

1 A population-level statistic for assessing Mendelian behavior  
2 of genotyping-by-sequencing data from highly duplicated  
3 genomes  
4

5 Lindsay V. Clark<sup>1\*</sup>, Wittney Mays<sup>23</sup>, Alexander E. Lipka<sup>2</sup>, and Erik J. Sacks<sup>2</sup>

6 <sup>1</sup>Roy J. Carver Biotechnology Center, University of Illinois at Urbana-Champaign, Urbana, IL  
7 61801, USA

8 <sup>2</sup>Department of Crop Sciences, University of Illinois at Urbana-Champaign, Urbana, IL 61801,  
9 USA

10 <sup>3</sup>Sandia National Laboratories, Livermore, CA 94551, USA

11 \*To whom correspondence should be addressed. Email: [lvclark@illinois.edu](mailto:lvclark@illinois.edu)

12 **Abstract**

13 **Background**

14 Given the economic and environmental importance of allopolyploids and other species with  
15 highly duplicated genomes, there is a need for methods to distinguish paralogs, i.e. duplicate  
16 sequences within a genome, from Mendelian loci, i.e. single copy sequences that pair at meiosis.  
17 The ratio of observed to expected heterozygosity is an effective tool for filtering loci but requires  
18 genotyping to be performed first at a high computational cost, whereas counting the number of  
19 sequence tags detected per genotype is computationally quick but very ineffective in inbred or  
20 polyploid populations. Therefore, new methods are needed for filtering paralogs.

## 21 **Results**

22 We introduce a novel statistic,  $H_{ind}/H_E$ , that uses the probability that two reads sampled from a  
23 genotype will belong to different alleles, instead of observed heterozygosity. The expected value  
24 of  $H_{ind}/H_E$  is the same across all loci in a dataset, regardless of read depth or allele frequency. In  
25 contrast to methods based on observed heterozygosity, it can be estimated and used for filtering  
26 loci prior to genotype calling. In addition to filtering paralogs, it can be used to filter loci with  
27 null alleles or high overdispersion, and identify individuals with unexpected ploidy and hybrid  
28 status. We demonstrate that the statistic is useful at read depths as low as five to 10, well below  
29 the depth needed for accurate genotype calling in polyploid and outcrossing species.

## 30 **Conclusions**

31 Our methodology for estimating  $H_{ind}/H_E$  across loci and individuals, as well as determining  
32 reasonable thresholds for filtering loci, is implemented in polyRAD v1.6, available at  
33 <https://github.com/lvclark/polyRAD>. In large sequencing datasets, we anticipate that the ability  
34 to filter markers and identify problematic individuals prior to genotype calling will save  
35 researchers considerable computational time.

## 36 **Keywords**

37 Polyploidy, single nucleotide polymorphism (SNP), heterozygosity, next generation DNA-  
38 sequencing (NGS), genome duplication

## 39 **Background**

40 Highly duplicated genome sequences are common throughout the plant kingdom. These include  
41 recent allopolyploids such as wheat, cotton, canola, strawberry, and coffee, as well as species  
42 with evidence of ancient whole genome duplication such as maize and legumes [1]. This

43 phenomenon is also present in the animal kingdom, for example allopolyploidy in the model frog  
44 *Xenopus*, as well as an ancient tetraploidization event followed by diploidization in salmonid  
45 fishes [2, 3]. For species in which paralogous sequences no longer pair at meiosis, accurate  
46 separation of paralogs in DNA and RNA sequence analysis, including reference genome  
47 assembly, remains challenging [4]. This separation of paralogs is especially important in variant  
48 calling, because SNPs and indels will not behave in a Mendelian fashion if the reads originate  
49 from more than one locus yet are erroneously attributed to a single locus [5]. Accurate variant  
50 calling therefore impacts all downstream analysis that assumes Mendelian inheritance, including  
51 linkage and QTL mapping, association studies, genomic selection, population genetics, and  
52 parentage analysis. For example, failure to remove paralogs from downstream analysis has been  
53 demonstrated to bias estimates of allele frequency and inbreeding as well as population structure  
54 [4, 6–8].

55 Due in part to the difficulty of assembling highly duplicated reference genomes, several methods  
56 have been published for filtering collapsed paralogous loci from genotyping-by-sequencing  
57 (GBS, including restriction-site associated DNA sequencing (RAD) approaches) datasets without  
58 the need for a reference genome. The most straightforward approach is to call genotypes and  
59 then determine if observed heterozygosity exceeds expected heterozygosity [9–11]. However,  
60 sampling error at low read depth can confound this filtering step by causing heterozygotes to be  
61 miscalled as homozygotes, lowering the observed heterozygosity. Moreover, estimating  
62 observed heterozygosity becomes complicated when polysomic inheritance is expected, due to  
63 the challenge of estimating allele copy number. Bayesian genotype calling methods mitigate the  
64 underestimation of observed heterozygosity, but at substantial computational cost [12–14].  
65 Another approach is to filter loci that have read depth above an arbitrary threshold [15], although

66 due to differences in amplification efficiency based on fragment size and GC content, this  
67 method could fail to filter some paralogs while filtering other non-paralogous loci. Peterson et  
68 al. [16] developed a method, extended by Willis et al. [17], that involved counting the number of  
69 unique haplotypes per individual for a putative locus, with the idea that in a collapsed paralog,  
70 the number of haplotypes would exceed the ploidy. However, this method can be confounded by  
71 sequencing error and cross-contamination among samples, and its sensitivity depends on allele  
72 frequencies, inbreeding, and ploidy. Other approaches have examined read depth ratios within  
73 individual genotypes [18] as well as read depth ratios in combination with observed  
74 heterozygosity [19]. Lastly, multiple methods identify putative paralogs based on networks of  
75 similarity among sequence tags [6, 20].

76 We present a novel statistic,  $H_{ind}/H_E$ , for evaluating marker quality, in particular for assessing  
77 whether a marker represents one Mendelian locus or multiple collapsed paralogous loci, based  
78 upon read depth distribution in a population. For a Mendelian locus, the statistic has the same  
79 expected value regardless of number of alleles, allele frequency, and total read depth. As a  
80 result, the distribution of the statistic can be visualized across loci in order to identify threshold  
81 values for filtering. The expected value can be calculated from ploidy (assuming disomic or  
82 polysomic inheritance) and the inbreeding coefficient, or the mode value of the statistic in a  
83 population can be used to estimate ploidy or inbreeding. Notably, because genotype calls are not  
84 needed in order to estimate this statistic, it can be used for filtering loci before any genotype  
85 calling is performed, saving computation time. Technical parameters such as sequencing error  
86 and overdispersion can influence estimates, but are explored here using simulated data so that  
87 they can be accounted for. We extend our Bayesian genotype calling software polyRAD [12] to  
88 implement the novel statistic and determine appropriate cutoffs.

## 89 **Results**

### 90 **The $H_{ind}$ statistic**

91 Here we describe a novel statistic,  $H_{ind}$ , that is based on sequence read depth across all alleles at  
92 a given locus and sample, and is agnostic of genotype calls, inheritance mode, and ploidy. It is  
93 related to observed heterozygosity,  $H_O$ , which in a diploid can be thought of as a matrix of ones  
94 and zeros indicating whether the genotype at each sample\*locus is heterozygous.  $H_{ind}$  is instead  
95 a number ranging from zero to one, indicating the probability that if two sequencing reads were  
96 sampled without replacement at that sample\*locus, they would represent different alleles. The  
97 abbreviation “ind” stands for “individual”, as it is calculated for each individual before averaging  
98 across a population. It can be calculated for SNP loci or for multiallelic haplotype- or tag-based  
99 loci, as long as allelic read depth is available.

100 The expected value for  $H_{ind}$  in a natural population of diploids or polysomic polyploids is:

101 Eqn. 1:  $\overline{H_{ind}} = \frac{k-1}{k} H_E (1 - F)$

102 where  $k$  is the ploidy,  $H_E$  is the expected heterozygosity at the same locus, and  $F$  is the  
103 inbreeding coefficient.  $H_E$  is the probability that two alleles drawn at random from the  
104 population will be different,  $(1 - F)$  is the probability that two alleles randomly drawn from an  
105 individual will not be identical by descent, and  $(k - 1)/k$  is the probability that two sequencing  
106 reads originated from different chromosome copies. Multiplied together, these three terms yield  
107 the probability that two sequence reads from one sample at one locus will be different from each  
108 other.

109 If we divide  $H_{ind}$  by  $H_E$ :

110 Eqn. 2:  $\overline{H_{ind}}/H_E = \frac{k-1}{k} (1 - F)$

111 we now have a statistic that is only dependent on ploidy and inbreeding, two parameters that we  
112 will assume to be consistent across samples and loci.

113 In a mapping population, the term  $H_E * (1 - F)$  must be replaced by the probability, for a given  
114 locus, that two locus copies in a progeny will be different alleles. This requires knowledge of the  
115 ploidy, parental genotypes, and population design including number of generations of  
116 backcrossing and self-fertilization. This probability, which we will call  $H_{E.map}$ , can be estimated  
117 by simulation of the cross. The expectation is then:

118 Eqn. 3:  $\overline{H_{ind}}/H_{E.map} = \frac{k-1}{k}$

119 Common factors that influence  $H_{ind}/H_E$  are listed in Table 1, and explored in subsequent sections.

120

121 **Table 1. Biological and technical parameters that influence the expected value and**  
122 **variance of  $H_{ind}/H_E$ .**

Parameter	Effect
Ploidy	Expected value increases with ploidy.
Inbreeding (including population structure)	Expected value decreases as inbreeding increases.
Hybridization	Value increases with increase in heterozygosity from hybridization across species or divergent populations.
Paralogy	Value increases if multiple loci are collapsed into one.
Overdispersion	Expected value decreases as read depth ratios deviate further from allelic dosage.
Sequencing error	Value is biased upward by sequencing error, especially at low minor allele frequencies.
Null alleles (e.g. restriction site polymorphisms, deletions)	Expected value decreases with increasing null allele frequency.
Minor allele frequency	Variance decreases at increased minor allele frequency. Overdispersion, sequencing error, or very low read depth in combination with low minor allele frequency bias the value upward.
Sample size	Variance decreases at increased sample size.
Number of alleles	Multiallelic loci have lower variance than biallelic SNPs.
Read depth	Low read depth loci tend to have low values due to the presence of null alleles. High read depth loci tend to have high values due to paralogy. Genome-wide increases in read depth (e.g. larger library size) reduce variance in the statistic, as well as reducing upward bias at low minor allele frequencies.
Polyploid mapping populations	Variance is lower at markers with higher heterozygosity in the progeny.

123

124 **Empirical estimation of  $H_{ind}/H_E$**

125 Say that we have sequence read depths,  $\{d_{1m} \dots d_{jm}\}$ , across a set of  $j$  alleles at a single locus in  
126 an individual  $m$ . Total read depth in one individual is

127 Eqn. 4:  $D_m = \sum_{i=1}^j d_{im}$

128 As long as there are two or more reads, we can estimate  $H_{ind}$  within that individual using the  
129 Gini-Simpson index [21]:

130 Eqn. 5:  $\hat{H}_{ind,m} = \left(1 - \sum_{i=1}^j \left(\frac{d_{im}}{D_m}\right)^2\right) \frac{D_m}{D_m - 1}$

131 For a population of  $n$  individuals with sequencing reads, allele frequencies are estimated from  
132 average within-individual read depth ratios:

133 Eqn. 6:  $\hat{p}_i = \frac{\sum_{m=1}^n \frac{d_{im}}{D_m}}{n}$

134 And expected heterozygosity is estimated as

135 Eqn. 7:  $\hat{H}_E = 1 - \sum_{i=1}^j \hat{p}_i^2$

136 Averaged across  $n$  individuals with two or more reads at a given locus in a natural population,  
137 the expectation is then:

138 Eqn. 8:  $\hat{H}_{ind}/\hat{H}_E = \frac{\sum_{m=1}^n \hat{H}_{ind,m}/\hat{H}_E}{n} \cong \frac{k-1}{k} (1 - F)$



139 In a mapping population,  $\hat{H}_{ind,m}$  is estimated in the same way.  $H_{E,map}$  is estimated from parental  
140 genotypes and population design, and the expected average ratio within a locus is given in Eqn.  
141 3.

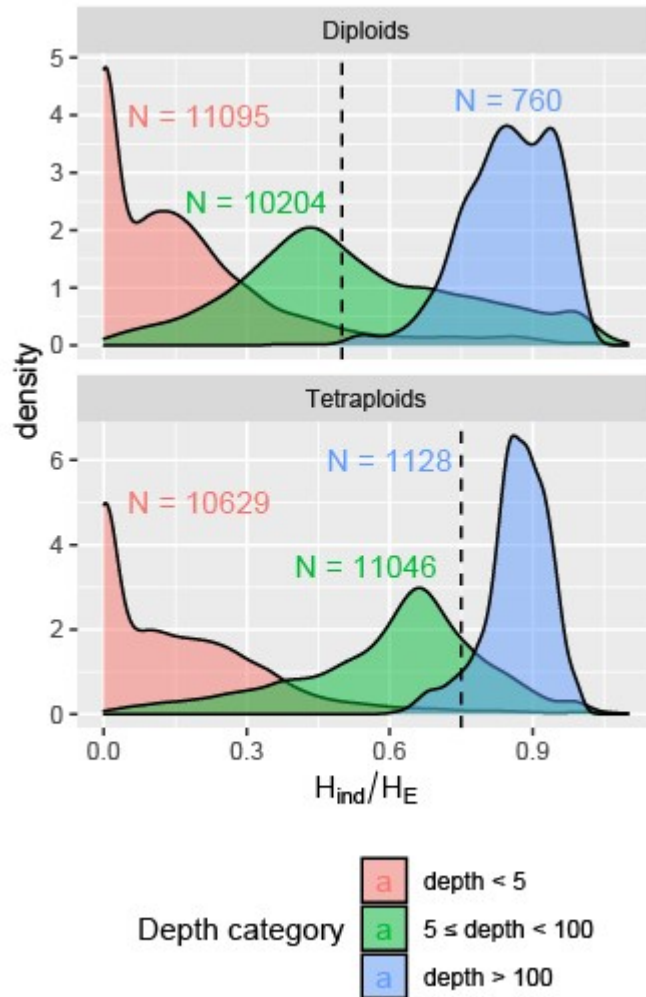
## 142 **Utility of $H_{ind}/H_E$ for detecting collapsed paralogs in a diversity panel**

143 To compare the distribution of  $H_{ind}/H_E$  values for Mendelian loci versus collapsed paralogs, we  
144 aligned *M. sacchariflorus* tag sequences to the *M. sinensis* reference genome, in which they  
145 should align to the correct paralog most of the time, and to the *S. bicolor* reference genome, in  
146 which two paralogs from *Miscanthus* correspond to one alignment location. We found that loci  
147 with a mean read depth less than five had very low estimates of  $H_{ind}/H_E$ , likely due to restriction  
148 site polymorphisms or other technical issues (Fig. 1 and Additional File 1: Fig. S1). As mean  
149 read depth increased above 100 in our dataset, however, loci tended to have  $H_{ind}/H_E$  values above  
150 the expectation for a Mendelian locus, suggesting that most loci at this depth and higher were in  
151 fact collapsed paralogs (Fig. 1, and Additional File 1: Figs. S1 and S2).

152 When a mean depth of five was used as a cutoff and the *M. sinensis* genome was used as a  
153 reference, the peak  $H_{ind}/H_E$  value was slightly below the expected values of 0.5 for diploids and  
154 0.75 for tetraploids (Fig. 2), indicating some inbreeding, likely due to population structure [22].  
155 A second peak was observed at a higher value of  $H_{ind}/H_E$  (Fig. 2), likely representing sets of tags  
156 that belonged to different Mendelian loci despite aligning to the same location (i.e.  
157 misalignments). When *S. bicolor* was used as the reference genome, the opposite trend was  
158 observed, where most loci had a  $H_{ind}/H_E$  above the expected value, indicating collapsed paralogs,  
159 but a second peak was observed closer to the expected value, indicating regions in the *S. bicolor*  
160 genome that may only have synteny with one region of the *M. sinensis* genome (Fig. 2).

161 Although the peaks overlapped somewhat, they were distinct enough that a reasonable threshold  
162 for identifying putative collapsed paralogs could be visually determined (Fig. 2). Moreover,  
163 although the diploid and tetraploid datasets were processed separately, they were largely in  
164 agreement about which loci were Mendelian and which were collapsed paralogs (Additional File  
165 1: Fig. S2), suggesting that the filtering performed in one population can be applied to another  
166 population, which could be especially useful for populations that are too small for accurate  
167 estimation of  $H_{ind}/H_E$ .

168 In both the diploid and tetraploid datasets, the distribution and peak values of  $H_{ind}/H_E$  were  
169 similar regardless of whether biallelic SNPs or multiallelic, haplotype-based markers were used  
170 (Additional File 1: Fig. S3). However, the variance of  $H_{ind}/H_E$  was approximately 20% higher  
171 when SNPs were used, suggesting that the higher information content of multiallelic markers  
172 improves the precision of  $H_{ind}/H_E$  estimates.



173

174 **Figure 1. Relationship between  $H_{ind}/H_E$  statistic and mean sequence read depth per locus.**

175 Loci were called across 356 diploid and 268 tetraploid *Miscanthus sacchariflorus* based on

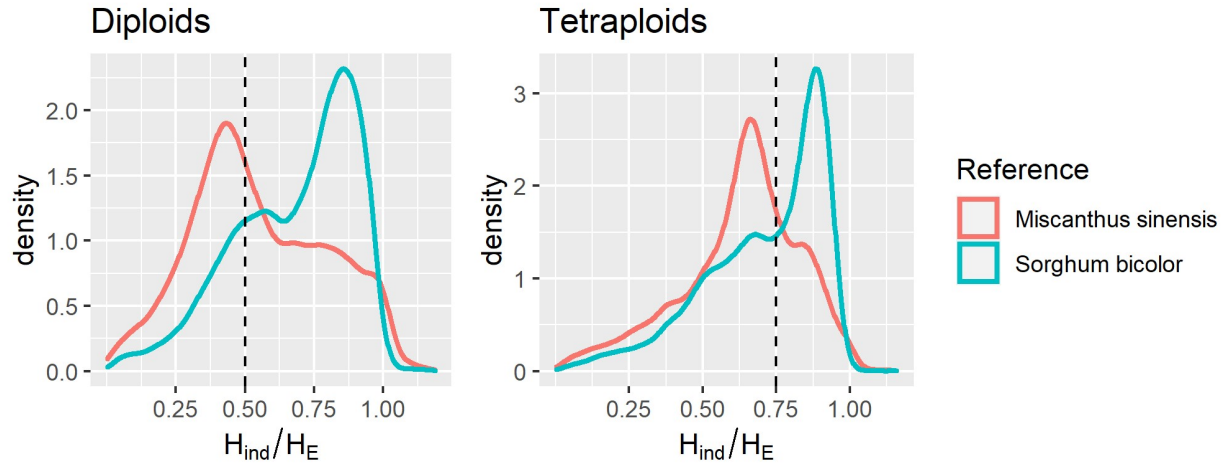
176 alignments to the *M. sinensis* reference genome. The number of loci in each depth category is

177 indicated. Fig. S1 in Additional File 1 provides justification for the depth thresholds for

178 categories. The expected value for a Mendelian locus in Hardy-Weinberg equilibrium is shown

179 with a dashed line.

180



181

182 **Figure 2. Effect of reference genome and ploidy on  $H_{ind}/H_E$  per locus in *Miscanthus***

183 *sacchariflorus*. Loci with a mean read depth below five were omitted, leaving 11,516 loci

184 aligned to the *M. sinensis* reference and 8,820 loci aligned to the *Sorghum bicolor* reference.

185 Expected values for Mendelian loci under Hardy-Weinberg equilibrium are shown with dashed

186 lines.

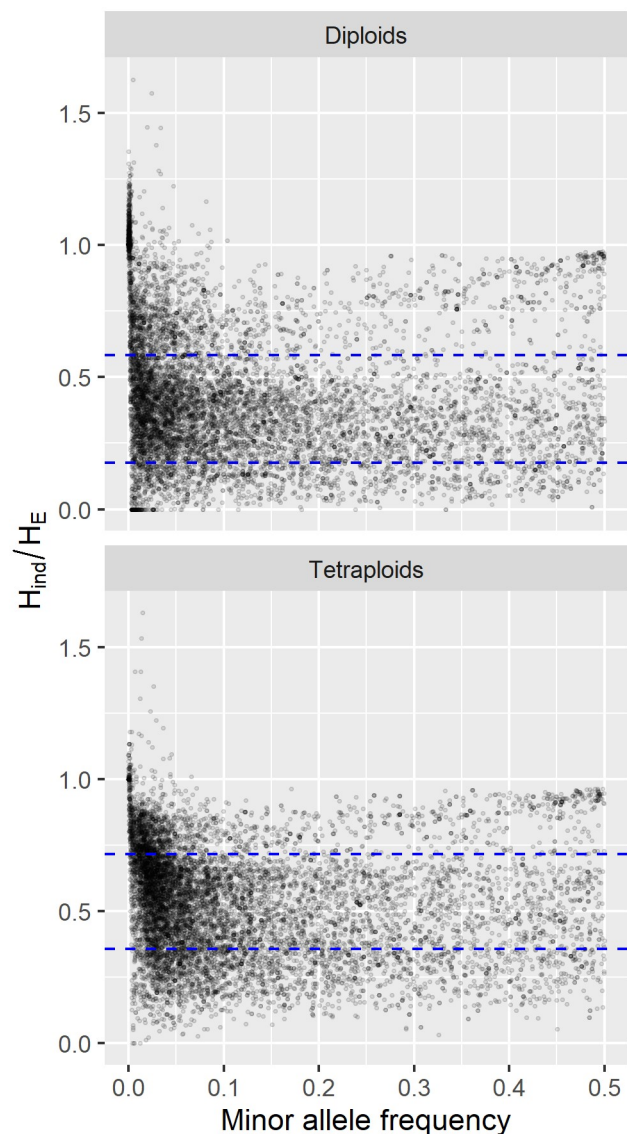
187

188 The *ExpectedHindHe* function in polyRAD was used to set thresholds for filtering the diploid  
189 and tetraploid datasets. Based on results from the *TestOverdispersion* function, the  
190 overdispersion parameter was set to 11 for diploids and 10 for tetraploids. Based on the  
191 observed distribution of  $H_{ind}/H_E$  in the dataset, the inbreeding coefficient was set to 0.35 for  
192 diploids and 0.25 for tetraploids. Based on these parameters, as well as read depth and allele  
193 frequencies in the datasets, the ranges for retaining 95% of Mendelian loci were 0.175 to 0.584  
194 in diploids and 0.356 to 0.716 in tetraploids as estimated by *ExpectedHindHe*, resulting in 40.2%  
195 and 42.3% of loci being filtered, respectively (Table 2). Markers within genes were  
196 underrepresented among markers that were filtered for having  $H_{ind}/H_E$  below the lower threshold,  
197 and overrepresented among markers that were filtered for having  $H_{ind}/H_E$  above the upper  
198 threshold, significant in Fisher's Exact Test at  $P < 0.0005$  (Table 2). Markers that were filtered  
199 having  $H_{ind}/H_E$  above the upper threshold tended to have minor allele frequencies that were very  
200 low, consistent with the markers representing sequencing error rather than true alleles, or very  
201 high, consistent with the markers representing collapsed paralogs (Fig. 3).

202 **Table 2. Contingency tables of number of markers retained and filtered for being above or**  
203 **below  $H_{ind}/H_E$  thresholds in *Miscanthus sacchariflorus*, by whether or not the marker was**  
204 **within a gene.**

	Diploids		Tetraploids	
	In a gene	Not in a gene	In a gene	Not in a gene
Filtered; too low	337	950	588	1727
Retained	2201	3654	2419	3500
Filtered; too high	1361	1287	1091	930

205



206

207 **Figure 3. Filtering by  $H_{ind}/H_E$  vs. minor allele frequency in *Miscanthus sacchariflorus*. A**

208 dataset of 10,458 SNP loci was tested across 356 diploid and 268 tetraploid individuals. Blue

209 dashed lines indicate filtering thresholds to retain 95% of Mendelian loci based on simulated

210 distributions.

211 **By individual,  $H_{ind}/H_E$  reflects ploidy and hybrid status**

212 In addition to evaluating the mean  $H_{ind}/H_E$  within loci, we also obtained the mean statistic within

213 individuals in order to assess the utility of the statistic for determining ploidy. We found that

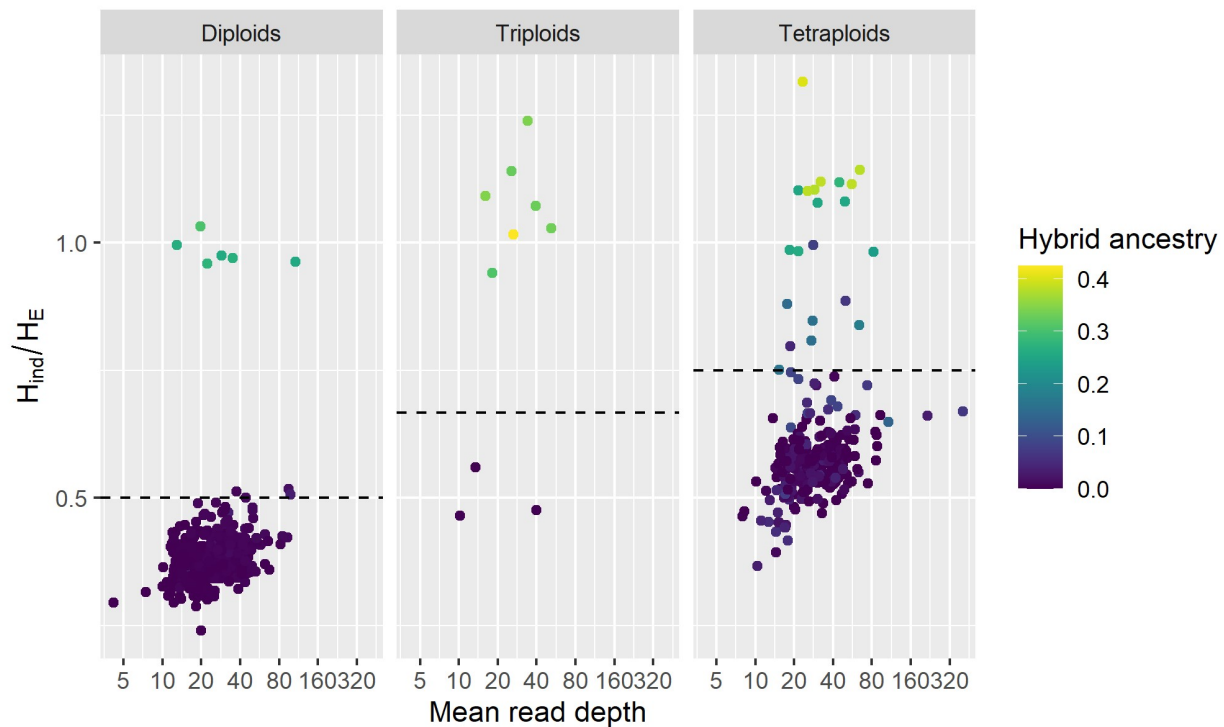
214  $H_{ind}/H_E$  increased with ploidy, largely independent of read depth (Fig. 4). Although the  
215 distributions overlapped too much for  $H_{ind}/H_E$  to be a conclusive indicator of ploidy, it could still  
216 potentially be used to identify outlier individuals whose ploidy should be confirmed by other  
217 means (e.g. flow cytometry). Additionally, because our empirical dataset included many natural  
218 interspecific (*M. sacchariflorus* × *M. sinensis*) F1 hybrid and backcross individuals, we were  
219 also able to observe that  $H_{ind}/H_E$  values were considerably higher in hybrids than in non-hybrids,  
220 reflecting higher heterozygosity (Fig. 4).

### 221 **Variance and bias in the $H_{ind}/H_E$ statistic using simulated data**

222 Using simulated data resembling a diversity panel or natural population, the mean  $H_{ind}/H_E$   
223 estimate decreased as inbreeding increased, with diploid and tetraploid loci being  
224 indistinguishable at an inbreeding coefficient of 0.8 or higher (Fig. 5). Sequencing error had  
225 little effect on the estimate at a minor allele frequency of 0.05, but caused an inflated estimate at  
226 a minor allele frequency of 0.01, particularly as inbreeding increased (Fig. 5). Variance and bias  
227 in the statistic were minimized if there were at least 500 samples, minor allele frequency was  
228 0.05 or higher, and read depth was at least 5 (Fig. 6). Ploidy had negligible impact on variance  
229 and bias (Fig. 6). Read depth and minor allele frequency influenced the estimates for collapsed  
230 paralogs, but not enough to interfere with distinguishing them from Mendelian markers (Fig. 6).  
231 As expected, overdispersion (deviation of read depth ratios from allelic dosage ratios) reduced  
232 the mean  $H_{ind}/H_E$  estimate, with the effect of overdispersion being greater at higher minor allele  
233 frequencies (Additional File 1: Fig. S4). The  $H_{ind}/H_E$  estimate also decreased linearly as null  
234 allele frequency increased (Additional File 1: Fig. S5).

235 In simulated F1 mapping populations, the standard deviation of the  $H_{ind}/H_E$  ranged from 0.012 to  
236 0.076 depending on the marker type (Fig. 7). In tetraploids, marker types with high expected

237 heterozygosity in the progeny, such as triplex x nulliplex and triplex x simplex, had lower  
238 variance in the estimate than marker types with lower expected heterozygosity in the progeny,  
239 such as simplex x nulliplex and simplex x simplex (Fig. 7). A few rare markers had  $H_{ind}/H_E$   
240 estimates that deviated very far from the expected value, indicating that the parents were  
241 incorrectly genotyped (Fig. 7).

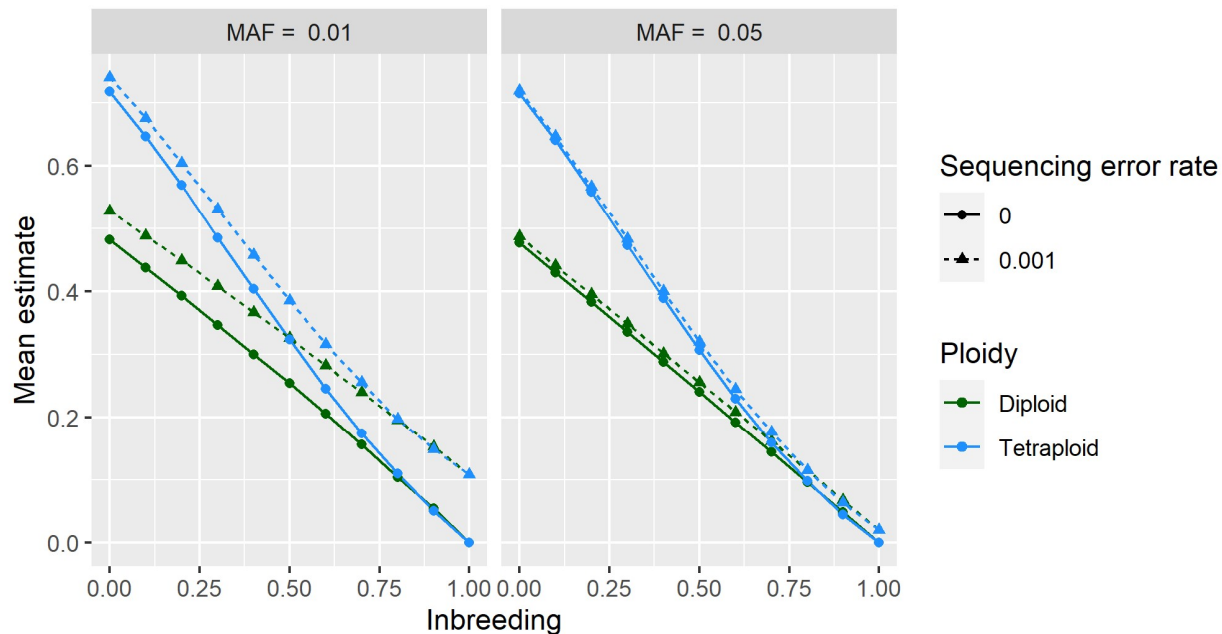


242

243 **Figure 4. Relationship between ploidy, sequence read depth, hybrid ancestry, and  $H_{ind}/H_E$**   
244 **among 620 *M. sacchariflorus* individuals.** Ploidy and proportion of ancestry from *M. sinensis*  
245 (hybrid ancestry) were determined previously [22]. Read depth and  $H_{ind}/H_E$  were averaged  
246 across 10,000 loci. The expected value for  $H_{ind}/H_E$  under Hardy-Weinberg equilibrium is shown  
247 with the dashed line.

248

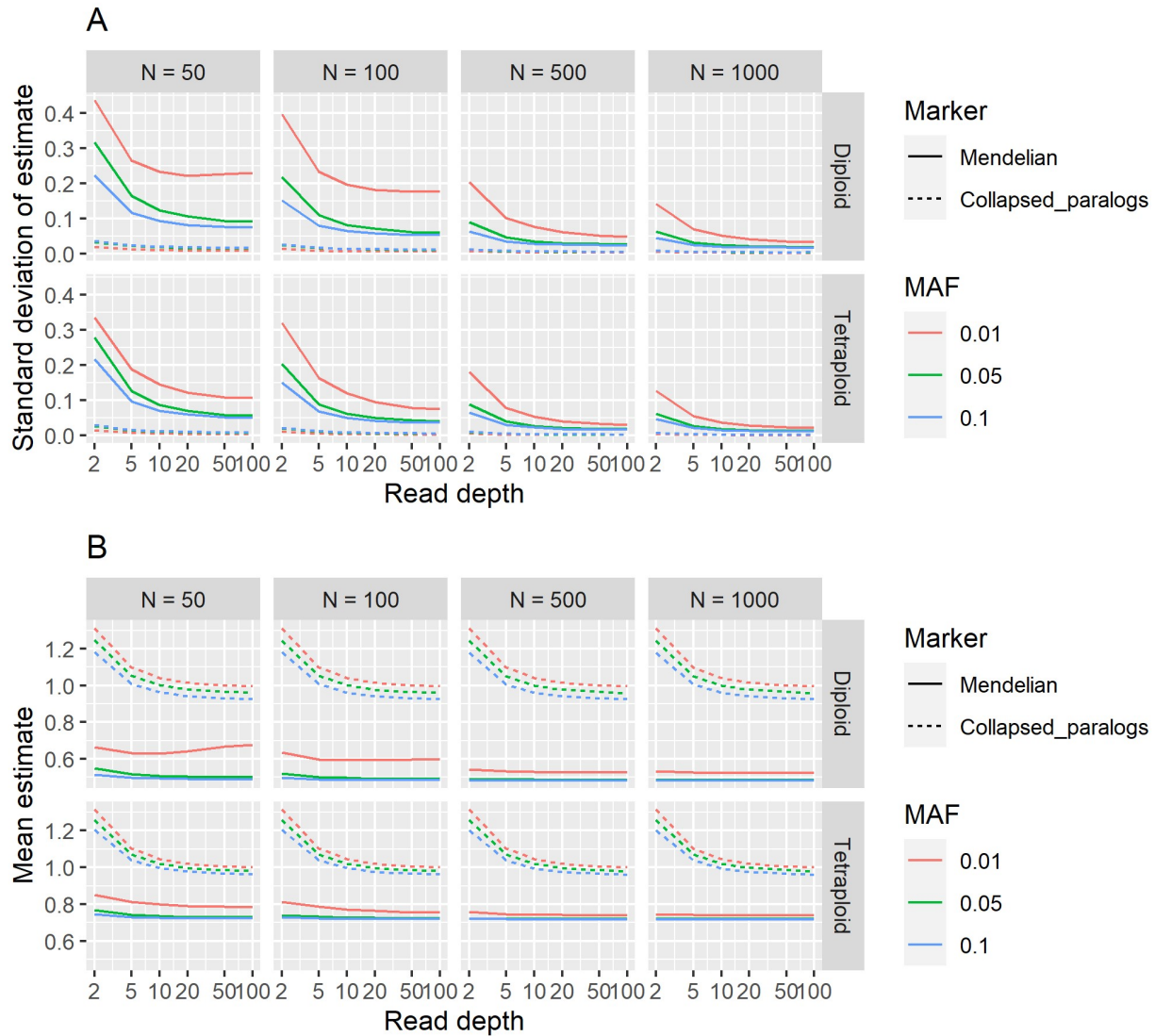




249

250 **Figure 5. Combined effects of inbreeding, ploidy, minor allele frequency (MAF), and**  
251 **sequencing error on mean estimates of  $H_{ind}/H_E$  using simulated data.** At each combination  
252 of parameters, 20,000 biallelic loci were simulated with a read depth of 20 and overdispersion  
253 parameter of 20. The x-axis indicates the inbreeding coefficient (the probability that two alleles  
254 in an individual are identical by descent) while the y-axis indicates the  $H_{ind}/H_E$  estimate.

255



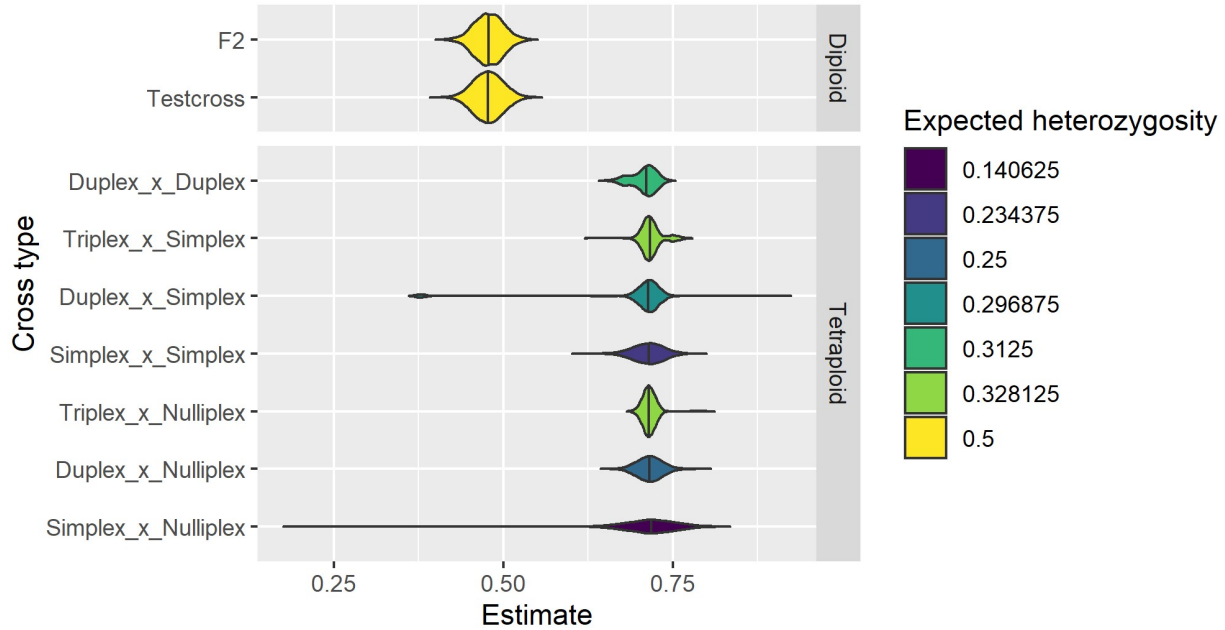
256

257 **Figure 6. Effect of sample size, read depth, and minor allele frequency on variance and**  
 258 **bias of estimates of  $H_{ind}/H_E$ .** For each combination of ploidy, sample size (N), read depth, and  
 259 minor allele frequency (MAF), 5000 biallelic Mendelian loci were simulated under Hardy-  
 260 Weinberg Equilibrium with an overdispersion parameter of 20 and sequencing error rate of  
 261 0.001. Additionally, 5000 collapsed paralogs, each consisting of two Mendelian loci, were  
 262 simulated under each set of the same parameters. (A) Standard deviation of  $H_{ind}/H_E$  estimates.

263 (B) Mean  $H_{ind}/H_E$  estimates. Expected values are 0.5 for diploids and 0.75 for tetraploids;  
264 deviations from these values indicate bias in estimation.

265

266



267

268 **Figure 7. Distribution of  $H_{ind}/H_E$  estimates in simulated F1 mapping populations.** For each  
269 cross type, 5000 biallelic loci with a read depth of 20, overdispersion parameter of 20, and  
270 sequencing error rate of 0.001 were simulated across 500 individuals.

271

## 272 **Comparison with other approaches**

273 To compare effectiveness at filtering paralogs between  $H_{ind}/H_E$  and other approaches, 1000  
274 Mendelian loci and 1000 collapsed paralogs were simulated in 200 diploid and 200 tetraploid  
275 individuals at three levels of inbreeding. The median allele frequency was 0.026 and median  
276 read depth per Mendelian locus was 21. For each statistic, the 95<sup>th</sup> percentile for Mendelian loci  
277 was determined, and the proportion of collapsed paralogs that would be filtered at that threshold  
278 was estimated. The  $H_{ind}/H_E$  approach and observed over expected heterozygosity ( $H_O/H_E$ )  
279 performed best, with  $H_O/H_E$  having the disadvantage that genotyping must be performed before it  
280 can be estimated, thus increasing processing time two orders of magnitude over  $H_{ind}/H_E$  (Table  
281 3). The  $H_{ind}/H_E$  thresholds used for filtering were 0.58, 0.41, and 0.17 in diploids and 0.76, 0.48,  
282 and 0.17 in tetraploids at inbreeding levels of 0.1, 0.5, and 0.9, respectively. The haplotype  
283 counting approach [17] and allelic depth ratio Z-score approach [19] both performed reasonably  
284 well in diploids but were much less effective in tetraploids, with haplotype counting being  
285 useless in tetraploids at high inbreeding, while the Z-score approach additionally suffered in  
286 terms of computational time due to the need for genotyping. However, haplotype counting used  
287 2- to 3-fold less computational time than  $H_{ind}/H_E$ , and thus could be advantageous in diploids  
288 when millions of loci are being processed. Lastly, filtering on read depth alone was not very  
289 effective given the variation in read depth among loci.

290

291 **Table 3. Effectiveness of various statistics for identifying paralogs, using simulated data**  
 292 **across three levels of inbreeding. Standard error is shown for proportion paralogs filtered.**

Statistic	Ploidy	Proportion paralogs filtered			Median processing time (s / 1000 loci)
		$F = 0.1$	$F = 0.5$	$F = 0.9$	
$H_{ind}/H_E$	Diploid	$0.988 \pm 0.003$	$0.998 \pm 0.001$	$1.000 \pm 0.000$	0.17
		$0.905 \pm 0.009$	$0.994 \pm 0.002$	$1.000 \pm 0.000$	
	Tetraploid	$0.009$	$0.002$	$0.000$	0.16
$H_O/H_E$	Diploid	$0.993 \pm 0.003$	$0.989 \pm 0.003$	$0.997 \pm 0.002$	17.90
		$0.976 \pm 0.005$	$0.935 \pm 0.008$	$0.897 \pm 0.010$	
	Tetraploid	$0.005$	$0.008$	$0.010$	50.42
Proportion individuals with more haplotypes than expected	Diploid	$0.985 \pm 0.004$	$0.982 \pm 0.004$	$0.948 \pm 0.007$	0.07
	Tetraploid	$0.564 \pm 0.016$	$0.163 \pm 0.012$	$0.001 \pm 0.001$	0.05
Absolute value of Z-score for read depth ratio	Diploid	$0.878 \pm 0.010$	$0.888 \pm 0.010$	$0.851 \pm 0.011$	18.00
	Tetraploid	$0.642 \pm 0.015$	$0.542 \pm 0.016$	$0.511 \pm 0.016$	52.15
Mean read depth	Both	$0.396 \pm 0.015$			0.00

293

## 294 **Discussion**

### 295 **Properties of the $H_{ind}/H_E$ statistic**

296 While  $H_{ind}/H_E$  can be used, in combination with other metrics, to assess locus quality, this should  
 297 be performed with an understanding of what biological and technical phenomena can cause it to  
 298 deviate from the expected value. Inbreeding from any source will lower the expected value  
 299 below  $(k - 1)/k$ , where  $k$  is the ploidy; this includes not only self-fertilization and preferential  
 300 mating with relatives, but also population structure, which is why we observed values below  $(k -$

301  $1/k$  even in self-incompatible, wind-pollinated *M. sacchariflorus* (Figs. 1-4). A benefit of this,  
302 however, is that as long as ploidy is known and overdispersion can be reasonably estimated (e.g.  
303 with the *TestOverdispersion* function in polyRAD),  $H_{ind}/H_E$  can be used to estimate inbreeding,  
304 either at the population or individual level, directly from sequence read depth. Given that we  
305 observed  $H_{ind}/H_E$  to be inflated at low minor allele frequencies, we recommend using the mode  
306  $H_{ind}/H_E$  at markers with minor allele frequency of at least 0.05 for estimating inbreeding.  
307 Additionally, individuals that are hybrids between species or between highly diverged  
308 populations, as well as DNA samples that are an accidental mix of two or more individuals, may  
309 have  $H_{ind}/H_E$  above the expected value (Fig. 4). Strong selection for homozygotes or  
310 heterozygotes at particular loci would be expected to lower and raise  $H_{ind}/H_E$ , respectively.

311 At the locus level, a  $H_{ind}/H_E$  that exceeds the expected value can be an indication that alleles are  
312 derived from paralogous loci rather than a true Mendelian locus. More broadly, if all alleles  
313 truly belong to a single locus, then the expected value is  $(1 - F)(k - 1)/k$ . However, if a set of  
314 random, independent alleles were assigned to one putative locus, the expected value of  $H_{ind}/H_E$   
315 would be one, because the probability of sampling reads from two different alleles within one  
316 individual would be the same as the probability of sampling reads from two different alleles in  
317 the general population. In the *M. sacchariflorus* dataset, markers within genes were  
318 overrepresented among markers that were filtered for having  $H_{ind}/H_E$  above the expected value,  
319 likely due to high sequence conservation between paralogs (Table 2). A  $H_{ind}/H_E$  of zero could  
320 indicate a cytoplasmic marker, because while there may be variation in the population, each  
321 individual would only be expected to possess reads from one allele. Loci with highly  
322 overdispersed read depth distributions due to technical issues such as differential fragment size  
323 or variation in library preparation would also be expected to have  $H_{ind}/H_E$  below expectations; it

324 may be advantageous to filter these from the dataset as they will tend to yield poor-quality  
325 genotype calls. Lastly, loci with common null alleles have lower than expected  $H_{ind}/H_E$   
326 (Additional File 1: Fig. S3), resulting in a tendency to filter loci that are not within genes as these  
327 regions are less conserved (Table 2). Null alleles can be the result of restriction cut site  
328 polymorphism in RAD-based techniques, primer binding site mutations in amplicon sequencing,  
329 or deletion mutations using any genotyping method. Because they are a common problem,  
330  $H_{ind}/H_E$  can be used to identify and filter loci with null alleles.

331 The expected value of  $H_{ind}/H_E$  is independent of read depth, number of individuals sampled, and  
332 the allele frequency. However, all of these factors influence the variance of the estimate, and  
333 low minor allele frequency especially can bias it upwards (Fig. 5-6). As there is no generalized  
334 formula to estimate the variance of a ratio, the variance of  $H_{ind}/H_E$  cannot be estimated  
335 mathematically. Moreover, sequencing error inflates the estimate at low minor allele frequency  
336 (Fig. 5), and polyRAD cannot account for sequence quality scores or alignment quality scores  
337 since it only imports allelic read depth. We therefore recommend simulating data for Mendelian  
338 loci given the ploidy, inbreeding, sample size, sequencing error rate, and distribution of read  
339 depth and allele frequency observed in the dataset of interest. The distribution of  $H_{ind}/H_E$  across  
340 simulated loci then can be used to determine cutoff values for filtering loci in the empirical  
341 dataset. The *ExpectedHindHe* and *ExpecteHindHeMapping* functions are available in polyRAD  
342 for this purpose, and suggest cutoffs for filtering loci in order to retain 95% of Mendelian loci.  
343 Depending on the downstream application, we recommend considering the number of markers  
344 needed versus the importance of marker quality when determining thresholds for read depth,  
345 allele frequency, and  $H_{ind}/H_E$ .

346  $H_{ind}/H_E$  is more useful for detecting paralogs when haplotypes are treated as alleles (i.e. loci can  
347 be multiallelic), as opposed to when all loci are treated as biallelic SNPs, simply due to the fact  
348 that multiallelic markers are more information-rich than biallelic markers for the same  
349 distribution of minor allele frequencies. We observed that, for the same set of SNPs in *M.*  
350 *sacchariflorus*, the median value of  $H_{ind}/H_E$  per locus was very similar regardless of whether they  
351 were phased into haplotypes within the span of a single RAD tag, but the variance in  $H_{ind}/H_E$  was  
352 about 20% higher for SNPs vs. haplotypes (Additional File 1: Fig. S3). This improved power  
353 and information content is why polyRAD generally imports multiallelic, haplotype-based  
354 genotypes rather than SNPs as the default. Other methods for marker calling in highly  
355 duplicated genomes have also benefitted from the use of haplotype information [11, 23], and  
356 multiallelic markers have been found to be advantageous over biallelic SNPs for linkage  
357 mapping in polyploids [24]. It should be noted that in this study we only phased SNPs that were  
358 certain to have originated from the same sequencing reads based on physical linkage and read  
359 depth. The  $H_{ind}/H_E$  statistic cannot be estimated using haplotypes spanning longer distances,  
360 given that read depth will vary from locus to locus within haplotype.

### 361 **Uses of the $H_{ind}/H_E$ statistic**

362 We anticipate locus-filtering to be the most common application of the  $H_{ind}/H_E$  statistic, with  
363 major advantages being that it is not biased by read depth or allele frequency and can be  
364 estimated prior to genotype calling. We demonstrate that it is similar to  $H_O/H_E$  in effectiveness  
365 for filtering paralogs, with substantial savings on computational time (Table 2). We should note  
366 that our  $H_O/H_E$  estimates used Bayesian genotype calls from polyRAD, which mitigate the  
367 underestimation of observed heterozygosity as compared to naïve genotype calls [12].  
368 Stringency of filtering should depend on the genotype quality needed for downstream analysis;



369 for example, parentage analysis and QTL mapping are sensitive to genotyping errors, whereas  
370 genome-wide association studies and estimations of population structure from principal  
371 components analysis are less sensitive. Missing data rate, median read depth, and minor allele  
372 frequency are common criteria that should be used in combination with  $H_{ind}/H_E$  to determine  
373 which loci to retain for downstream analysis. In our empirical dataset, we found the loci ranging  
374 in depth from five to 100 had the best distribution of  $H_{ind}/H_E$  (Fig. 1), but a higher minimum  
375 depth may be required for applications that require accurate genotype calling, and the optimal  
376 maximum depth used in filtering depends on the overall depth of the dataset. The use of  
377 observed heterozygosity, read depth ratios within genotypes, and number of haplotypes per  
378 individual are redundant with  $H_{ind}/H_E$  and unnecessary if it has already been used for filtering. In  
379 addition to its use for detecting paralogs in highly duplicated genomes,  $H_{ind}/H_E$  can be used for  
380 marker filtering in less duplicated genomes where occasional paralogs are still an issue.  
381 Additionally, in any species, markers with low values of  $H_{ind}/H_E$  (e.g. below the 95% confidence  
382 interval generated by simulated data) are likely to have null alleles, high overdispersion, or other  
383 technical issues and should generally be removed from the dataset. We found that using  $H_{ind}/H_E$   
384 to filter our *M. sacchariflorus* dataset impacted minor allele frequency and proportion of markers  
385 in genes in ways consistent with the removal of markers with null alleles, collapsed paralogs, or  
386 false alleles due to sequencing error (Table 2 and Fig. 3).

387 Although less accurate for determining ploidy than techniques such as flow cytometry, when  
388 averaged within individuals,  $H_{ind}/H_E$  can be used to identify individuals whose ploidy might  
389 deviate from expectations and should be confirmed. If flow cytometry is not an option, several  
390 other tools exist for the estimation of ploidy directly from next-generation sequencing data [25].  
391 Lastly,  $H_{ind}/H_E$  could be potentially useful for improving reference genome assemblies,

392 increasing the value of complementing a de novo assembly with a resequencing or genotyping-  
393 by-sequencing effort in a large population or diversity panel. Regions of the reference genome  
394 that contain collapsed paralogs are expected to have inflated  $H_{ind}/H_E$  values, which could be  
395 visualized in a smoothed plot of  $H_{ind}/H_E$  vs. alignment position.

396 At a minor allele frequency of 0.05, a read depth of five or higher is sufficient to estimate  
397  $H_{ind}/H_E$  with minimal variance (Fig. 6). It is notable that a read depth of five is too low to call  
398 genotypes with confidence, to some extent in diploids but especially in polyploids. However,  
399 using the  $H_{ind}/H_E$  statistic, such low depth data are useful for a variety of applications such as  
400 identification of outlier individuals in terms of ploidy and hybridity, estimation of inbreeding,  
401 identification of loci with technical issues, and assessment of reference genome quality. This in  
402 turn can enable researchers to reduce sequencing costs by generating preliminary, low-depth  
403 datasets to evaluate these issues before (or instead of) sequencing more deeply.

## 404 **Conclusions**

405 Here we introduce the  $H_{ind}/H_E$  statistic, which can be used for evaluating marker and sample  
406 quality in genotyping-by-sequencing datasets for a variety of downstream applications. We  
407 demonstrate that reads from paralogous loci cause the statistic to be above the expected value,  
408 whereas technical issues such as overdispersion and null alleles cause the statistic to be below  
409 the expected value. In typical datasets (hundreds of individuals, read depth above five) the  
410 statistic has sufficiently low variance to be useful for filtering loci. The polyRAD R package can  
411 estimate  $H_{ind}/H_E$ , suggest filtering cutoffs based on simulated data, and perform genotyping after  
412 filtering.

## 413 **Materials and Methods**

### 414 **Implementation in polyRAD**

415 Functions for estimating  $H_{ind}/H_E$  and  $H_{ind}/H_{E.map}$  are available in polyRAD v1.2 and later, and are  
416 named *HindHe* and *HindHeMapping*, respectively. Both utilize an internal Rcpp function for  
417 fast calculation, take a *RADdata* object as input, and return a matrix of values, with samples in  
418 rows and loci in columns. The mean value across rows can then be used to get a per-sample  
419 estimate, for identifying individuals that are interspecies hybrids or unexpected ploidies. The  
420 mean value across columns can be used to get a per-locus estimate for filtering loci.

421 Additionally, polyRAD v1.5 and later includes the *ExpectedHindHe* and  
422 *ExpectedHindHeMapping* functions, which simulate data to emulate the sample size, allele  
423 frequency distribution or parental genotypes, and read depth distribution of an empirical dataset,  
424 and return the distribution of  $H_{ind}/H_E$  as if all loci were Mendelian, giving the user reasonable  
425 thresholds to use for filtering loci.

426 PolyRAD v1.6 is currently available on CRAN, and can be installed using  
427 `install.packages("polyRAD")`.

### 428 **Datasets for testing**

429 Two types of datasets were used to test  $H_{ind}/H_E$ : (1) empirical data from a diversity panel of  
430 *Miscanthus sacchariflorus*, and (2) simulated datasets of diversity panels and of biparental F1  
431 mapping populations. Previously published RAD-seq data for an *M. sacchariflorus* diversity  
432 panel [22] were used for the empirical tests. All species in the *Miscanthus* genus share an  
433 ancient genome duplication, increasing the chromosome number to 19 from the base of 10 in the  
434 Andropogoneae tribe [26–28]. Moreover, some populations of *M. sacchariflorus* display

435 autotetraploidy in addition to this genome duplication ( $4x = 76$ ) [22, 29], allowing us to test our  
436 algorithm in situations where tetrasomic inheritance is expected, in addition to the more typical  
437 disomic inheritance. *Miscanthus* is also highly heterozygous due to being wind-pollinated and  
438 self-incompatible [30], thus heterozygosity cannot be used to identify paralogs as easily as it  
439 could in an inbred crop species. Together, these factors make *M. sacchariflorus* an ideal test  
440 case.

441 To compare values of  $H_{ind}/H_E$  in putatively Mendelian markers versus collapsed paralogs,  
442 markers were called from the same dataset using either *Miscanthus sinensis* or *Sorghum bicolor*  
443 as a reference because *M. sinensis* has a whole genome duplication with respect to *S. bicolor*.  
444 Raw sequence reads from *M. sacchariflorus* were processed by the TASSEL-GBSv2 pipeline  
445 [31] to identify unique tag sequences and their depths in all individuals. Tag sequences were  
446 then aligned to the *Miscanthus sinensis* v7.1 reference genome [32] and the *Sorghum bicolor*  
447 v3.1.1 reference genome [33] using Bowtie 2 [34]. The tag manager feature of TagDigger [35]  
448 was used to process the SAM files, recording the alignment location of each tag in both reference  
449 genomes. Tag alignment locations within the *S. bicolor* reference were retained for further  
450 analysis if they corresponded to two alignment locations in the *M. sinensis* reference matching  
451 the known synteny between chromosomes. Under this filtering, 239,501 tags were retained at  
452 18,402 *S. bicolor* alignment locations corresponding to 36,804 *M. sinensis* alignment locations,  
453 in a set of 356 diploid and 268 tetraploid individuals.  $H_{ind}/H_E$  was then estimated per-locus in  
454 polyRAD for both the *M. sinensis* and *S. bicolor* alignments.

455 To compare the variance of  $H_{ind}/H_E$  when biallelic SNPs were used versus multiallelic,  
456 haplotype-based markers, the TASSEL-GBSv2 pipeline was used to call SNP variants from *M.*  
457 *sacchariflorus* and export them to VCF. Markers from chromosome 1 were imported to

458 polyRAD using *VCF2RADdata*, with and without the option to phase SNPs into haplotypes,  
459 yielding 3710 and 10,458 loci, respectively. The phasing performed by *VCF2RADdata* only  
460 phases SNPs that are certain to have originated from the same reads based on allelic read depth  
461 and physical distance.  $H_{ind}/H_E$  was then estimated by locus in polyRAD separately for diploids  
462 and tetraploids.

463 Simulated diversity panel datasets were generated in order to assess the effect of minor allele  
464 frequency, sample size, read depth, sequencing error, overdispersion, inbreeding, ploidy, and null  
465 alleles on variance and bias of the  $H_{ind}/H_E$  statistic, using the *SimGenotypes* and *SimAlleleDepth*  
466 functions in polyRAD v1.6. See Clark et al. [12] (Eqn. 2) for a definition of the overdispersion  
467 parameter; lower values result in allelic read depths that deviate further from the ratios expected  
468 based on allelic dosage. Three sets of data were simulated. (1) Minor allele frequencies of 0.01,  
469 0.05, and 0.1; sample sizes of 100, 500, and 1000; and genotype read depths of 2, 5, 10, 20, 50,  
470 and 100 were simulated in all combinations under diploidy and tetraploidy, with no inbreeding, a  
471 sequencing error rate of 0.001, and an overdispersion parameter of 20. For each combination,  
472 5000 biallelic loci were simulated, as well as 5000 collapsed paralogs that each consisted of two  
473 Mendelian loci combined. (2) Minor allele frequencies of 0.01 and 0.05, overdispersion  
474 spanning all integers from 5 to 20, sequencing error rates of 0 and 0.001, and inbreeding ( $F$ ; the  
475 probability that two locus copies in an individual are identical by descent) spanning all intervals  
476 of 0.1 from 0 to 1 were simulated in all combinations under diploidy and tetraploidy, with a  
477 sample size of 500 and a read depth of 20. For each combination, 20,000 biallelic loci were  
478 simulated. (3) Minor non-null allele frequencies of 0.01 and 0.05 and null allele frequencies of  
479 0.01, 0.05, 0.1, and 0.2 were simulated in all combinations under diploidy and tetraploidy, with a  
480 sample size of 500, a read depth of 20, a sequencing error rate of 0.001, overdispersion of 20,

481 and no inbreeding. For each combination, 5000 triallelic (with one allele being null, i.e. having  
482 all of its reads discarded) loci were simulated.

483 Simulated F1 mapping population datasets were generated in order to assess the effect of ploidy  
484 and marker type on variance of the  $H_{ind}/H_E$  statistic. For diploids, testcross (homozygote x  
485 heterozygote) and F2 (heterozygote x heterozygote) markers were evaluated. For tetraploids,  
486 simplex x nulliplex (AAAB x AAAA), duplex x nulliplex (AABB x AAAA), triplex x nulliplex  
487 (ABBB x AAAA), simplex x simplex (AAAB x AAAB), simplex x duplex (AAAB x AABB),  
488 simplex x triplex (AAAB x ABBB), and duplex x duplex (AABB x AABB) markers were  
489 evaluated. For each marker type, 5000 biallelic markers were simulated in a population with 500  
490 offspring, with a read depth of 20, a sequencing error rate of 0.001, and overdispersion parameter  
491 of 20.

492 To evaluate effectiveness of various approaches for filtering paralogs, 1000 Mendelian loci and  
493 1000 collapsed paralogs were simulated in 200 diploid and 200 tetraploid individuals each at  
494 three levels of inbreeding. Number of alleles was evenly distributed from two to eight in  
495 Mendelian loci. Allele frequency was sampled from a gamma distribution with shape of 0.3 and  
496 scale of 1, divided by 10 and added to 0.01 to ensure a minimum minor allele frequency, given  
497 that allele frequency filtering is typically performed during variant calling and/or data import.  
498 One allele frequency at each locus was generated as one minus the sum of all other allele  
499 frequencies, to emulate the typical situation of one common allele and one or more rare alleles.  
500 Genotypes were simulated from the allele frequencies assuming an inbreeding coefficient ( $F$ ) of  
501 0.1, 0.5, or 0.9. Mean read depth per locus was drawn from a gamma distribution with a shape  
502 of 3.2 and scale of 8. Read depth at individual genotypes was then drawn from a gamma  
503 distribution with the locus depth / 10 as the shape, and a scale of 10. Allelic read depth was

504 simulated assuming an overdispersion parameter of 20 and a sequencing error rate of 0.001.  
505 Collapsed paralogs were simulated in the same way, but with number of alleles per locus ranging  
506 from one to eight, and two random loci being combined to form a collapsed paralog.

## 507 **Comparison with other approaches**

508 To call genotypes for the  $H_O/H_E$  and Z-score [19] approaches, the *IterateHWE* function in  
509 polyRAD was used with default parameters to obtain genotype probabilities, and then  
510 *GetProbableGenotypes* was used to get discrete genotypes, with genotypes set to missing if  
511 allele copy numbers did not add up to the ploidy. To extend its use to polyploids,  $H_O$  was  
512 estimated as the probability that two alleles sampled from a genotype without replacement would  
513 be different from each other, averaged across individuals within a locus. The Z-score approach  
514 [19] was originally only defined for biallelic markers in diploids. To extend it for multiallelic  
515 markers and polyploid species, for each marker genotypes with  $ploidy - 1$  copies of the most  
516 common allele (i.e. the heterozygous genotype class expected to be most common) were  
517 identified, and allelic read depth summed across those samples. Deviation of read depth of the  
518 most common allele from the expected ratio was then estimated as a Z-score:

$$519 \text{ Eqn 9: } Z = \frac{\frac{ploidy-1}{ploidy} * N - N_A}{\sqrt{N * \frac{ploidy-1}{ploidy} * \frac{1}{ploidy}}}$$

520 Where  $N$  is the total read depth across all samples in the given heterozygous genotype class, and  
521  $N_A$  is the read depth of the common allele summed across those same samples. The number of  
522 haplotypes per genotype was counted as the number of haplotypes with read depth of three or  
523 higher, following Willis et al. [17].

524 **Declarations**

525 **Ethics approval and consent to participate**

526 Not applicable

527 **Consent for publication**

528 Not applicable

529 **Availability of data and materials**

530 Raw sequence reads for *M. sacchariflorus* are available on the NCBI Sequence Read Archive

531 under accessions SRP026347, SRP048207, SRP063572, and SRP087645. Genotype calls and

532 read depths for *M. sacchariflorus* are available on the Illinois Data Bank at

533 [https://doi.org/10.13012/B2IDB-8170405\\_V1](https://doi.org/10.13012/B2IDB-8170405_V1). Tag alignments and counts for *M. sacchariflorus*

534 using the *M. sinensis* and *S. bicolor* reference genomes are available on the Illinois Data Bank at

535 [https://doi.org/10.13012/B2IDB-4814898\\_V1](https://doi.org/10.13012/B2IDB-4814898_V1). All scripts for testing the  $H_{ind}/H_E$  statistic are

536 available on GitHub at [https://github.com/lvclark/paralog\\_id](https://github.com/lvclark/paralog_id), archived on Zenodo at

537 <https://doi.org/10.5281/zenodo.5425343>.

538 Project name: polyRAD

539 Project home page: <https://github.com/lvclark/polyRAD>

540 Archived version: <https://doi.org/10.5281/zenodo.1143744>

541 Operating system: Platform independent

542 Programming language: R, C++ via Rcpp

543 Other requirements:  $R \geq 3.5.0$ ; CRAN packages fastmatch, Rcpp, and stringi; Bioconductor

544 package pcaMethods



545 License: GNU GPL ( $\geq 2$ )

546 Any restrictions to use by non-academics: None

## 547 **Competing interests**

548 The authors declare they have no competing interests.

## 549 **Funding**

550 This material is based upon work supported by the National Science Foundation under Grant No.  
551 1661490. The funding body was not involved in the design, analysis, or interpretation of the  
552 study.

## 553 **Authors' contributions**

554 LVC designed the  $H_{ind}/H_E$  statistic, wrote the polyRAD software, performed the analysis, and  
555 wrote the manuscript. WM performed literature review and tested the software and statistic.  
556 AEL gave statistical advice. EJS provided the *M. sacchariflorus* datasets. All authors read and  
557 approved the final manuscript.

## 558 **Acknowledgements**

559 We thank Jiale He for testing various methods to detect paralogous loci in *Miscanthus*.

## 560 **References**

- 561 1. Renny-Byfield S, Wendel JF. Doubling down on genomes: Polyploidy and crop plants. *Am J*  
562 *Bot.* 2014;101:1711–25. doi:10.3732/ajb.1400119.
- 563 2. Gregory TR, Mable BK. Polyploidy in Animals. In: Gregory TR, editor. *The Evolution of the*  
564 *Genome*. San Diego: Elsevier; 2005. p. 427–517.
- 565 3. Session AM, Uno Y, Kwon T, Chapman JA, Toyoda A, Takahashi S, et al. Genome evolution  
566 in the allotetraploid frog *Xenopus laevis*. *Nature*. 2016;538:336–43. doi:10.1038/nature19840.
- 567 4. Dufresne F, Stift M, Vergilino R, Mable BK. Recent progress and challenges in population

- 568 genetics of polyploid organisms: an overview of current state-of-the-art molecular and statistical  
569 tools. *Mol Ecol.* 2014;23:40–69. doi:10.1111/mec.12581.
- 570 5. Kaur S, Francki MG, Forster JW. Identification, characterization and interpretation of single-  
571 nucleotide sequence variation in allopolyploid crop species. *Plant Biotechnol J.* 2012;10:125–38.  
572 doi:10.1111/j.1467-7652.2011.00644.x.
- 573 6. Nadukkalam Ravindran P, Bentzen P, Bradbury IR, Beiko RG. PMERGE: Computational  
574 filtering of paralogous sequences from RAD-seq data. *Ecol Evol.* 2018;8:7002–13.  
575 doi:10.1002/ece3.4219.
- 576 7. Verdu CF, Guichoux E, Quevauvillers S, De Thier O, Laizet Y, Delcamp A, et al. Dealing  
577 with paralogy in RADseq data: in silico detection and single nucleotide polymorphism validation  
578 in *Robinia pseudoacacia* L. *Ecol Evol.* 2016;6:7323–33. doi:10.1002/ece3.2466.
- 579 8. Meirmans PG, Van Tienderen PH. The effects of inheritance in tetraploids on genetic diversity  
580 and population divergence. *Heredity (Edinb).* 2013;110:131–7. doi:10.1038/hdy.2012.80.
- 581 9. Hohenlohe PA, Amish SJ, Catchen JM, Allendorf FW, Luikart G. Next-generation RAD  
582 sequencing identifies thousands of SNPs for assessing hybridization between rainbow and  
583 westslope cutthroat trout. *Mol Ecol Resour.* 2011;11 Suppl 1:117–22. doi:10.1111/j.1755-  
584 0998.2010.02967.x.
- 585 10. Arruda MP, Brown P, Brown-Guedira G, Krill AM, Thurber C, Merrill KR, et al. Genome-  
586 wide association mapping of Fusarium head blight resistance in wheat using genotyping-by-  
587 sequencing. *Plant Genome.* 2016;9. doi:10.3835/plantgenome2015.04.0028.
- 588 11. Tinker NA, Bekele WA, Hattori J. Haplotag: Software for Haplotype-Based Genotyping-by-  
589 Sequencing Analysis. *G3.* 2016;6:857–63. doi:10.1534/g3.115.024596.
- 590 12. Clark LV, Lipka AE, Sacks EJ. polyRAD: Genotype Calling with Uncertainty from  
591 Sequencing Data in Polyploids and Diploids. *G3.* 2019;9:663–73. doi:10.1534/g3.118.200913.
- 592 13. Gerard D, Ferrão LFV, Garcia AAF, Stephens M. Genotyping Polyploids from Messy  
593 Sequencing Data. *Genetics.* 2018;210 November:789–807. doi:10.1534/genetics.118.301468.
- 594 14. Blischak PD, Kubatko LS, Wolfe AD. SNP genotyping and parameter estimation in  
595 polyploids using low-coverage sequencing data. *Bioinformatics.* 2018;34:407–15.  
596 doi:10.1093/bioinformatics/btx587.
- 597 15. Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH. Stacks: building and  
598 genotyping loci de novo from short-read sequences. *G3.* 2011;1:171–82.  
599 doi:10.1534/g3.111.000240.
- 600 16. Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. Double Digest RADseq: An  
601 Inexpensive Method for De Novo SNP Discovery and Genotyping in Model and Non-Model  
602 Species. *PLoS One.* 2012;7:e37135. doi:10.1371/journal.pone.0037135.
- 603 17. Willis SC, Hollenbeck CM, Puritz JB, Gold JR, Portnoy DS. Haplotyping RAD loci: an  
604 efficient method to filter paralogs and account for physical linkage. *Mol Ecol Resour.*  
605 2017;17:955–65. doi:10.1111/1755-0998.12647.

- 606 18. Gayral P, Melo-Ferreira J, Glémin S, Bierne N, Carneiro M, Nabholz B, et al. Reference-  
607 Free Population Genomics from Next-Generation Transcriptome Data and the Vertebrate-  
608 Invertebrate Gap. *PLoS Genet.* 2013;9. doi:10.1371/journal.pgen.1003457.
- 609 19. McKinney GJ, Waples RK, Seeb LW, Seeb JE. Paralogs are revealed by proportion of  
610 heterozygotes and deviations in read ratios in genotyping-by-sequencing data from natural  
611 populations. *Mol Ecol Resour.* 2017;17:656–69. doi:10.1111/1755-0998.12613.
- 612 20. Lu F, Lipka AE, Glaubitz J, Elshire R, Cherney JH, Casler MD, et al. Switchgrass genomic  
613 diversity, ploidy, and evolution: novel insights from a network-based SNP discovery protocol.  
614 *PLoS Genet.* 2013;9:e1003215. doi:10.1371/journal.pgen.1003215.
- 615 21. Simpson EH. Measurement of diversity. *Nature.* 1949;163:688. doi:10.1038/163688a0.
- 616 22. Clark LV, Jin X, Petersen KK, Anzoua KG, Bagmet L, Chebukin P, et al. Population  
617 structure of *Miscanthus sacchariflorus* reveals two major polyploidization events, tetraploid-  
618 mediated unidirectional introgression from diploid *M. sinensis*, and diversity centred around the  
619 Yellow Sea. *Ann Bot.* 2019;124:731–48. doi:10.1093/aob/mcy161.
- 620 23. Clevenger JP, Ozias-Akins P. SWEEP: A Tool for Filtering High-Quality SNPs in Polyploid  
621 Crops. *G3.* 2015;5:1797–803. doi:10.1534/g3.115.019703.
- 622 24. Mollinari M, Olukolu BA, Da Pereira GS, Khan A, Gemenet D, Craig Yench G, et al.  
623 Unraveling the hexaploid sweetpotato inheritance using ultra-dense multilocus mapping. *G3*  
624 *Genes, Genomes, Genet.* 2020;10:281–92. doi:10.1534/g3.119.400620.
- 625 25. Kyriakidou M, Tai HH, Anglin NL, Ellis D, Strömviik M V. Current strategies of polyploid  
626 plant genome sequence assembly. *Front Plant Sci.* 2018;9:1660. doi:10.3389/fpls.2018.01660.
- 627 26. Ma X-F, Jensen E, Alexandrov N, Troukhan M, Zhang L, Thomas-Jones S, et al. High  
628 resolution genetic mapping by genome sequencing reveals genome duplication and tetraploid  
629 genetic structure of the diploid *Miscanthus sinensis*. *PLoS One.* 2012;7:e33821.  
630 doi:10.1371/journal.pone.0033821.
- 631 27. Swaminathan K, Chae WB, Mitros T, Varala K, Xie L, Barling A, et al. A framework genetic  
632 map for *Miscanthus sinensis* from RNAseq-based markers shows recent tetraploidy. *BMC*  
633 *Genomics.* 2012;13:142. doi:10.1186/1471-2164-13-142.
- 634 28. Kim C, Zhang D, Auckland SA, Rainville LK, Jakob K, Kronmiller B, et al. SSR-based  
635 genetic maps of *Miscanthus sinensis* and *M. sacchariflorus*, and their comparison to sorghum.  
636 *Theor Appl Genet.* 2012;124:1325–38. doi:10.1007/s00122-012-1790-1.
- 637 29. Clark LV, Stewart JR, Nishiwaki A, Toma Y, Kjeldsen JB, Jørgensen U, et al. Genetic  
638 structure of *Miscanthus sinensis* and *Miscanthus sacchariflorus* in Japan indicates a gradient of  
639 bidirectional but asymmetric introgression. *J Exp Bot.* 2015;66:4213–25.  
640 doi:10.1093/jxb/eru511.
- 641 30. Hirayoshi I, Nishikawa K, Kato R. Cytogenetical Studies on forage plants. (IV) Self-  
642 incompatibility in *Miscanthus*. *Japanese J Plant Breed.* 1955;5:167–70.  
643 doi:10.1270/jsbbs1951.5.167.
- 644 31. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. TASSEL:

- 645 software for association mapping of complex traits in diverse samples. *Bioinformatics*.  
646 2007;23:2633–5. doi:10.1093/bioinformatics/btm308.
- 647 32. Mitros T, Session AM, James BT, Wu GA, Belaffif MB, Clark L V., et al. Genome biology  
648 of the paleotetraploid perennial biomass crop *Miscanthus*. *Nat Commun*. 2020;11:5442.  
649 doi:10.1038/s41467-020-18923-6.
- 650 33. McCormick RF, Truong SK, Sreedasyam A, Jenkins J, Shu S, Sims D, et al. The *Sorghum*  
651 *bicolor* reference genome: improved assembly, gene annotations, a transcriptome atlas, and  
652 signatures of genome organization. *Plant J*. 2018;93:338–54. doi:10.1111/tpj.13781.
- 653 34. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.  
654 2012;9:357–9. doi:10.1038/nmeth.1923.
- 655 35. Clark LV, Sacks EJ. TagDigger: user-friendly extraction of read counts from GBS and RAD-  
656 seq data. *Source Code Biol Med*. 2016;11:11. doi:10.1186/s13029-016-0057-7.
- 657