More traditional 3RAD with size selection (3RAD_SizeSelect_Plate5) yielded
460 15% of reads on target. Similarly, the optimized and two-primer treatments had the highest
461 average coverage, at 764 and 830 reads per million reads per sample, respectively (Table 2), but
462 we note that the two-primer coverage is inflated with undetected duplicate sequences from

463 multiple rounds of PCR. The one-primer and low-template treatments had slightly lower average
464 coverages, at 712 and 612 reads per million reads per sample, respectively. The size-selected
465 treatment had the lowest average coverage at 142 reads per million reads per sample.

466 The coverage per locus per million reads was much higher among the RADcap samples
467 than traditional 3RAD size-selected samples (Figure 3). The increased performance of RADcap
468 is also apparent when the coverage per locus is plotted as a density distribution (Supplementary
469 Figure 2). Among RADcap samples, the variation in coverage per locus per million reads
470 sequenced per sample was lower for RADcap_optimized_Plate5 than 3RAD_SizeSelect_Plate5
471 and RADcap_1cycle_Plate1 ($p < 0.0083$ in both cases), but the variation in coverage for the
472 RADcap_optimized_Plate5 did not differ from the RADcap_5cycle_Plate1 or the
473 RADcap_Low_Template_Plate1 ($p > 0.37$ in both cases; Figure 3).

474

475 RADcap Effectively and Consistently Enriched Target Loci and Produces Dense SNP Matrices

476 We consistently recovered more targeted loci within RADcap treatments than traditional
477 3RAD with size selection (Figure 5; Supplementary Table 12). Specifically, the optimized
478 treatment performed the best, with 912 loci recovered at 50% matrix occupancy, 880 recovered
479 at 80% occupancy, and 820 recovered at 90% occupancy. The 96 samples analyzed from the
480 one-primer and two-primer treatments performed slightly poorer, with 840 and 874 loci
481 recovered at 50% matrix occupancy, 697 and 823 loci recovered at 80% matrix occupancy, and
482 552 and 764 loci recovered at 90% matrix occupancy, respectively. Traditional 3RAD with size
483 selection returned 821, 642, and 510 loci at the same levels of matrix occupancy. As expected,
484 RADcap_Low_Template_Plate1 showed the poorest performance, returning 730, 338, and 155

485 loci at the same levels of occupancy. The number of SNPs called within loci showed similar
486 patterns (Fig 5b), with RADcap_optimized_Plate5 performing better than all other treatments.
487 The effects of population size on the number of SNPs called is apparent in the differences we
488 observed between RADcap_5cycle_Plate1 and RADcap_5cycle_Plate1-4 and
489 RADcap_1cycle_Plate1 and RADcap_1cycle_Plate1-4.

490 Although the effectiveness of enrichment within a plate is one metric of consistency, the
491 more important metric for most researchers is the consistency with which reduced representation
492 approaches collect data across plates or from all individuals in a population. At 90% occupancy
493 for 4x coverage, more than half of the loci (516 of 964; 54%) were shared between all treatments
494 except low-template, an additional 286 loci (30%) were shared among all three RADcap
495 treatments, and 125 loci (13%) were shared among the RADcap_5cycle_Plate1,
496 RADcap_optimized_Plate5, and 3RAD_SizeSelect_Plate5 (Figure 4). Impressively,
497 RADcap_optimized_Plate5 contained data at 4x coverage for 962 of the 964 loci (99.8%; Figure
498 4). Only 34 loci (3.5%) were present in only two treatments, and only 2 loci (0.02%) were
499 present in a single treatment (Figure 4). Thus, most loci were present in most samples no matter
500 from which treatment they originated.

501 In both RADcap_optimized_Plate5 and RADcap_5cycle_Plate1, we recovered most of
502 the 964 loci in most samples regardless of the sequencing depth (Figure 6). By comparison, in
503 the 3RAD with size selection treatment, even samples with the largest number of reads did not
504 include as many loci as these RADcap treatments. When we modeled a reduced number of reads
505 over all samples for each locus in the RADcap_optimized_Plate5 treatment, we found that
506 20,000 to 30,000 reads were sufficient to capture all loci with at least 4x coverage, and 60,000

507 reads per sample were sufficient to achieve 10x coverage at all loci (Figure 7). To achieve 20x
508 or higher coverage at all loci, we estimated that $\geq 100,000$ reads per sample were required,
509 although it may not be economical to sequence all loci at 20x coverage or higher (see
510 Discussion).

511

512 Error rate and Genetic Diversity of *Wisteria*

513 The amount of missing data in samples replicated within a treatment was effectively
514 equal between 1-cycle and 5-cycle treatments (7.20% and 7.61%, respectively).
515 Because the optimized and size-selected treatments had the same samples and were filtered to
516 have the same occupancy, we show the effect of estimating diversity with a smaller dataset. In
517 the 80% occupancy matrices, we recovered 3744 SNPs in the optimized treatment and 2554
518 SNPs in the size-selected treatment. On average, two more samples were genotyped in the
519 optimized treatment than size-selected treatment for each SNP (Table 3). The number of alleles
520 number of effective alleles per SNP, Shannon's information index, and F_{IS} were higher for size-
521 selected samples than optimized samples. Although observed heterozygosity was the same,
522 expected heterozygosity was higher for size-selected samples.

523

524 **Discussion:**

525 Our overall goal was to develop a simple system that would efficiently sample a
526 consistent portion of the genome from large numbers of individuals at low cost. Our optimized
527 protocol achieves $\geq 4x$ coverage for $\geq 90\%$ of samples at 99.8% of targeted loci with $< 187,000$
528 reads per sample (Figure 4). This high sequencing depth means that the optimized protocol

529 could achieve $\geq 4x$ coverage for $>90\%$ of 96 samples at >900 loci with only an average of 20,000
530 reads per sample (Figure 7). Furthermore, a single researcher was able to process >1000 DNA
531 samples within a week and to sequence that pool with less than one-third of a NextSeq high
532 output run (i.e., $\sim \$1$ /sample in sequencing costs). RADcap performs exceedingly well and
533 achieves our overall goal. It is well known, however, that $4x$ coverage will often lead to
534 inaccurate genotypes and that deeper sequencing is needed for consistent and accurate
535 genotyping (DePristo *et al.* 2011; Sims *et al.* 2014). Fortunately, RADcap is also sufficiently
536 efficient that $10-20x$ coverage can be obtained for 90% complete matrices with affordable
537 amounts of sequencing (Figure 7).

538 We discuss each major facet of our RADcap approach, the results we have obtained, and
539 we note additional ways to implement and expand upon the RADcap technique (Table 4). Other
540 groups with different goals have also combined RADseq with sequence capture (Jones & Good
541 2015), such as Suchan *et al.*'s (2015) use of RADseq fragments as baits. While completing this
542 manuscript, a separate group published a similar method of sequence capture (Rapture; Ali *et al.*
543 2015), which shares the same general goal as our approach, and we discuss some similarities and
544 differences between RADcap and Rapture.

545

546 Initial 3RAD SNP Discovery and Bait Design

547 RADcap is not limited to 1000 loci. Our initial goal was to obtain SNP discovery data
548 for >2000 polymorphic loci, to design baits for the best ~ 1000 loci, and to retain a set of ~ 500
549 loci after data filtering that consistently produced high-quality genotypes. Stacks performed well
550 for the task of identifying SNPs and polymorphic loci - using four individuals for initial SNP

551 discovery yielded the desired number of polymorphic loci (2573). However, we recognize that
552 using only four individuals is *not* ideal because it limits the ability to identify polymorphic loci
553 that result from biological variation instead of errors. In addition, polymorphic loci with rare
554 alleles were likely not attained in this small sample due to ascertainment bias (Clark *et al.* 2005;
555 Nielsen 2000). For future RADcap projects, we will use 16-96 individuals for SNP discovery.
556 One constraint of our current approach is that the pilot-scale 3RAD experiment requires a
557 significant amount of time to complete, including a potential queue for the Illumina machine and
558 several weeks to synthesize baits. However, if a genome sequence is available for the focal
559 organisms, the genome could be digested *in silico* and loci with mapped SNPs could be used for
560 bait design.

561

562 Random Tagging at the i5 Index Position Allows Removal of PCR Duplicates

563 We were motivated to develop a system to remove PCR duplicates from RADseq
564 libraries to achieve several goals: 1) reduce artificial confidence in genotypes resulting from
565 undetected duplicates in ddRAD-type data, 2) satisfy the assumptions for data input to
566 probabilistic base callers such as GATK, and 3) develop a general approach to identify PCR
567 duplicates that would be easy to implement and optimize in a variety of experimental conditions.
568 To achieve these goals, we implemented a new iTru5 primer with random i5 index sequences
569 (iTru5-8N) for the Adapterama system (Glenn *et al.* 2016a; Glenn *et al.* 2016b; Glenn *et al.*
570 2016c) and a single-cycle of strand extension to incorporate iTru5 sequences into new library
571 strands. The advantages of our approach include its: a) low-cost, b) simplicity, and c) freedom
572 from requiring changes to standard Illumina sequencing or data processing protocols.

573 We used the iTru5-8N tag to successfully remove PCR duplicates from our data using
574 new additions to the Stacks codebase, which is desirable for methods such as Stacks that rely on
575 coverage to determine whether a variant is real or an artifact of PCR or sequencing (Casbon *et al.*
576 2011). Although there are $4^8=65,536$ possible iTru5-8N sequence tags for each locus, false
577 duplicates (i.e., independent DNA molecules with the same iTru5-8N sequence tag) will be
578 encountered at much lower coverage (Schweyen *et al.* 2014), similar to how a relatively small
579 group of people is likely to have a pair that share the same birthday (McKinney 1966). The
580 number of iTru5-8N sequence tags that we used to identify duplicates is much larger than tag
581 pools used in the past (e.g., Schweyen *et al.* 2014; Tin *et al.* 2015), allowing more than sufficient
582 depth of coverage after duplicate removal. The approach we created does not require researchers
583 to anneal complementary oligos within a large pool of oligos containing degenerate tags (cf.
584 Schweyen *et al.* 2014), which will produce a preponderance of double-stranded adapters with
585 mismatches.

586 While our approach is simple and powerful, it comes with several limitations. First, there
587 is an upper limit of 65,536 possible iTru5-8N sequence tags, thus the method we implemented in
588 Stacks uses the iTru5-8N sequence, plus the sequences to which any given iTru5-8N sequence
589 tag is incorporated to define duplicates, otherwise only 65,536 reads would be retained from any
590 library. Second, it is critical to use conditions that promote high efficiency of first strand
591 synthesis (i.e., our optimized treatment, and not the one-primer, two-primer, or low-template
592 treatments) to avoid high levels of PCR duplication. Finally, Stacks is currently the only software
593 that has been created and optimized to remove duplicates from these types of data.

594 Casbon *et al.* (2011) found that reduced amounts of template molecules going into PCR
595 increased the rate of PCR duplication. In contrast, our treatment with the least amount of input
596 DNA, RADcap_Low_Template_Plate1, had fewer duplicates than RADcap_1cycle_Plate1.
597 Thus, the specific conditions used for first strand extension are critical to producing a diverse
598 RAD library and can be even more important than the amount of template used. Comparing the
599 number of duplicates in the one-primer and two-primer treatments further illustrates the
600 importance of the conditions used for first strand extension. The strand extension conditions as
601 well as the starting DNA quality and quantity in plates 1-4 of the one-primer and two-primer
602 treatments were identical, thus the reduced number of duplicates identified in the two-primer
603 treatment relative to the one-primer treatment is due to hidden duplicates in the two-primer
604 treatment. The low-level of duplicates in the size-selected and optimized protocols (which have
605 only one cycle of first strand extension) demonstrates that high levels of duplicates are not
606 inevitable, and that careful optimization of reaction conditions can keep duplicates to quite low
607 proportions. However, the only way to know what percentage of the reads are duplicates is to
608 implement a strategy to detect duplicates, which also facilitates their removal. Thus, tagging and
609 removing duplicates is prudent for all RADcap and RADseq experiments.

610

611 Optimizing RADcap Efficiency

612 Our modifications of the 1-cycle treatment to the optimized (also single-primer)
613 treatment included increasing the PCR volume and adding locus-specific bait blockers for the
614 few loci that were over-abundant in the first set of RADcap reads. These modifications
615 decreased the variation in coverage among loci and increased the number of loci we recovered.

616 Although the bait-blockers produced an effect (the percentage of reads attributed to the blocked
617 loci decreased from 14.7% in the 1-cycle treatment and 8.1% in the 5-cycle treatment to 3.6% in
618 the optimized treatment), the size of this effect was modest. Thus, we surmise that increasing the
619 PCR volume used for first-strand synthesis was far more important. Unfortunately, because we
620 tested the one-primer and two-primer treatments on plates 1-4 while we tested the increased PCR
621 volume using plate 5, it is unclear whether the initial library quality of plate 5 was higher than
622 plates 1-4 or if the increased volume of the PCR reaction decreased the rate of PCR duplication.
623 We do not expect or have evidence that the DNA varied in any significant way among plates, but
624 we cannot rule out that possibility. Thus, although the optimized treatment conditions appear to
625 have produced a very significant reduction in duplicates, additional experiments are necessary to
626 definitively reach that conclusion.

627

628 RADcap Captures Nearly All Loci Targeted

629 We sequenced most of the targeted loci in RADcap_optimized_Plate5 and only slightly
630 fewer loci in RADcap_5cycle_Plate1 (Figure 4). The 99.8% overlap between
631 RADcap_optimized_Plate5 and RADcap_5cycle_Plate1 illustrates the strength of using
632 sequence capture to collect RAD loci: we were able to recover most of the same loci across at
633 least 90% of 192 samples that we prepared several weeks apart.

634 Although RADseq protocols reduce the genome being studied, the number of loci in a
635 typical analysis can still be in the tens of thousands. Our RADcap method allows for sequencing
636 efforts to be further focused, permitting higher levels of multiplexing and locus recovery in an
637 experiment while avoiding some of the problems such as biases in PCR amplification, variation

638 at restriction enzyme cut sites, and purification and size selection methods (DaCosta & Sorenson
639 2014; Gautier *et al.* 2013) that can occur with ddRAD or traditional RADseq experiments
640 implemented at the same scale.

641

642 RADcap Produces Dense Matrices

643 We found that the probabilistic base-caller, GATK, worked well on these data,
644 recovering and retaining large numbers of loci and SNPs at 50%, 80%, and 90% matrix
645 occupancy (Figure 5). The number of loci and SNPs called by GATK follows predictable
646 patterns within and among datasets from all treatments. We selected GATK because it is in
647 common use across a variety of genotyping studies and because it performs well for moderate- to
648 large-scale data sets. However, GATK is unsuited to most ddRAD data sets because of read
649 duplication, thus our system for removing duplicates was critical for meeting the assumptions of
650 GATK. Obviously, we could have used Stacks or other SNP-calling software packages (e.g.
651 FreeBayes, pyRAD, SAMTools), and subsequent work will provide a detailed comparison
652 among SNP-calling software packages.

653

654 RADcap Adds Relatively Few Errors to Illumina Sequences

655 Errors in RADseq data derive from library preparation and sequencing methods. Errors
656 introduced by PCR may be common in RADseq data because even extremely high fidelity DNA
657 polymerases introduce errors on the order of 2.8×10^{-7} per nucleotide incorporated (KAPA,
658 Boston, MA, USA). If these errors occurred in a single cycle of PCR, it would result in 4,683
659 errors in the 223 million reads in the present dataset. Because there can be many PCR cycles in

660 RADseq library preparation, an incorrect base incorporated during early cycles be amplified to
661 high coverage in the dataset. A much larger problem is the 0.1% substitution error rate made by
662 Illumina machines (Glenn 2011), which results in an additional 28,544,000 incorrect bases in a
663 dataset of 223 million PE64 reads. Even very small biases away from randomness in the
664 Illumina sequencing error distribution can create significant downstream problems in data
665 analysis. Decloning facilitated by the random iTru5-8N primer tags does not prevent PCR or
666 sequencing errors, but the use of probabilistic base calling algorithms as implemented in GATK
667 or other SNP-calling software can help reduce the likelihood of a base introduced by these errors
668 from being called as a true variant.

669

670 RADcap Works with Mixtures of Baits from Different Organisms Diluted to 1x Concentration

671 We used a bait set containing baits from two organisms – 2328 baits from our focal
672 *Wisteria* groups and 2624 baits from an unrelated organism (kudzu; *Pueraria montana* var.
673 *lobata*). As a result, baits for both species were present in all captures, despite DNA from only
674 one species being present in any given capture. We also synthesized fewer than the maximum
675 allowed number of baits (i.e. 4,952 baits in synthesis scale meant for up to 20,000 baits). The
676 large number of loci we captured and the dense genotype matrices we created suggest that there
677 was no meaningful interference from the additional baits during sequence capture and that the
678 concentration of baits we applied to each sample pool was sufficient.

679 By mixing baits for two different projects, we were able to reduce baits costs by 50% for
680 each project. The MYbaits-1 synthesis allows ~20,000 baits and the smallest synthesis scale is
681 sufficient for 12 captures. If fewer than 20,000 baits are needed, then the concentration of baits

682 is increased proportionately (e.g., 20,000 baits at 1x or 10,000 baits at 2x). Baits can also be
683 divided among multiple projects. For example, researchers could pool four different projects,
684 each with 1000 baits, into a single MYbaits-1 synthesis to achieve baits at 5x normal
685 concentration. This would mean workers could achieve 60 (12×5) captures of 96 samples, rather
686 than 12. Divided evenly, this would allow 15×96 (1440) samples from each of the 4 projects to
687 be enriched when purchasing the smallest possible single MYbaits-1 synthesis. Obviously,
688 different numbers of taxa with varying numbers of samples could be accommodated (e.g., 5,000
689 samples from one taxon and 760 samples from the other three; see Heyduk *et al.* (2016) for
690 additional examples). This flexibility creates many opportunities for collecting RADcap data at
691 low cost.

692

693 Comparison of RADcap to Rapture

694 Rapture is a similar, enrichment-based, RAD sequencing approach that uses a two-step
695 protocol to capture RADseq loci. In the first step, researchers digest, randomly shear, and ligate
696 biotinylated adapters to DNA, which can then be separated from other genomic fragments with
697 magnetic beads before library preparation continues. In the second step, similar to RADcap,
698 custom library-specific baits are used to capture loci of interest. RADcap and Rapture are
699 similar in that they require DNA isolation, restriction enzyme digests, ligation of adapters,
700 pooling, clean-up, capture, and sequencing. Both methods are significant advances that increase
701 the density and consistency of genotype matrices while simultaneously reducing costs for large-
702 scale projects.

703 There are, however, significant differences in cost and flexibility between RADcap and
704 Rapture as a result of RADcap's integration with 3RAD and the Adapterama system. The 3RAD
705 adapters require 8 phosphorylated oligos and 32 plain oligos to achieve 96 combinations (Glenn
706 *et al.* 2016b), whereas Rapture requires 96 biotinylated oligos plus 96 phosphorylated oligos,
707 making Rapture adapters 10 times more expensive (\$370 for RADcap vs. \$3750 for Rapture).
708 Adding or switching to different enzymes in Rapture requires additional sets of adapters at \$3750
709 per enzyme, whereas 3RAD facilitates the use of many different possible enzymes and
710 combinations of enzymes (Glenn *et al.* 2016b) with fewer sets of interchangeable adapters at
711 \$370 per set. In addition, RADcap does not require any commercial library preparation kits,
712 whereas Rapture makes use of commercial library preparation kits.

713 In addition to cost differences, duplicate detection has fewer false positives in RADcap
714 than Rapture. Rapture detects duplicates based on the starting position of Read 2, which may be
715 anywhere along ~500 bases (following shearing and size selection). Thus, RADcap has 65,536
716 tags, whereas Rapture has ~500. Additionally, dual digest RADcap increases coverage at both
717 ends of the library molecule, and we show that fewer reads per sample are required to achieve
718 the same coverage with RADcap than Rapture (20,000 versus 50,000 for at least 4x coverage,
719 respectively). On the other hand, Rapture's use of random shearing increases the length of the
720 genomic region that is sequenced, which may be an advantage worth the trade-off in decreased
721 coverage, depending upon the goals of the project.

722

723 Future Improvements and Extensions

724 RADcap (and Rapture) open the door to a variety of additional research opportunities.
725 One of the most important is the option of using the capture baits from RAD loci on randomly-
726 sheared genomic libraries (i.e., standard genomic libraries). Such work will facilitate direct
727 comparisons between RAD loci and other loci commonly used for sequence capture (exons,
728 UCEs, anchored loci, etc.). Although preparing randomly sheared genomic libraries for RADcap
729 increases the cost per sample, it will allow: 1) assembling contigs at captured loci so that more
730 sequence is available to facilitate a) ortholog identification in other species, b) identifying
731 additional linked SNPs, c) phasing SNPs within captured loci, and d) better understanding of the
732 sequence context for the RADcap loci; 2) investigating rates of divergence at restriction sites; 3)
733 collecting RAD loci from samples with deeper divergences than is feasible with restriction sites
734 (i.e., for phylogenetics), and 4) using PHYLUCE (Faircloth 2015) and other analytical tools that
735 have been developed for UCEs and other sequence capture systems. Capture baits also facilitate
736 using RADseq for degraded and contaminated samples (cf. Graham *et al.* 2015) and focusing on
737 microsatellite loci present in RADseq libraries, either through the use of locus-specific baits that
738 target the flanking regions or via generic baits to the repeats (cf. Glenn & Schable 2005). Thus,
739 the baits identified with RADcap will serve a variety of purposes in future work.

740 Given our high efficiency with two baits per locus, future work should investigate the
741 efficiency of capturing loci with a single bait per locus where we suspect the success rate will
742 remain high because of the reduced diversity in RAD libraries. We have presented datasets and
743 results obtained from the first sets of samples we have conducted with RADcap. Although we
744 are bullish on the future of RADcap, we expect researchers will explore and discover how

745 changes to the protocols (often unintentional), as well as their use on different organisms, impact
746 the outcome.

747

748 Summary

749 We present a novel protocol to cheaply sequence a specific set of hundreds to thousands
750 of loci in hundreds to thousands of samples. We demonstrate a generalizable method for
751 identifying PCR duplicates in Illumina libraries, even double-digest, RADseq-type libraries. We
752 also show that it is possible to reduce PCR duplicates to 5% of the total library, to routinely
753 achieve >80% on-target reads, and to achieve dense matrices of genotypes from hundreds of
754 individuals. Our method is sufficiently efficient that researchers can choose to sequence less
755 deeply to achieve commonly observed coverage in RADseq-type datasets, or they can affordably
756 sequence much more deeply to obtain high-quality genotypes. Molecular ecology research is
757 filled with choices about where to spend limited resources (time and money). We strongly
758 recommend that researchers adopt methods that yield high coverage and dense matrices of high-
759 confidence genotypes, and we hope that RADcap allows other scientists to obtain high-quality
760 data and make more robust conclusions about their study systems.

761

762 **Acknowledgements**

763 We thank Todd Pierson, Kerin Bentley, and Natalia Bayona. This work was supported by grants
764 DEB-1242260 and DEB-1146440 from the U.S. National Science Foundation and the U.S.
765 National Science Foundation Partnership for International Research and Education (PIRE)
766 program (OISE 0730218). This study was also supported in part by resources and technical

767 expertise from the Georgia Advanced Computing Resource Center, a partnership between the
768 University of Georgia's Office of the Vice President for Research and Office of the Vice
769 President for Information Technology.

770

771 **Competing Interests**

772 The authors declare competing interests. AD is employed by MYcroarray, a for-profit business
773 that sells customized MYbaits kits. TJK and TCG are partially supported through cost-recovery
774 projects in the EHS DNA lab, including projects that use 3RAD, and both are likely to use
775 RADcap in the future.

776

777 **Author Contributions**

778 Conceived RADcap: TCG, BCF; conceived *Wisteria* project: SLH, RM; conceived and
779 implemented decloning software: JMC; prepared samples: SLH, TJK; designed baits: AD;
780 analyzed data: SLH, BCF; wrote paper: SLH, BCF, TCG, RM; contributed funding and other
781 resources: TCG, RM; all authors edited and approved of the final version.

782

783

784 **Literature Cited**

- 785 Ali OA, O'Rourke SM, Amish SJ, *et al.* (2015) RAD Capture (Rapture): Flexible and efficient
786 sequence-based genotyping. *bioRxiv*, 028837.
- 787 Altshuler D, Pollara VJ, Cowles CR, *et al.* (2000) An SNP map of the human genome generated
788 by reduced representation shotgun sequencing. *Nature* **407**, 513-516.
- 789 Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA (2016) Harnessing the power of
790 RADseq for ecological and evolutionary genomics. *Nat Rev Genet* **17**, 81-92.
- 791 Andrews KR, Hohenlohe PA, Miller MR, *et al.* (2014) Trade-offs and utility of alternative
792 RADseq methods: Reply to Puritz *et al.* 2014. *Molecular ecology* **23**, 5943-5946.
- 793 Baird NA, Etter PD, Atwood TS, *et al.* (2008) Rapid SNP Discovery and Genetic Mapping
794 Using Sequenced RAD Markers. *PLoS ONE* **3**, e3376.
- 795 Bansal V, Harismendy O, Tewhey R, *et al.* (2010) Accurate detection and genotyping of SNPs
796 utilizing population sequencing data. *Genome research* **20**, 537-545.
- 797 Bi K, Vanderpool D, Singhal S, *et al.* (2012) Transcriptome-based exon capture enables highly
798 cost-effective comparative genomic data collection at moderate evolutionary scales. *BMC*
799 *Genomics* **13**, 1-14.
- 800 Boratyn GM, Camacho C, Cooper PS, *et al.* (2013) BLAST: a more efficient report with
801 usability improvements. *Nucleic Acids Research* **41**, W29-W33.
- 802 Cao H, Wu J, Wang Y, *et al.* (2013) An Integrated Tool to Study MHC Region: Accurate SNV
803 Detection and HLA Genes Typing in Human MHC Region Using Targeted High-
804 Throughput Sequencing. *PLoS ONE* **8**, e69388.
- 805 Cariou M, Duret L, Charlat S (2013) Is RAD-seq suitable for phylogenetic inference? An in
806 silico assessment and optimization. *Ecology and Evolution* **3**, 846-852.
- 807 Casbon JA, Osborne RJ, Brenner S, Lichtenstein CP (2011) A method for counting PCR
808 template molecules with application to next-generation sequencing. *Nucleic Acids*
809 *Research* **39**.
- 810 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA (2013) Stacks: an analysis tool
811 set for population genomics. *Molecular ecology* **22**, 3124-3140.
- 812 Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011) Stacks: Building and
813 Genotyping Loci De Novo From Short-Read Sequences. *G3-Genes Genomes Genetics* **1**,
814 171-182.
- 815 Chen H, Boutros PC (2011) VennDiagram: a package for the generation of highly-customizable
816 Venn and Euler diagrams in R. *BMC bioinformatics* **12**, 35.
- 817 Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R (2005) Ascertainment bias in
818 studies of human genome-wide polymorphism. *Genome research* **15**, 1496-1502.
- 819 Craig DW, Pearson JV, Szelinger S, *et al.* (2008) Identification of genetic variants using bar-
820 coded multiplexed sequencing. *Nature Methods* **5**, 887-893.
- 821 DaCosta JM, Sorenson MD (2014) Amplification Biases and Consistent Recovery of Loci in a
822 Double-Digest RAD-seq Protocol. *PLoS ONE* **9**, e106713.
- 823 Danecek P, Auton A, Abecasis G, *et al.* (2011) The variant call format and VCFtools.
824 *Bioinformatics* **27**, 2156-2158.
- 825 Davey JL, Blaxter MW (2010) RADSeq: next-generation population genetics. *Briefings in*
826 *Functional Genomics* **9**, 416-423.

- 827 Davey JW, Hohenlohe PA, Etter PD, *et al.* (2011) Genome-wide genetic marker discovery and
828 genotyping using next-generation sequencing. *Nature Reviews Genetics* **12**, 499-510.
- 829 DePristo MA, Banks E, Poplin R, *et al.* (2011) A framework for variation discovery and
830 genotyping using next-generation DNA sequencing data. *Nat Genet* **43**, 491-498.
- 831 Faircloth BC (2015) PHYLUCE is a software package for the analysis of conserved genomic
832 loci. *Bioinformatics*.
- 833 Faircloth BC, McCormack JE, Crawford NG, *et al.* (2012) Ultraconserved Elements Anchor
834 Thousands of Genetic Markers Spanning Multiple Evolutionary Timescales. *Systematic*
835 *biology* **61**, 717-726.
- 836 Fountain ED, Pauli JN, Reid BN, Palsbøll PJ, Peery MZ (2016) Finding the right coverage: the
837 impact of coverage and sequence quality on single nucleotide polymorphism genotyping
838 error rates. *Molecular Ecology Resources*, n/a-n/a.
- 839 Gautier M, Gharbi K, Cezard T, *et al.* (2013) The effect of RAD allele dropout on the estimation
840 of genetic variation within and between populations. *Molecular ecology* **22**, 3165-3178.
- 841 Glenn TC (2011) Field guide to next-generation DNA sequencers. *Molecular Ecology Resources*
842 **11**, 759-769.
- 843 Glenn TC, Faircloth BC, Nilsen R, *et al.* (2016a) Adapterama I: Universal stubs and primers for
844 thousands of dual-indexed Illumina libraries. *in prep.*
- 845 Glenn TC, Pierson T, Kieran TJ, *et al.* (2016b) Adapterama III: Quadruple-indexed triple-
846 enzyme RADseq libraries from picograms of DNA (3RAD). *in prep.*
- 847 Glenn TC, Pierson TW, Kieran TJ, *et al.* (2016c) Adapterama II: Universal amplicon
848 sequencing on Illumina platforms (TaggiMatrix). *in prep.*
- 849 Glenn TC, Schable NA (2005) Isolating Microsatellite DNA Loci. In: *Methods in Enzymology*,
850 pp. 202-222. Academic Press.
- 851 Gnirke A, Melnikov A, Maguire J, *et al.* (2009) Solution hybrid selection with ultra-long
852 oligonucleotides for massively parallel targeted sequencing. *Nature Biotechnology* **27**,
853 182-189.
- 854 Gordon A, Hannon G (2010) Fastx-toolkit. *FASTQ/A short-reads preprocessing tools*
855 (*unpublished*) http://hannonlab.cshl.edu/fastx_toolkit.
- 856 Graham CF, Glenn TC, McArthur AG, *et al.* (2015) Impacts of degraded DNA on restriction
857 enzyme associated DNA sequencing (RADSeq). *Molecular Ecology Resources*.
- 858 Harvey MG, Smith BT, Glenn TC, Faircloth BC, Brumfield RT (2013) Sequence capture versus
859 restriction site associated DNA sequencing for phylogeography. *arXiv preprint*
860 *arXiv:1312.6439*.
- 861 Heyduk K, Trapnell DW, Barrett CF, Leebens-Mack J (2016) Phylogenomic analyses of species
862 relationships in the genus *Sabal* (Arecaceae) using targeted sequence capture. *Biological*
863 *Journal of the Linnean Society* **117**, 106-120.
- 864 Hiatt JB, Pritchard CC, Salipante SJ, O'Roak BJ, Shendure J (2013) Single molecule molecular
865 inversion probes for targeted, high-accuracy detection of low-frequency variation.
866 *Genome research* **23**, 843-854.
- 867 Jabara CB, Jones CD, Roach J, Anderson JA, Swanstrom R (2011) Accurate sampling and deep
868 sequencing of the HIV-1 protease gene using a Primer ID. *Proc Natl Acad Sci U S A* **108**,
869 20166-20171.

- 870 Jones MR, Good JM (2015) Targeted capture in evolutionary and ecological genomics.
871 *Molecular ecology*.
- 872 Kenny EM, Cormican P, Gilks WP, *et al.* (2010) Multiplex Target Enrichment Using DNA
873 Indexing for Ultra-High Throughput SNP Detection. *DNA Research*.
- 874 Kivioja T, Vaharautio A, Karlsson K, *et al.* (2012) Counting absolute numbers of molecules
875 using unique molecular identifiers. *Nat Meth* **9**, 72-74.
- 876 Lemmon AR, Emme SA, Lemmon EM (2012) Anchored Hybrid Enrichment for Massively
877 High-Throughput Phylogenomics. *Systematic biology* **61**, 727-744.
- 878 Li C, Hofreiter M, Straube N, Corrigan S, Naylor GJ (2013) Capturing protein-coding genes
879 across highly divergent species. *Biotechniques* **54**, 321-326.
- 880 Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows–Wheeler transform.
881 *Bioinformatics* **25**, 1754-1760.
- 882 Li H, Handsaker B, Wysoker A, *et al.* (2009) The Sequence Alignment/Map format and
883 SAMtools. *Bioinformatics* **25**, 2078-2079.
- 884 Lischer HEL, Excoffier L (2012) PGDSpider: an automated data conversion tool for connecting
885 population genetics and genomics programs. *Bioinformatics* **28**, 298-299.
- 886 Mastretta-Yanes A, Arrigo N, Alvarez N, *et al.* (2015) Restriction site-associated DNA
887 sequencing, genotyping error estimation and de novo assembly optimization for
888 population genetic inference. *Molecular Ecology Resources* **15**, 28-41.
- 889 McKenna A, Hanna M, Banks E, *et al.* (2010) The Genome Analysis Toolkit: a MapReduce
890 framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-
891 1303.
- 892 McKinney EH (1966) Generalized Birthday Problem. *The American Mathematical Monthly* **73**,
893 385-387.
- 894 Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA (2007) Rapid and cost-effective
895 polymorphism identification and genotyping using restriction site associated DNA
896 (RAD) markers. *Genome research* **17**, 240-248.
- 897 Miner BE, Stöger RJ, Burden AF, Laird CD, Hansen RS (2004) Molecular barcodes detect
898 redundancy and contamination in hairpin-bisulfite PCR. *Nucleic Acids Research* **32**,
899 e135-e135.
- 900 Nielsen R (2000) Estimation of population parameters and recombination rates from single
901 nucleotide polymorphisms. *Genetics* **154**, 931-942.
- 902 Novaes E, Drost DR, Farmerie WG, *et al.* (2008) High-throughput gene and SNP discovery in
903 *Eucalyptus grandis*, an uncharacterized genome. *BMC Genomics* **9**, 1-14.
- 904 Okou DT, Steinberg KM, Middle C, *et al.* (2007) Microarray-based genomic selection for high-
905 throughput resequencing. *Nature Methods* **4**, 907-909.
- 906 Pante E, Abdelkrim J, Viricel A, *et al.* (2015) Use of RAD sequencing for delimiting species.
907 *Heredity* **114**, 450-459.
- 908 Peakall R, Smouse PE (2006) GENALEX 6: genetic analysis in Excel. Population genetic
909 software for teaching and research. *Molecular Ecology Notes* **6**, 288-295.
- 910 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012) Double Digest RADseq: An
911 Inexpensive Method for De Novo SNP Discovery and Genotyping in Model and Non-
912 Model Species. *PLoS ONE* **7**.

- 913 Pröll J, Danzer M, Stabentheiner S, *et al.* (2011) Sequence Capture and Next Generation
914 Resequencing of the MHC Region Highlights Potential Transplantation Determinants in
915 HLA Identical Haematopoietic Stem Cell Transplantation. *DNA Research* **18**, 201-210.
- 916 Puritz JB, Matz MV, Toonen RJ, *et al.* (2014) Demystifying the RAD fad. *Molecular ecology*
917 **23**, 5937-5942.
- 918 Rabenstein H, Bangol B, Linke S, *et al.* (2015) Large fragment target enrichment and sequencing
919 of the 4.4 Mb major histocompatibility complex by region-specific extraction and HiSeq
920 2500. *Human Immunology* **76, Supplement**, 57.
- 921 Raposo do Amaral F, Neves LG, Resende MFR, Jr., *et al.* (2015) Ultraconserved Elements
922 Sequencing as a Low-Cost Source of Complete Mitochondrial Genomes and
923 Microsatellite Markers in Non-Model Amniotes. *PLoS ONE* **10**, e0138446.
- 924 Rogers Y-H, Venter JC (2005) Genomics: Massively parallel sequencing. *Nature* **437**, 326-327.
- 925 Rohland N, Reich D (2012) Cost-effective, high-throughput DNA sequencing libraries for
926 multiplexed target capture. *Genome research* **22**, 939-946.
- 927 Saintenac C, Jiang D, Akhunov ED (2011) Targeted analysis of nucleotide and copy number
928 variation by exon capture in allotetraploid wheat genome. *Genome Biology* **12**, R88-R88.
- 929 Schweyen H, Rozenberg A, Leese F (2014) Detection and Removal of PCR Duplicates in
930 Population Genomic ddRAD Studies by Addition of a Degenerate Base Region (DBR) in
931 Sequencing Adapters. *Biological Bulletin* **227**, 146-160.
- 932 Shiroguchi K, Jia TZ, Sims PA, Xie XS (2012) Digital RNA sequencing minimizes sequence-
933 dependent bias and amplification noise with optimized single-molecule barcodes. *Proc*
934 *Natl Acad Sci U S A* **109**, 1347-1352.
- 935 Signorell A (2015) DescTools: Tools for descriptive statistics, p. R package version 0.99.15.
- 936 Sims D, Sudbery I, Illott NE, Heger A, Ponting CP (2014) Sequencing depth and coverage: key
937 considerations in genomic analyses. *Nature Reviews Genetics* **15**, 121-132.
- 938 Smith E, Jepsen K, Khosroheidari M, *et al.* (2014) Biased estimates of clonal evolution and
939 subclonal heterogeneity can arise from PCR duplicates in deep sequencing experiments.
940 *Genome Biology* **15**, 420.
- 941 Smouse PE, Peakall R (2012) GenAlEx 6.5: genetic analysis in Excel. Population genetic
942 software for teaching and research—an update. *Bioinformatics* **28**, 2537-2539.
- 943 Suchan T, Pitteloud C, Gerasimova N, *et al.* (2015) Hybridization capture using RAD probes
944 (hyRAD), a new tool for performing genomic analyses on museum collection specimens.
945 *bioRxiv*, 025551.
- 946 Tautz D, Ellegren H, Weigel D (2010) Next Generation Molecular Ecology. *Molecular ecology*
947 **19**, 1-3.
- 948 Tin MMY, Rheindt FE, Cros E, Mikheyev AS (2015) Degenerate adaptor sequences for
949 detecting PCR duplicates in reduced representation sequencing data improve genotype
950 calling accuracy. *Molecular Ecology Resources* **15**, 329-336.
- 951 Trusty JL, Johnson KJ, Lockaby BG, Goertzen LR (2007) Bi-Parental Cytoplasmic DNA
952 Inheritance in Wisteria (Fabaceae): Evidence from a Natural Experiment. *Plant and Cell*
953 *Physiology* **48**, 662-665.
- 954 Trusty JL, Lockaby BG, Zipperer WC, Goertzen LR (2008) Horticulture, hybrid cultivars and
955 exotic plant invasion: a case study of Wisteria (Fabaceae). *Botanical Journal of the*
956 *Linnean Society* **158**, 593-601.

- 957 Valder P (1995) *Wisterias: A comprehensive guide* Timber Press, Potland, Oregon.
958 Wiedmann RT, Smith TP, Nonneman DJ (2008) SNP discovery in swine by reduced
959 representation and high throughput pyrosequencing. *BMC Genetics* **9**, 1-7.
960 Wilson EH (1916) The Wisterias of China and Japan. *The Gardeners' Chronicle* **1545**, 61-62.
961 Wyman D (1949) The Wisterias. *Arnoldia* **9**, 17-28.
- 962

963 **Data Accessibility**

964 Raw Reads will be at NCBI SRA: xxxxxxxx.

965
966
967
968
969
970
971

972 Table 1: Overview of treatments, DNA plates used with each treatment, and how plates were
973 grouped for analyses. The iTru5-8N reaction volume for the size selected and optimized
974 treatment represents the same reactions, as these treatments were split after the single-primer
975 PCR and clean up.

Treatment	Plate ID's	Cycles with iTru5-8N	iTru5-8N reaction volume (μ l)	All PCR duplicates tagged?	Captured?	Analysis groups (plate ID's)
RADcap_1cycle	1 - 4	1	100	Yes	Yes	1; 1-4
RADcap_5cycle	1 - 4	5	100	No	Yes	1; 1-4
RADcap_Low_Template	1	1	25	Yes	Yes	1
RADcap_optimized	5	1	300	Yes	Yes	5
3RAD_SizeSelect	5	1	300	Yes	No	5

976
977
978

979 Table 2: The total reads, percent retained after quality filtering in *process_radtags*, percent
 980 retained after decloning with *clone_filter*, percent mapped with BWA mem algorithm for each
 981 plate, and the average coverage per million reads sequenced per sample of all loci.

Plate	Number of reads	% retained after quality filtering	% retained after decloning	% reads that map to reference	Average Coverage
RADcap_1cycle_Plate1	19,397,440	94.9	25.1	79.6	712
RADcap_1cycle_Plate2	14,703,752	95.0	20.4	85.7	-
RADcap_1cycle_Plate3	15,865,294	94.8	67.2	81.5	-
RADcap_1cycle_Plate4	17,907,048	95.0	63.3	83.8	-
RADcap_5cycle_Plate1	18,045,032	95.0	83.3	84.8	830
RADcap_5cycle_Plate2	17,968,264	95.0	86.9	85.5	-
RADcap_5cycle_Plate3	2,332,154	94.1	95.7	84.3	-
RADcap_5cycle_Plate4	18,455,900	95.1	86.1	84.1	-
RADcap_Low_Template_Plate1	3,285,096	95.3	41.0	65.8	612
RADcap_optimized_Plate5	17,929,096	95.5	94.2	90.1	764
3RAD_SizeSelect_Plate5	39,543,602	95.9	94.8	15.1	142

982

983 Table 3: The number of samples genotyped (N), number of alleles (Na), effective number of
 984 alleles (Ne), Shannon's Information Index (I), observed (Ho) and expected (He) heterozygosity,
 985 and fixation index (F_{IS}) as calculated in GenAlEx for plate 5 prepared via RADcap and 3RAD.

Plate	N	Na	Ne	I	Ho	He	F_{IS}
RADcap_optimized_Plate5	93.807 ± 0.061	2.014 ± 0.003	1.143 ± 0.004	0.169 ± 0.003	0.103 ± 0.003	0.095 ± 0.002	0.180 ± 0.007
3RAD_SizeSelect_Plate5	91.575 ± 0.106	2.018 ± 0.003	1.159 ± 0.005	0.184 ± 0.004	0.104 ± 0.004	0.104 ± 0.003	0.227 ± 0.008

986
 987
 988
 989

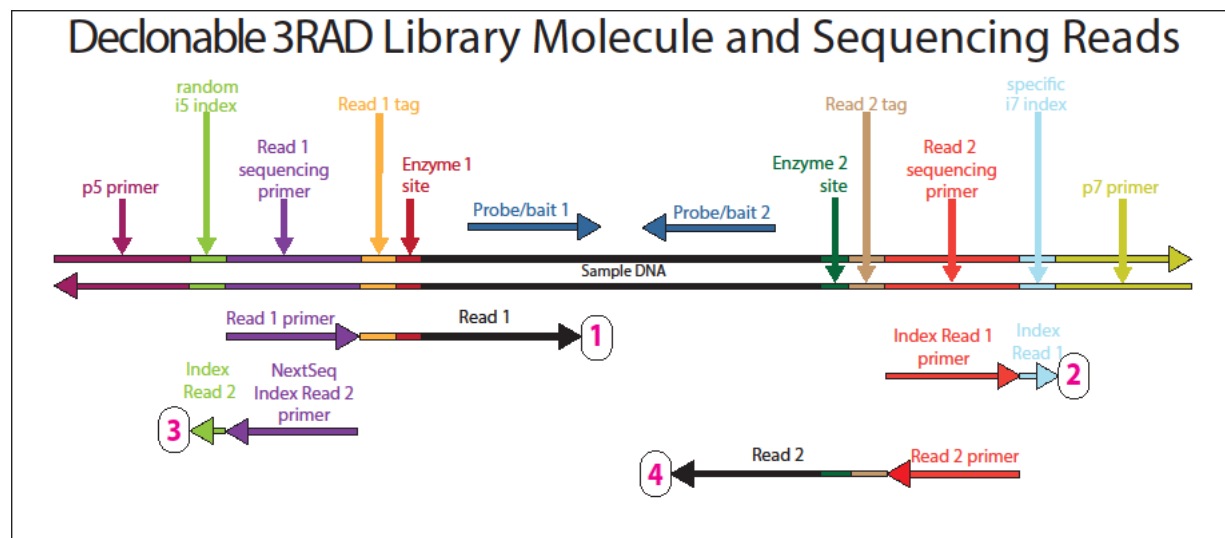
990 Table 4. Major processes and reagents of RADcap, current costs, and potential improvements to
 991 reduce costs and/or increase throughput. Current costs of processes are calculated on a per-
 992 sample basis for the methods used herein and assume full 96-well plates. Current costs of major
 993 reagents are the batch cost (per project).

Process	Reagent cost (\$US)	Alternatives and Potential Improvements
Isolating DNA	2	Speed-bead based DNA isolation
Normalizing DNA	0.2	Robots, Sequel-Prep or similar
Digesting DNA & Ligating Adapters	1	Reduce amount of enzymes used
Single cycle degenerate iTru5	0.2	
Size Selection (SNP/locus discovery)	0.1 [¶]	Speed-beads or gel-cut
NextSeq PE150	2 [¶]	HiSeq PE100 – PE150
Sequence Capture	1.2	-
NextSeq PE75	1.2	
RADcap genotyping total per sample	5.8	
Major Reagents		
MYcroarray MYbaits	2400*	Single-baits per locus, Increase number of projects pooled for baits
3RAD Adapters	370 [§]	Purchase aliquots or share
iTru Primers	345 [§]	Purchase aliquots or share

994 [¶]Excluded in the total genotyping cost per sample; *included or [§]excluded in per sample cost

995
996 Figure 1: Sequencing reads that can be obtained from full length 3RAD library molecules with
997 iTru5-8N sequence tags. The top double stranded molecule shows a 3RAD library molecule
998 prepared as described in the text. The color-scheme follows those of Glenn *et al.* (2016a; 2016b;
999 2016c) and Figure 2. The horizontal arrows above the text indicate positions on baits. The
1000 horizontal arrows beneath the library molecule indicate Illumina sequencing primers (binding to
1001 the complementary strand of the library molecules). The tip of the arrowhead indicates the 3'
1002 end of the primer and the direction of elongation for sequencing. Four sequencing reads are
1003 shown for each library prepared molecule, with one read for each index and each strand of the
1004 genomic DNA, including internal indexes. Reads are arranged 1 to 4 (numbered in magenta)
1005 from top to bottom, respectively. The arrow immediately 3' of the primers, indicates the data
1006 that are obtained from that primer.

1007



1008

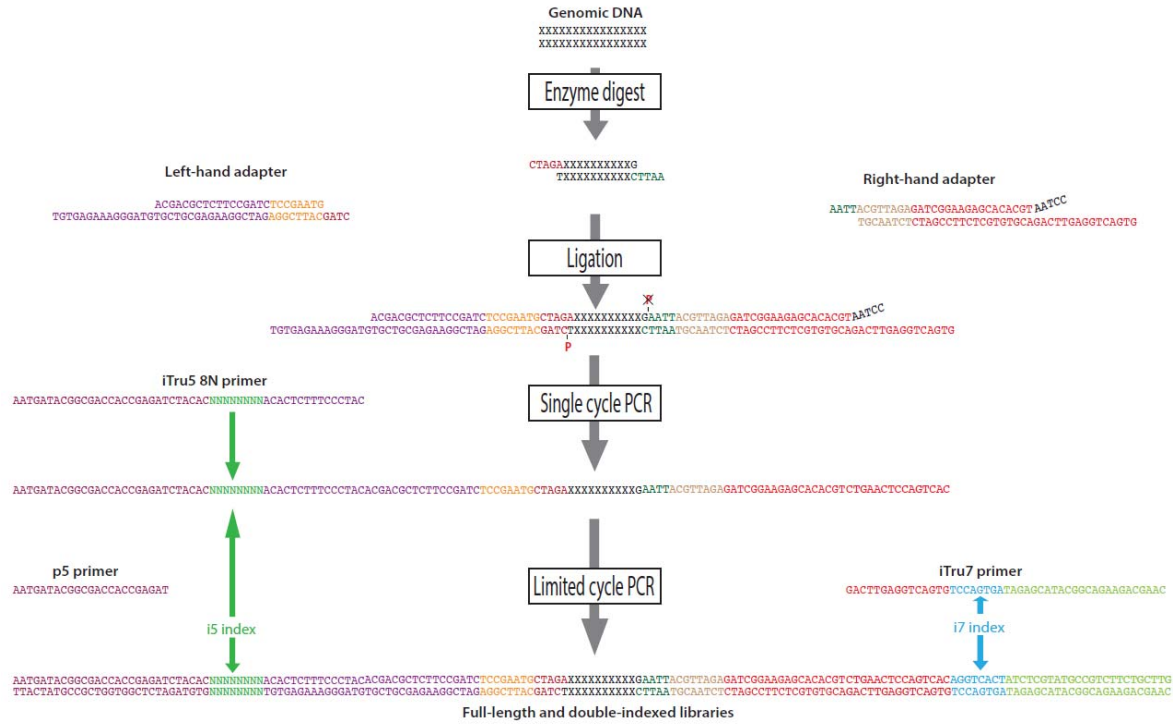
1009

1010

1011 Figure 2: The components of the library molecule added in different steps of the protocol and the
1012 sequence of the ends of the molecules. Genomic DNA is digested with enzymes that leave
1013 enzyme-specific sticky ends, to which we ligate adapters. The left hand adapter is comprised of
1014 four bases that bind to the XbaI restriction site overhang (dark red), a sample-specific internal
1015 sequence tag, used to identify the sample (orange) and a Read 1 sequencing primer that is
1016 partially single stranded to facilitate annealing of the iTru5 primer (purple). The right-hand
1017 adapter is a y-yoke adapter composed of the four bases that bind to the EcoRI restriction site
1018 overhang (dark green), a sample-specific internal sequence tag (tan), and the Read 2 sequencing
1019 primer (red). During the single cycle PCR, the iTru5 primer is added to the library: the partial
1020 library is denatured, the primer anneals to the Read 1 sequencing primer overhang (purple), and
1021 extends, thereby adding the degenerate barcode with 8 N bases (green), and the P5 primer
1022 (maroon) which anneals to the Illumina flowcell. After cleaning up the reaction, a limited cycle
1023 PCR is performed to add the iTru7 primer, comprised of the Read 2 sequencing primer (red)
1024 which anneals to the single stranded adapter added earlier, a sample-specific barcode (blue), and
1025 P7 primer (light green) which anneals to the Illumina flowcell.
1026

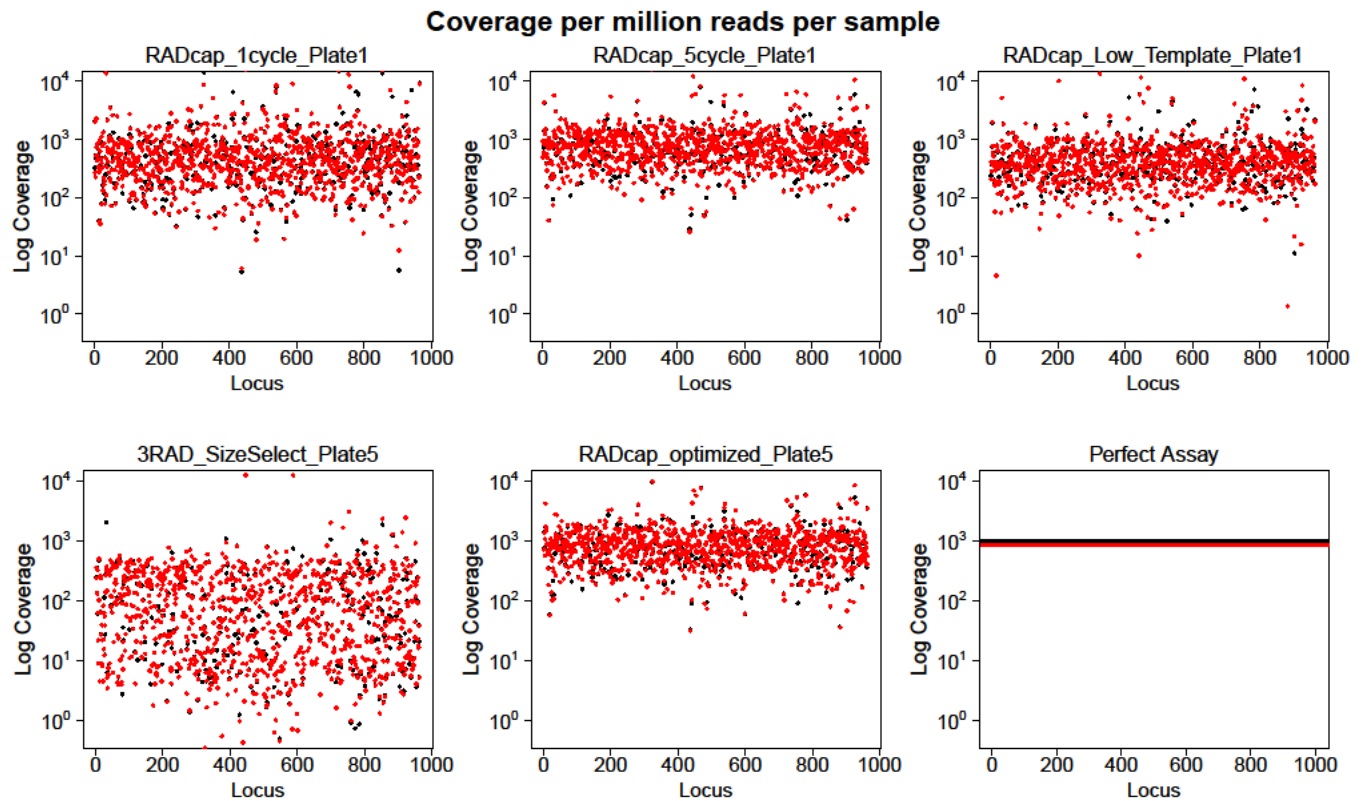
1027

3RAD Decloning-Index TruSeq GBS



1028

1029 Figure 3: The average coverage (log-transformed) per million reads per sample for 96 samples
1030 plotted on a log scale for the y-axis. Coverage for Read 1 is in black and Read 2 is in red. The
1031 bottom right panel represents a perfect assay, where all loci have the same coverage of 1000x per
1032 million reads per sample. The optimized treatment has a lower variance in coverage across loci
1033 than the size-selected or one-primer treatments ($p < 0.0083$), but variance in coverage was
1034 similar to the two-primer or low-template treatments ($p > 0.37$).

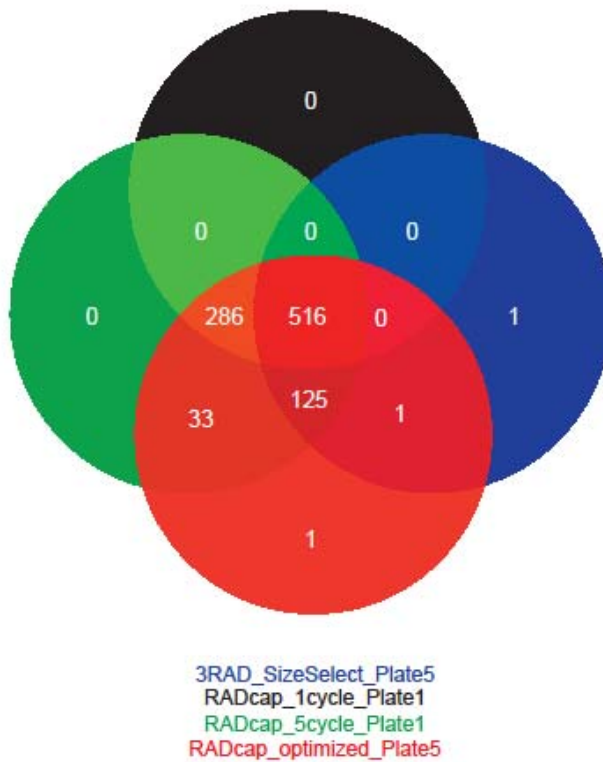


1035
1036
1037
1038
1039

1040

1041 Figure 4: Venn diagram of the number of loci with at least 4x coverage shared across at least
1042 90% (86) of samples in a single plate of the one-primer, two-primer, optimized, and size-selected
1043 treatments. Almost all of the 964 loci were recovered in the optimized and 5-cycle treatments,
1044 while the fewest loci were recovered in the size-selected treatment. Over half of the loci were
1045 shared among all four treatments, and only 3.7% of loci were found in one or two treatments.

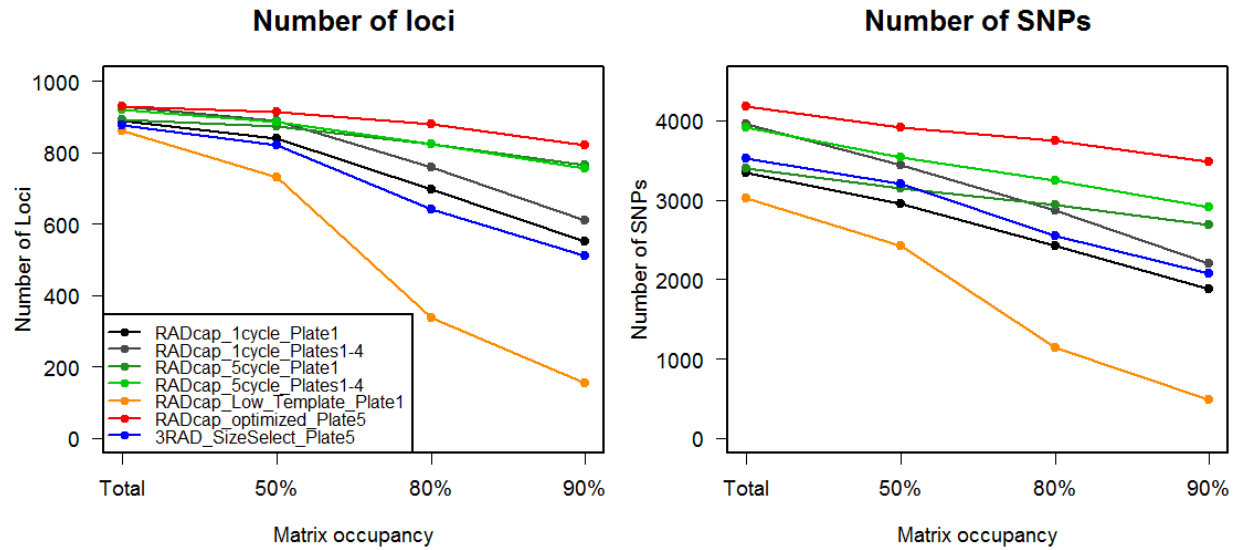
Number of loci with at least 4x coverage
in at least 90% of samples in four treatments



1046

1047

1048 Figure 5: The number of loci and SNPs retained at various levels of matrix occupancy for
1049 different treatments, analyzed with GATK. The number of loci and SNPs is the highest and most
1050 consistent in the optimized treatment.



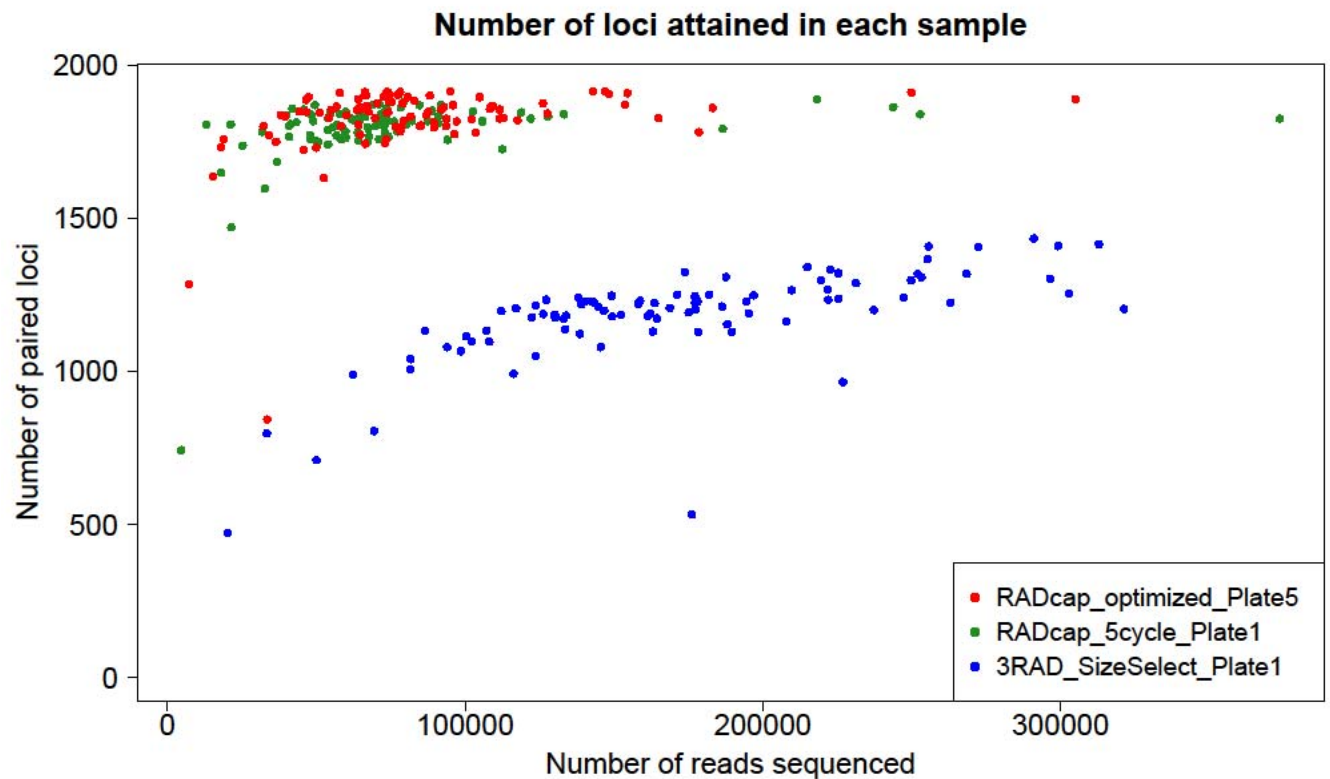
1051

1052

1053

1054

1055 Figure 6: A scatterplot of the number of paired loci sequenced to $\geq 4x$ coverage in all samples
1056 and the number of reads sequenced for each sample. All samples in plates
1057 RADcap_optimized_Plate5, RADcap_5cycle_Plate1, and 3RAD_SizeSelect_Plate1 had between
1058 approximately 20,000 and 350,000 reads. Most samples in the optimized and two-primer
1059 treatments with at least 50,000 reads had all paired loci, whereas most size-selected samples,
1060 even those with 300,000 reads, did not have all paired loci.



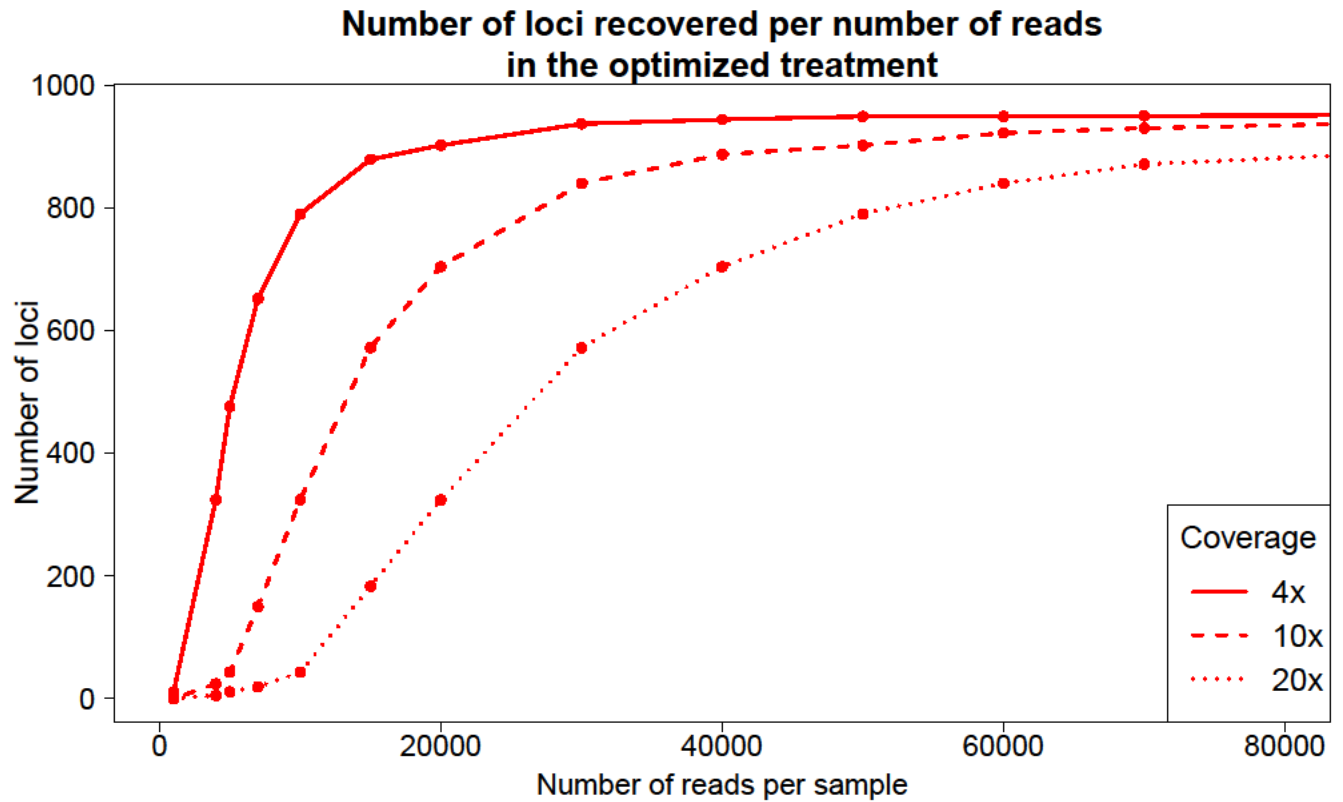
1061

1062

1063

1064

1065 Figure 7: The number of loci that should be recovered at various read depths for a minimum
1066 coverage of 4x, 10x, and 20x. For 4x coverage, 30,000 reads is enough to sequence all loci,
1067 whereas, 60,000 reads per sample is required for 10x coverage. For 20x coverage, over 80,000
1068 reads are required, and it may not be practical to sequence all loci at 20x coverage or higher.
1069



1070