

1 Resolving microsatellite genotype ambiguity in populations of al-
2 lopolyploid and diploidized autopolyploid organisms using negative
3 correlations between allelic variables

4 Lindsay V. Clark^{1*}, Andrea Drauch Schreier²

5 ¹Department of Crop Sciences, University of Illinois, Urbana-Champaign, 1201 W. Gre-
6 gory Drive, Urbana, IL 61801, USA; ²Department of Animal Science, University of Cali-
7 fornia – Davis, Davis, CA 95616, USA

8 *Correspondence: Lindsay V. Clark, E-mail: lvclark@illinois.edu

9 Keywords: polyploidy; R package; simple sequence repeat (SSR); sturgeon

10 Running title: Assigning alleles to isoloci in polyploids

11 **Abstract**

12 A major limitation in the analysis of genetic marker data from polyploid organisms is
13 non-Mendelian segregation, particularly when a single marker yields allelic signals from
14 multiple, independently segregating loci (isoloci). However, with markers such as mi-
15 crosatellites that detect more than two alleles, it is sometimes possible to deduce which
16 alleles belong to which isoloci. Here we describe a novel mathematical property of codom-
17 inant marker data when it is recoded as binary (presence/absence) allelic variables: under
18 random mating in an infinite population, two allelic variables will be negatively correlated
19 if they belong to the same locus, but uncorrelated if they belong to different loci. We
20 present an algorithm to take advantage of this mathematical property, sorting alleles into
21 isoloci based on correlations, then refining the allele assignments after checking for consis-
22 tency with individual genotypes. We demonstrate the utility of our method on simulated
23 data, as well as a real microsatellite dataset from a natural population of octoploid white
24 sturgeon (*Acipenser transmontanus*). Our methodology is implemented in the R package
25 POLYSAT version 1.5.

26 Introduction

27 Polyploidy, both recent and ancient, is pervasive throughout the plant kingdom (Udall
28 & Wendel, 2006), and to a lesser extent, the animal kingdom (Gregory & Mable, 2005).
29 However, genetic studies of polyploid organisms face considerable limitations, given that
30 most genetic analyses were designed under the paradigm of diploid Mendelian segregation.
31 In polyploids, molecular markers typically produce signals from all copies of duplicated
32 loci, causing difficulty in the interpretation of marker data (Dufresne *et al*, 2014). If signal
33 (e.g. fluorescence in a SNP assay, or peak height of microsatellite amplicons in capillary
34 electrophoresis) is not precisely proportional to allele copy number, partial heterozygotes
35 may be impossible to distinguish from each other (*e.g.* AAAB vs. AABB vs. ABBB)
36 (Clark & Jasieniuk, 2011; Dufresne *et al*, 2014). However, under polysomic inheritance (all
37 copies of a locus having equal chances of pairing with each other at meiosis), it is possible
38 to deal with allele copy number ambiguity using an iterative algorithm that estimates
39 allele frequencies, estimates genotype probabilities, and re-estimates allele frequencies
40 until convergence is achieved (De Silva *et al*, 2005; Falush *et al*, 2007). Genotypes cannot
41 be determined with certainty using such methods, but population genetic parameters can
42 be estimated.

43 The situation is further complicated when not all copies of a locus pair with each other
44 with equal probability at meiosis. “Disomic inheritance” refers to situations in which the
45 locus behaves as multiple independent diploid loci (Obbard *et al*, 2006); similarly, one
46 could refer to an octoploid locus as having “tetrasomic inheritance” if it behaved as two
47 tetrasomic loci. In this manuscript we will refer to duplicated loci that do not pair with
48 each other at meiosis (or pair infrequently) as “isoloci” after Obbard *et al* (2006). When
49 a genetic marker consists of multiple isoloci, it is not appropriate to analyze that marker
50 under the assumption of polysomic inheritance; for example, if allele A can be found at
51 both isoloci but allele B is only found at one isolocus in a population, the genotypes
52 AAAB and AABB are possible but AB BB is not (excluding rare events of meiotic pairing
53 between isoloci). Markers from autopolyploids that have undergone diploidization are

54 likely to behave as multiple isoloci; a locus may still exist in multiple duplicated copies,
55 but the chromosomes on which those copies reside may have diverged so much that they
56 no longer pair at meiosis, or pair with different probabilities (Obbard *et al*, 2006). This
57 segregation pattern is also typically the case in allopolyploids, in which homeologous
58 chromosomes from two different parent species might not pair with each other during
59 meiosis. Further, meiotic pairing in allopolyploids may occur between both homologous
60 and homeologous chromosome pairs, but at different rates based on sequence similarity
61 (Gaeta & Pires, 2010; Obbard *et al*, 2006), which often differs from locus to locus even
62 within a species (Dufresne *et al*, 2014). Waples (1988) proposed a method for estimating
63 allele frequencies in polyploids under disomic inheritance, although it requires that allele
64 dosage can be determined in heterozygotes (in his example, by intensity of allozyme bands
65 on a gel) and allows a maximum of two alleles per locus, with both isoloci possessing both
66 alleles. De Silva *et al* (2005) describe how their method for estimating allele frequencies
67 under polysomic inheritance, allowing for multiple alleles, can be extended to cases of
68 disomic inheritance, but require that isoloci have non-overlapping allele sets, and do not
69 address the issue of how to determine which alleles belong to which isolocus.

70 Given that marker data do not follow straightforward Mendelian laws in polyploid
71 organisms, they are often recoded as a matrix of ones and zeros reflecting the presence
72 and absence of alleles (sometimes referred to as “allelic phenotypes”; Obbard *et al*, 2006).
73 In mapping populations such binary data can be useful if one parent is heterozygous for
74 a particular allele and the other parent lacks that allele, in which case segregation may
75 follow a 1:1 ratio and can be analyzed under the diploid testcross model (Swaminathan
76 *et al*, 2012; Rousseau-Gueutin *et al*, 2008) (other ratios are possible, in which case the
77 testcross model does not apply). However, in natural populations, inheritance of dominant
78 (presence/absence) markers typically remains ambiguous, and such markers are treated as
79 binary variables that can be used to assess similarity among individuals and populations
80 but are inappropriate for many population genetic analyses, *e.g.* tests that look for
81 departures from or make assumptions of Hardy-Weinberg Equilibrium (Clark & Jasieniuk,

82 2011).

83 Microsatellites are a special case given that they have multiple alleles, allowing for the
84 possibility of assigning alleles to isoloci, which would drastically reduce the complexity
85 of interpreting genotypes in allopolyploids and diploidized autopolyploids. For example,
86 if an allotetraploid individual has alleles A, B, and C, and if A and B are known to
87 belong to one isolocus and C to the other, the genotype can be recoded as AB at one
88 isolocus and CC at the other isolocus, and the data can be subsequently processed as if
89 they were diploid. If two isoloci are sufficiently diverged from each other, they may have
90 entirely different sets of alleles. This is in contrast to other markers such as SNPs and
91 AFLPs that only have two alleles (except in rare cases of multi-allelic SNPs), in which case
92 isoloci must share at least one allele (or be monomorphic, and therefore uninformative).
93 With microsatellites, one could hypothetically examine all possible combinations of allele
94 assignments to isoloci and see which combination was most consistent with the genotypes
95 observed in the dataset, but this method would be impractical in terms of computation
96 time and so alternative methods are needed. Catalán *et al* (2006) proposed a method
97 for assigning microsatellite alleles to isoloci based on the inspection of fully homozygous
98 genotypes in natural populations. In their example with an allotetraploid species, any
99 genotype with just two alleles was assumed to be homozygous at both isoloci, and therefore
100 those two alleles could be inferred to belong to different isoloci. With enough unique
101 homozygous genotypes, all alleles could be assigned to one isolocus or the other, and both
102 homozygous and heterozygous genotypes could be resolved. However, their method made
103 the assumption of no null alleles, and would fail if it encountered any homoplasmy between
104 isoloci (alleles identical in amplicon size, but belonging to different isoloci). Moreover, in
105 small datasets or datasets with rare alleles, it is likely that some alleles in the dataset will
106 never be encountered in a fully homozygous genotype. The method of Catalán *et al* (2006)
107 was never implemented in any software to the best of our knowledge, despite being the
108 only published methodology for splitting polyploid microsatellite genotypes into diploid
109 isoloci.

110 In this manuscript, we present a novel methodology for assigning microsatellite alle-
111 les to isoloci based on the distribution of alleles among genotypes in the dataset. Our
112 method is appropriate for natural populations, as long as the dataset can be split into
113 reasonably-sized groups of individuals (~ 100 individuals or more) lacking strong pop-
114 ulation structure. It is also appropriate for certain mapping populations, including F_2 ,
115 recombinant inbred lines, and doubled haploids. It can be used on organisms of any
116 ploidy as long as each subgenome has the same ploidy, for example octoploid species with
117 four diploid subgenomes or two tetraploid subgenomes, but not two diploid subgenomes
118 and one tetraploid subgenome. Negative correlations between allelic variables are used
119 to cluster alleles into putative isolocus groups, which are then checked against individual
120 genotypes. If necessary, alleles are swapped between clusters or declared homoplasious
121 so that the clusters agree with the observed genotypes within a certain error tolerance.
122 Genotypes can then be recoded, with each marker split into two or more isoloci, such that
123 isoloci can then be analyzed as diploid or polysomic markers. Our method works when
124 there are null alleles, homoplasmy between isoloci, or occasional meiotic recombination be-
125 tween isoloci, albeit with reduced power to find the correct set of allele assignments. We
126 test our methodology on simulated allotetraploid, allohexaploid, and allo-octoploid (hav-
127 ing two tetrasomic genomes) data, and compare its effectiveness to that of the method
128 of Catalán *et al* (2006). We also demonstrate the utility of our method on a real dataset
129 from a natural population of octoploid white sturgeon (*Acipenser transmontanus*). Our
130 methodology, as well as a modified version of the Catalán *et al* (2006) methodology, are
131 implemented in the R package POLYSAT version 1.5.

132 **Materials and Methods**

133 **Rationale**

134 Say that a microsatellite dataset is recoded as an “allelic phenotype” matrix, such that
135 each row represents one individual, and each allele becomes a column (or an “allelic

136 variable”) of ones and zeros indicating whether that allele is present in that individual
137 or not. Under Hardy-Weinberg equilibrium and in the absence of linkage disequilibrium,
138 these allelic variables are expected to be independent if the alleles belong to different
139 loci or different isoloci. However, if two alleles belong to the same locus (or isolocus),
140 the allelic variables should be negatively correlated. This is somewhat intuitive given
141 that the presence of a given allele means that there are fewer locus copies remaining in
142 which the other allele might appear (Fig. 1A). The negative correlation can also be proved
143 mathematically (Supplementary Materials and Methods). We use “correlation” in a broad
144 sense here; “negative correlation” means that the presence of one allele is associated with
145 the absence of another allele or vice versa.

146 **Algorithm for clustering alleles into isoloci**

147 **Preliminary clusters: the alleleCorrelations function**

148 An overview of our algorithm is presented in Fig. 2. To test independence of two binary
149 allelic variables, we use Fisher’s exact test since it is appropriate for small sample sizes,
150 which are likely to occur in typical population genetics datasets when rare alleles are
151 present. A 2-by-2 contingency table is generated for the test, with rows indicating presence
152 or absence of the first allele, columns indicating presence or absence of the second allele,
153 and each cell indicating the number of individuals in that category (Fig. 1B). A one-tailed
154 Fisher’s exact test is used, with the alternative hypothesis being that more individuals just
155 have one allele of the pair than would be expected if the allelic variables were independent,
156 *i.e.* the alternative hypothesis is that the odds ratio is less than one, indicating a negative
157 association between the presence of the first allele and the presence of the second allele.
158 This alternative hypothesis corresponds to the two alleles belonging to the same isolocus,
159 whereas the null hypothesis is that they belong to different isoloci and therefore assort
160 independently. The P-values from Fisher’s exact test on each pair of allelic variables
161 from a single microsatellite marker are then stored in a symmetric square matrix. We
162 expect to see clusters of alleles with low P-values between them; alleles within a cluster

163 putatively belong to the same isolocus. For clustering algorithms, zeros are inserted along
164 the diagonal of the matrix, since the P-values are used as a dissimilarity statistic. The
165 function `alleleCorrelations` in POLYSAT 1.5 produces such a matrix of P-values for a
166 single microsatellite marker. The same function also produces two sets of preliminary
167 assignments of alleles to isoloci, using UPGMA and the Hartigan & Wong (1979) method
168 of K-means clustering, respectively. The `n.subgen` argument is used to specify how many
169 subgenomes the organism has, *i.e.* into how many isoloci each locus should be split.

170 Population structure can also cause correlation between allelic variables, for example
171 if two alleles are both common in one subpopulation and rare in another. Because correla-
172 tion caused by population structure can potentially obscure the correlations that are used
173 by our method, the `alleleCorrelations` function checks for significant positive correla-
174 tions (after Holm-Bonferroni multiple testing correction) between allelic variables, which
175 could only be caused by population structure, scoring error (such as stutter peaks being
176 mis-called as alleles, and therefore tending to be present in the same genotypes as their
177 corresponding alleles), or linkage disequilibrium (if two isoloci are part of a tandem du-
178 plication on the same chromosome, as opposed to duplication resulting from polyploidy),
179 and prints a warning if such correlations are found.

180 If one or more alleles are present in all genotypes in a dataset, it is not possible to
181 perform Fisher's exact test to look for correlations between those fixed allelic variables and
182 any others. The function `alleleCorrelations` therefore checks for fixed alleles before
183 performing Fisher's exact test. Each fixed allele is assigned to its own isolocus. If only
184 one isolocus remains, all remaining alleles are assigned to it. If no isoloci remain (*e.g.*
185 in an allotetraploid with two fixed alleles and several variable alleles), then all remaining
186 alleles are assigned as homoplasious to all isoloci. If multiple isoloci remain (*e.g.* in
187 an allohexaploid with one fixed allele), then Fisher's exact test, K-means clustering and
188 UPGMA are performed to assign the alleles to the remaining isoloci. It is possible that an
189 allele with a very high frequency may be present in all genotypes but not truly fixed (*i.e.*
190 some genotypes are heterozygous). However, allele swapping performed by `testAlGroups`

191 (see below) can assign alleles to an isolocus even if that isolocus already has an allele
192 assigned to it that is present in all individuals.

193 **Corrected clusters: the `testAlGroups` function**

194 Although K-means was more accurate overall than UPGMA using simulated data (Supple-
195 mentary Table 1), UPGMA sometimes assigned alleles correctly when K-means assigned
196 them incorrectly. To choose between K-means and UPGMA when they give different
197 results, the function `testAlGroups` in POLYSAT checks every genotype in the dataset
198 against both results. Assuming no null alleles or homoplasy (which are dealt with later
199 in the algorithm), a genotype is consistent with a set of allele assignments if it has at
200 least one allele belonging to each isolocus, and no more alleles belonging to each isolocus
201 than the ploidy of that isolocus (*e.g.* two in an allotetraploid). The ploidy of isoloci is
202 specified using the `SGploidy` argument. The set of results that is consistent with the
203 greatest number of genotypes is selected, or K-means in the event of a tie. Selecting the
204 best results out of K-means and UPGMA improved the accuracy of allele assignments at
205 all ploidies, particularly hexaploids (Supplementary Table 1).

206 We expected that rarer alleles would be more likely to be assigned incorrectly, given
207 that they would be present in fewer genotypes and therefore there would be less statistical
208 power to detect correlations between their variables and other allelic variables. To cor-
209 rect the allele assignments, an algorithm was added to the `testAlGroups` function that
210 individually swaps the assignment of each rare allele to the other isolocus (or isoloci) and
211 then checks whether the new set of assignments is consistent with a greater number of
212 genotypes than the old set of assignments. If an allele is successfully swapped, then every
213 other rare allele is checked once again, until no more swaps are made. The maximum
214 number of genotypes in which an allele must be present to be considered a rare allele is
215 adjusted using the `rare.al.check` argument to the `testAlGroups` function. We found
216 that swapping alleles present in $\leq 50\%$ of genotypes (`rare.al.check = 0.5`) improved
217 the accuracy of the algorithm (Supplementary Table 1), so we used that setting in all

218 evaluations of the algorithm except where noted otherwise. Note that the frequency of
219 genotypes with a given allele will always be higher than the allele frequency itself, al-
220 though a 50% threshold is still much higher than the cutoff for considering an allele to be
221 “rare” in most population genetic analyses.

222 Although our algorithm attempts primarily to sort alleles into non-overlapping groups,
223 there is always a possibility that different isoloci have some alleles with identical amplicon
224 sizes (homoplasmy). Therefore, we introduced an algorithm to the `testAlGroups` function
225 to check whether any genotypes were still inconsistent with the allele assignments after
226 the allele swapping step, and assign alleles to multiple isoloci until all genotypes (or a
227 particular proportion that can be adjusted with the `threshold` argument, to allow for
228 meiotic or scoring error) are consistent with the allele assignments. The allele that could
229 correct the greatest number of inconsistent genotypes (or in the event of a tie, the one
230 with the lowest P-values from Fisher’s exact test between it and the alleles in the other
231 isolocus) is made homoplasious first, then all genotypes are re-checked and the cycle is
232 repeated until the desired level of agreement between allele assignments and genotypes is
233 met.

234 Mutations in primer annealing sites are a common occurrence with microsatellite mark-
235 ers, and result in alleles that produce no PCR product, known as null alleles. One po-
236 tential issue with null alleles is that, when homozygous, they can result in genotypes
237 that do not appear to have any alleles from one isolocus. Such genotypes are used by
238 the `testAlGroups` function as an indicator that alleles should be swapped or made ho-
239 moplasious, which would be incorrect actions if the genotype resulted from a null allele
240 rather than inaccuracy of allele assignment. We therefore added an argument to the
241 `testAlGroups` function, `null.weight`, to indicate how genotypes with no apparent alleles
242 for one isolocus should be prioritized for determining which alleles to assign as homopla-
243 sious. If null alleles are expected to be common, `null.weight` can be set to zero so that
244 genotypes with no apparent alleles for one isolocus are not used for assigning homoplasmy.
245 The default value of 0.5 for `null.weight` will cause `testAlGroups` to use genotypes with

246 no apparent alleles for one isolocus as evidence of homoplasmy, but with lower priority than
247 genotypes with too many alleles per isolocus. (No argument was added to adjust the allele
248 swapping algorithm, since it only swaps alleles if the overall agreement with the dataset
249 is improved.)

250 **Recoding datasets based on allele assignments: the `processDatasetAllo` and** 251 **`recodeAllopoly` functions**

252 The function `processDatasetAllo` is a wrapper function that runs `alleleCorrelations`
253 and `testAlGroups` in sequence on every marker in the dataset. It tests several parameter
254 sets for `testAlGroups`. If the dataset was divided into subpopulations to prevent bias from
255 population structure, allele assignments from the same parameter set are merged across
256 subpopulations using the `mergeAlleleAssignments` function. `processDatasetAllo` gen-
257 erates a series of plots to indicate assignment quality, and selects a suggested best param-
258 eter set for each locus by first selecting the parameter set that results in the least amount
259 of missing data when the genotypes are recoded, or in the case of a tie the parameter set
260 that results in the fewest homoplasious alleles.

261 The list of allele assignments (output by `processDatasetAllo`) and the original
262 dataset are then passed to the `recodeAllopoly` function, which produces a new dataset
263 in which each marker is split into multiple isoloci. Missing data are substituted for geno-
264 types that cannot be resolved due to homoplasmy in the allele assignments. (For example,
265 if alleles A and B belong to different isoloci, and C belongs to both, the genotype ABC
266 could be AA BC, AC BB, or AC BC, assuming no null alleles.) An argument called
267 `allowAneuploidy` lets the user specify whether to allow for apparent meiotic error. If
268 `allowAneuploidy = TRUE`, for genotypes with too many alleles for one isolocus, the func-
269 tion will adjust the recorded ploidy for the relevant samples and isoloci. (Ploidy is used
270 by other POLYSAT functions, such as those that estimate allele frequency.) Otherwise,
271 missing data are inserted where there are too many alleles per isolocus.

272 **Implementation of the Catalán method: the `catalanAlleles` function**

273 POLYSAT 1.5 also includes an implementation of the algorithm of Catalán *et al* (2006).
274 One difference between our implementation and the original is that we allow ploidies
275 higher than tetraploid, *e.g.* in a hexaploid, a genotype with three alleles is assumed to
276 be fully homozygous. Additionally, after fully homozygous genotypes are examined, fully
277 heterozygous genotypes are also examined if necessary for assigning alleles that were not
278 present in any fully homozygous genotypes. The output of `catalanAlleles` can be passed
279 directly to `recodeAllopolypoly`.

280 **Simulated datasets**

281 The function `simAllopolypoly` was added to POLYSAT in order to generate simulated datasets
282 for testing the accuracy of allele assignment methods. It simulates one locus at a time, and
283 allows for adjustment of the number of isoloci, the ploidy of each isolocus, the number of
284 alleles for each isolocus, the number of alleles that are homoplasious between isoloci, the
285 number of null alleles (producing no amplicon), allele frequencies in the population, the
286 meiotic error rate (frequency at which different isoloci pair with each other at meiosis),
287 and the number of individual genotypes to output. By default, alleles from the first
288 isolocus are labeled A1, A2, etc., alleles from the second isolocus labeled B1, B2, etc., and
289 homoplasious alleles labeled H1, H2, etc.

290 For initial evaluation of clustering methods (Supplementary Table 1), 10,000 simulated
291 markers were generated for 100 individuals each for allotetraploid, allohexaploid, and allo-
292 octoploid (two tetrasomic isoloci) species under Hardy-Weinberg Equilibrium. Although
293 not included in the simulated datasets, note that it is also possible for an octoploid to
294 possess four diploid subgenomes, as in strawberry. Each isolocus had a randomly chosen
295 number of alleles between two and eight, and allele frequencies were generated randomly.
296 A set of allele assignments for one marker was considered to be correct if no alleles were
297 assigned incorrectly.

298 To evaluate the effect of sample size on assignment accuracy, 1000 additional markers

299 were simulated for populations of 50, 100, 200, 400, and 800 individuals for allotetraploid,
300 allohexaploid, and allo-octoploid species.

301 To simulate population structure, 5000 simulated markers were generated for two pop-
302 ulations of 50 allotetraploid individuals. Allele frequencies differed by five fixed amounts
303 (Table 1) between the two populations, with 1000 markers simulated for each amount.

304 The effect of homoplasmy on allele assignment methods was evaluated by simulating
305 1000 allotetraploid markers each for sample sizes of 50, 100, 200, 400, and 800, and
306 homoplasious allele frequencies of 0.1, 0.2, 0.3, 0.4, and 0.5.

307 To evaluate allele assignment when null alleles were present, 5000 markers were sim-
308 ulated for 100 allotetraploid individuals, with 1000 simulated markers at each null allele
309 frequency of 0.1, 0.2, 0.3, 0.4, and 0.5.

310 Occasional pairing between homeologous (in an allopolyploid) or paralogous (in an
311 autopolyploid) chromosomes may occur during meiosis. As a result, offspring may be
312 aneuploid, having too many or too few chromosomes from either homologous pair, or may
313 have translocations between homeologous or paralogous chromosomes. Most commonly,
314 the aneuploidy or translocations will occur in a compensated manner (Chester *et al*, 2015),
315 meaning that for a given pair of isoloci, the total number of copies will be the same as in a
316 non-aneuploid, but one isolocus will have more copies than expected and the other isolocus
317 will have fewer (*e.g.* three copies of one isolocus and one copy of the other isolocus in an
318 allotetraploid). To evaluate the accuracy of allele assignment for isoloci that occasionally
319 pair at meiosis, 4000 markers were simulated for 100 allotetraploid individuals, with 1000
320 simulated markers at each meiotic error rate of 0.01, 0.05, 0.10, and 0.20.

321 A custom script was written to simulate genotypes in allopolyploid mapping popula-
322 tions. Allotetraploid, allohexaploid, and allo-octoploid (with two tetrasomic subgenomes)
323 populations were simulated, with 200 individuals in each population. For each ploidy,
324 1000 loci were simulated for each generation spanning F_2 to F_8 , assuming completely
325 homozygous parents. Allele assignments were performed with the `alleleCorrelations`
326 and `testAlGroups` functions, with `null.weight = 1` and `rare.al.check = 0.25`.

327 Empirical dataset

328 To demonstrate the usefulness of our allele assignment method on a real dataset, we used
329 previously published data from natural populations of octoploid white sturgeon (*Acipenser*
330 *transmontanus*; Drauch Schreier *et al*, 2012). Previous studies of inheritance patterns in
331 this species suggested that it possesses two tetrasomic subgenomes, at least for portions
332 of its genome (Rodzen & May, 2002; Drauch Schreier *et al*, 2011). We selected for allele
333 assignment the eight microsatellite markers that, based on number of alleles per genotype,
334 appeared to be present in eight copies rather than four.

335 Because population structure can impact allele clustering, we first performed a prelim-
336 inary analysis of population structure using the `Lynch.distance` dissimilarity statistic in
337 POLYSAT and principal coordinates analysis (PCoA) using the `cmdscale` function in R.
338 Thirteen microsatellite markers were used for PCoA, including the eight used for allele
339 assignment and five tetrasomic (present in four copies rather than eight) markers. Al-
340 lele assignment methods were then tested on the whole dataset and on a subpopulation
341 identified by PCoA.

342 The `testAllGroups` function was run on the sturgeon dataset with and without allele
343 swapping (`rare.al.check` set to 0.5 and 0, respectively). In checking for homoplasmy, we
344 allowed up to 5% of genotypes to disagree with allele assignments in anticipation of meiotic
345 error, scoring error, or genotypes homozygous for null alleles (`tolerance` = 0.05), and
346 to allow for null alleles at low frequency we set `null.weight` = 0.5 so that genotypes
347 with too many alleles per isolocus would be used for assignment of homoplasmy first, before
348 genotypes with no alleles for one of their isoloci.

349 To evaluate the accuracy and usefulness of allele assignments, we compared G_{ST} (Nei
350 & Chesser, 1983) estimates using the five tetrasomic loci to estimates using the putatively
351 tetrasomic recoded isoloci. Pairs of isoloci were excluded from G_{ST} estimates if they had
352 any homoplasious alleles. Allele frequencies for tetrasomic loci and isoloci were estimated
353 using the method of De Silva *et al* (2005) using the `deSilvaFreq` function in POLYSAT
354 with the selfing rate set to 0.0001. Pairwise G_{ST} between sampling regions was then

355 estimated with the `calcPopDiff` function in POLYSAT.

356 Results

357 Simulated natural populations

358 For all ploidies, we found that the accuracy of both our method and the Catalán *et al*
359 (2006) method was dependent on sample size, and that our method performed better
360 than the Catalán *et al* (2006) method at all sample sizes (Fig. 3). For tetraploids and
361 hexaploids, the effect of sample size was greater on the Catalán *et al* (2006) method than
362 on our method, particularly at small sample sizes. For octoploids, the success of the
363 Catalán *et al* (2006) method was near zero even with 800 individuals in the dataset (due
364 to the low probability of producing fully homozygous genotypes at tetrasomic isoloci),
365 whereas our method had an accuracy of 93% with 800 octoploid individuals.

366 Both negative and positive correlations between allelic variables at different loci can
367 occur when the assumption of random mating is violated by population structure, con-
368 founding the use of negative correlations for assigning alleles to isoloci. We found that
369 accuracy of our method remained high ($\sim 90\%$) even at moderate levels of F_{ST} (~ 0.2 ;
370 Table 1). Interestingly, low levels of population structure ($F_{ST} \approx 0.02$) improved the
371 accuracy of our method to 99%, compared to 94% when $F_{ST} = 0$ (Table 1), probably as
372 a result of an increase in the number of double homozygous genotypes, which would have
373 been informative during the allele swapping step. For this same reason, the Catalán *et al*
374 (2006) method, which depends on double homozygous genotypes, had an improved success
375 rate as population structure increased, and exceeded our method in accuracy at moderate
376 levels of F_{ST} (Table 1). However, accuracy of our method decreased with increasing F_{ST}
377 when $F_{ST} > 0.02$ (Table 1), likely because correlations between alleles caused by popu-
378 lation structure outweighed the benefits of increased homozygosity. In our simulations,
379 significant positive correlations between allelic variables were found in most datasets that
380 had moderate population structure (Table 1).

381 One advantage of our method over that of Catalán *et al* (2006) is that our method
382 allows for alleles belonging to different isoloci to have identical amplicon sizes (homoplasmy).
383 We tested the accuracy of allele assignments across several sample sizes and frequencies
384 of homoplasious alleles, with and without the allele swapping algorithm (Fig. 4). Allele
385 assignments were most accurate when allele swapping was not performed before testing for
386 homoplasious alleles, and when the homoplasious allele was at a frequency of 0.3 in both
387 isoloci. When allele assignments were correct, we tested the mean proportion of genotypes
388 that were resolvable, given several frequencies of a homoplasious allele (Table 2). Although
389 accuracy of assignment had been highest with a homoplasious allele frequency of 0.3, only
390 57% of genotypes could be resolved in such datasets (Table 2).

391 To test the effect of null alleles on the accuracy of our allele assignment method,
392 we simulated datasets in which one isolocus had a null allele (Fig. 5). We found that,
393 when null alleles were present, the accuracy of the algorithm was greatly improved when
394 genotypes lacking alleles for one isolocus were not used as evidence of homoplasmy. We
395 also found that the allele swapping algorithm improved the accuracy of allele assignments
396 when the null allele was at a frequency of 0.1 in the population. However, at higher null
397 allele frequencies (≥ 0.4) allele assignments were more accurate without allele swapping.

398 We simulated datasets in which gametes resulting in compensated aneuploidy (meiotic
399 error) occurred at a range of frequencies from 0.01 to 0.2 (Fig. 6). At all meiotic error
400 rates, the allele swapping algorithm from `testAlGroups` improved the accuracy of allele
401 assignment (Fig. 6). Meiotic error did not have a large impact on the success of our
402 method; even at a meiotic error rate of 0.2 (where 0.5 would be fully autopolyploid), our
403 algorithm still had an accuracy of 62% on datasets of 100 individuals with no homoplasmy,
404 null alleles, or population structure (Fig. 6).

405 We also examined the effect of number of alleles on the accuracy of our method. Ac-
406 curacy was highest when the number of alleles was similar among isoloci (Supplementary
407 Table 2).

408 Assignment of alleles to isoloci in octoploid sturgeon

409 When using principal coordinates analysis to test for genetic structure prior to perform-
410 ing allele assignment, we identified two major genetic groups (Supplementary Table 3,
411 Supplementary Fig. 1) that were similar to the population structure previously observed
412 (Drauch Schreier *et al*, 2012). The smaller group (Pop 2) consisted of only 66 individuals
413 and, likely due to small sample size, produced poor quality allele assignments with high
414 levels of homoplasmy when analyzed by itself (data not shown). We therefore tested our
415 method on Pop 1 (183 individuals) and on the combined set of 249 individuals.

416 For five out of eight loci, our algorithm found allele assignments devoid of homoplasmy
417 when only Pop 1 was used for assignment and when the allele swapping algorithm was
418 used (Table 3). Eliminating the allele swapping algorithm or using the whole dataset for
419 allele assignment increased the number of apparent homoplasious alleles in most cases,
420 and did not reduce the number of apparent homoplasious alleles for any locus (Table 3).
421 For the three loci with homoplasmy, most genotypes in the dataset could not be assigned
422 unambiguously (Table 3). For the five loci with no apparent homoplasmy, nearly all geno-
423 types in Pop 1 could be assigned unambiguously, and approximately one half to three
424 quarters of the genotypes in Pop 2 (which was not used for creating the assignments)
425 could be assigned unambiguously (Table 3). Despite the fact that Pop 1 was previously
426 determined to consist of three subpopulations with pairwise Phi-PT [an F_{ST} analog that
427 can be used on both dominant and codominant markers (Peakall *et al*, 1995)] values rang-
428 ing from 0.06 to 0.17 (Drauch Schreier *et al*, 2012), allelic variable correlations resulting
429 from population structure did not appear to prevent us from obtaining reasonable allele
430 assignments for the five loci without homoplasmy. Significant positive correlations between
431 allelic variables were found at one and two out of eight loci when Pop 1 and the whole
432 dataset were used to make assignments, respectively (data not shown).

433 By recoding allo-octoploid markers as tetrasomic isoloci, we were able to estimate
434 allele frequencies, which would not have been possible otherwise. We were then able to
435 use allele frequencies to estimate pairwise G_{ST} between white sturgeon sampling regions.

436 G_{ST} estimates using recoded isoloci were very similar to estimates obtained using known
437 tetrasomic microsatellite markers (Supplementary Fig. 2), suggesting that allele assign-
438 ments were accurate. Out of the ten recoded isoloci, only one (Atr117_1) was consistently
439 an outlier in terms of G_{ST} estimates, giving especially high estimates between sampling
440 regions corresponding to Pop 1 and Pop 2 (Supplementary Fig. 2, Supplementary Table
441 3). Atr117_1 had especially low genotype variability due to an allele that was present
442 in all Pop 1 genotypes (Table 3), which likely accounted for the unusual G_{ST} estimates
443 at that isolocus. Otherwise, G_{ST} estimates appeared unaffected by the large amounts
444 of missing data introduced into Pop 2 by our method (Table 3, Supplementary Fig. 2),
445 suggesting any bias in allele frequencies caused by the missing data was negligible.

446 **Simulated mapping populations**

447 Negative correlations between allelic variables at the same isolocus can also occur in cer-
448 tain types of mapping populations, enabling the use of our algorithm to assign alleles
449 to isoloci in these populations. There are several requirements that must be met how-
450 ever. 1) To prevent correlations between unlinked allelic variables, all individuals in the
451 population must be equally related to each other. Pedigrees, nested association mapping
452 (NAM) populations, and multiple-cross mating designs are therefore not appropriate. 2)
453 No allele should be present in all individuals in the population. Our method therefore
454 cannot be used on backcross or near isogenic line (NIL) populations, which are expected to
455 segregate only AB and BB genotypes. 3) All alleles belonging to one isolocus should have
456 had the opportunity to pair with each other at meiosis. This eliminates F1 populations,
457 where an individual with genotype AB might be crossed to an individual with genotype
458 CD. However, allele assignments in F₂ populations, as well as related populations such
459 as recombinant inbred line (RIL) and doubled haploid (DH), can be performed with very
460 high accuracy using our algorithm.

461 Accuracy of allele assignment was 100% for allotetraploids and allohexaploids for all
462 population types tested (F₂ to F₈; Table 4). Due to the highly heterozygous nature of

463 tetrasomic loci, accuracy was 14% for allo-octoploids in the F_2 generation. However, accu-
464 racy for allo-octoploids increased to 91% in the F_3 and 100% in F_4 and higher populations,
465 due to increased homozygosity from selfing.

466 Discussion

467 Here we introduce the R package POLYSAT version 1.5, with several new functions ap-
468 plicable to the analysis of allopolyploids and diploidized autopolyploids. These include
469 `simAllopoly`, which generates simulated datasets; `catalanAlleles`, which uses the the
470 Catalán *et al* (2006) method to assign alleles to isoloci; `alleleCorrelations`, which per-
471 forms Fisher’s exact test between each pair of allelic variables from a marker, and then
472 uses K-means clustering and UPGMA to make initial assignments of alleles to isoloci;
473 `testAlGroups`, which checks the consistency of allele assignments with individual geno-
474 types, chooses between the K-means and UPGMA method, swaps alleles to different isoloci
475 if it improves consistency, and identifies homoplasious alleles; `mergeAlleleAssignments`,
476 which merges the allele assignments from two different populations using the same mi-
477 crosatellite marker; `processDatasetAllo`, which runs `alleleCorrelations`, `testAlGroups`
478 (with multiple parameter sets), and `mergeAlleleAssignments` on an entire dataset; and
479 `recodeAllopoly`, which uses allele assignments to recode the dataset, splitting each mi-
480 crosatellite marker into multiple isoloci. An overview of the data analysis workflow is
481 given in Fig. 2. Previous versions of POLYSAT (1.3 and earlier) were restricted in that
482 estimation of allele frequency and certain inter-individual distance metrics could only be
483 performed on autopolyploids. With the ability to assign alleles to isoloci, these parameters
484 may now be estimated for allopolyploids as well.

485 We found that, with simulated data, the accuracy of our allele assignment algorithm
486 was impacted by issues such as homoplasmy and null alleles, and that the optimal param-
487 eters for the algorithm depended on which of these issues were present in the dataset.
488 This suggests, since most users will not know whether their dataset has homoplasmy or

489 null alleles, that the `testAlGroups` function should initially be run with several different
490 parameter sets, and for each locus, the results with the fewest homoplasious alleles should
491 be chosen. A heatmap of the P-values generated from Fisher's exact test can also serve
492 as a qualitative visual indicator of how well the alleles can be separated into isolocus
493 groups. We also found that, although our allele assignment algorithm was negatively
494 impacted by meiotic error (pairing of non-homologous chromosomes during meiosis) and
495 moderate population structure, its accuracy remained fairly high in both cases. Assuming
496 correct allele assignments in a population with meiotic error, `recodeAllopolly` is able to
497 identify some but not all individuals with meiotic error, for example if alleles A, B, and
498 C belonged to one isolocus and D to another, an ABC D individual would be correctly
499 recoded, where as an ABB D individual would be incorrectly recoded as AB DD. Other-
500 wise, `recodeAllopolly` should give 100% accurate results if allele assignments are correct.
501 Sensitivity to population structure is the biggest drawback of our method in comparison
502 to that of Catalán *et al* (2006), which actually has improved results as population struc-
503 ture increases. However, even low frequencies of null alleles, homoplasmy, or meiotic error
504 can cause the method of Catalán *et al* (2006) to fail completely.

505 When discussing homoplasmy with respect to our algorithm, we have referred exclu-
506 sively to homoplasmy between alleles belonging to different isoloci. It is important to note
507 that homoplasmy between alleles within an isolocus is also possible, meaning that two or
508 more alleles belonging to one isolocus are identical in amplicon size but not identical by
509 descent. Although such homoplasmy is an important consideration for analyses that deter-
510 mine similarity between individuals and populations, homoplasmy within isoloci does not
511 affect the allele assignment methods described in this manuscript. Additionally, when
512 discussing null alleles, we have assumed that non-null alleles still exist for all isoloci. It is
513 also possible for an entire isolocus to be null. This is often apparent when a marker has
514 fewer alleles per genotype than expected, *e.g.* a maximum of two alleles per individual in
515 a tetraploid. Such loci should be excluded from the allele assignment analysis described
516 in this manuscript. If they are included in an analysis accidentally, they can be identified

517 by weak K-means/UPGMA clustering of alleles (which can be evaluated from the graph-
518 ical output of `processDatasetAllo`) and by a high proportion of alleles appearing to be
519 homoplasious.

520 Using a real microsatellite dataset from natural populations of white sturgeon, we
521 found that our method was useful for recoding over half of the markers into two inde-
522 pendently segregating isoloci each. Given that white sturgeon are octoploid with two
523 tetrasomic subgenomes (Drauch Schreier *et al*, 2011), we expected this dataset to be
524 problematic; having tetrasomic isoloci as opposed to disomic isoloci would reduce the
525 magnitude of the negative correlations between allelic variables, and was observed in sim-
526 ulations to reduce the accuracy of assignment using our method, although not nearly as
527 severely as the reduction in efficacy of the Catalán *et al* (2006) method (Supplementary
528 Table 1, Fig. 3). In population genetic studies, we expect that microsatellite markers
529 that can be recoded using our method could be used for analyses requiring polysomic or
530 disomic inheritance [for example, estimation of allele frequency and population differenti-
531 ation (Supplementary Fig. 2), Structure (Falush *et al*, 2007), or tests of Hardy-Weinberg
532 Equilibrium], while the remaining markers will still be useful for other analysis (for ex-
533 ample, Mantel tests using simple dissimilarity statistics). Additionally, we found that the
534 allele assignments that we made were still moderately useful for recoding genotypes in a
535 population that was not used for making the assignments. Despite the introduction of
536 missing data into Pop 2 when its genotypes were recoded, G_{ST} estimates were similar
537 to those obtained from non-recoded tetrasomic microsatellites in the same populations
538 (Supplementary Fig. 2). We do however recommend caution when interpreting results
539 from loci where our method has introduced missing data for a large portion of individuals.
540 Such results can be confirmed by comparison to results from loci with little or no missing
541 data.

542 Although inappropriate for biallelic marker systems such as single nucleotide poly-
543 morphisms (SNPs) and dominant marker systems such as AFLPs, the method that we
544 have described could theoretically be used to assign alleles to isoloci in any marker sys-

545 tem in which multiple alleles are the norm. Allozymes, although rarely used in modern
546 studies, are one such system. Although data from genotyping-by-sequencing (GBS, and
547 the related technique restriction site-associated DNA sequencing, or RAD-seq) are typi-
548 cally processed to yield biallelic SNP markers, in the future as typical DNA sequencing
549 read lengths increase, it may become common to find multiple SNPs within the physical
550 distance covered by one read. In that case, haplotypes may be treated as alleles, and
551 negative correlations between haplotypes may be used to assign them to isoloci.

552 **Obtaining polysat 1.5**

553 To obtain POLYSAT, first install the most recent version of R (available at <http://www.r-project.org>),
554 launch R, then at the prompt type:

```
555     install.packages("polysat")
```

556 In the “doc” subdirectory of the package installation, PDF tutorials are available for
557 POLYSAT as a whole and for the methodology described in this manuscript. Source code
558 is available at <https://github.com/lvclark/polysat/> under a GNU GPL-2 license.

559 **Acknowledgements**

560 Author LC was supported by the DOE Office of Science, Office of Biological and En-
561 vironmental Research (grant number DE-SC0012379). We thank Subject Editor Fred-
562 eric Austerlitz and three anonymous reviewers for feedback on an earlier version of this
563 manuscript.

564 **Supporting Information**

- 565 • polysatsupplementary.pdf: Supplementary materials and methods, tables, and fig-
566 ures.
- 567 • allopolyVignette.pdf: Tutorial for creating and using allele assignments in POLYSAT.

- 568 • tables_figs.R, sturgeontest.R: R scripts for reproducing the analyses in this manuscript.
- 569 • sturgeon.txt: White sturgeon microsatellite dataset used in sturgeontest.R.

570 References

- 571 Catalán P, Segarra-Moragues JG, Palop-Esteban M, Moreno C, González-Candelas F
572 (2006) A Bayesian approach for discriminating among alternative inheritance hypothe-
573 ses in plant polyploids: the allotetraploid origin of genus *Bordera* (Dioscoreaceae).
574 *Genetics*, **172**, 1939–1953.
- 575 Chester M, Riley RK, Soltis PS, Soltis DE (2015) Patterns of chromosomal variation in
576 natural populations of the neoallotetraploid *Tragopogon mirus* (Asteraceae). *Heredity*,
577 **114**, 309–317.
- 578 Clark LV, Jasieniuk M (2011) Polysat: an R package for polyploid microsatellite analysis.
579 *Molecular Ecology Resources*, **11**, 562–566.
- 580 De Silva HN, Hall AJ, Rikkerink E, McNeilage MA, Fraser LG (2005) Estimation of allele
581 frequencies in polyploids under certain patterns of inheritance. *Heredity*, **95**, 327–334.
- 582 Drauch Schreier A, Gille D, Mahardja B, May B (2011) Neutral markers confirm the
583 octoploid origin and reveal spontaneous autopolyploidy in white sturgeon, *Acipenser*
584 *transmontanus*. *Journal of Applied Ichthyology*, **27**, 24–33.
- 585 Drauch Schreier A, Mahardja B, May B (2012) Hierarchical patterns of population struc-
586 ture in the endangered fraser river white sturgeon (*Acipenser transmontanus*) and im-
587 plications for conservation. *Canadian Journal of Fisheries and Aquatic Sciences*, **69**,
588 1968–1980.
- 589 Dufresne F, Stift M, Vergilino R, Mable BK (2014) Recent progress and challenges in
590 population genetics of polyploid organisms: an overview of current state-of-the-art
591 molecular and statistical tools. *Molecular Ecology*, **23**, 40–69.
- 592 Falush D, Stephens M, Pritchard JK (2007) Inference of population structure using mul-
593 tilocus genotype data: dominant markers and null alleles. *Molecular Ecology Notes*, **7**,
594 574–578.
- 595 Gaeta RT, Pires JC (2010) Homeologous recombination in allopolyploids: the polyploid
596 ratchet. *New Phytologist*, **186**, 18–28.
- 597 Gregory TR, Mable BK (2005) Polyploidy in animals. In *The Evolution of the Genome*
598 (edited by TR Gregory), chap. 8, pp. 427–517. Elsevier, San Diego.
- 599 Hartigan JA, Wong MA (1979) A K-means clustering algorithm. *Journal of the Royal*
600 *Statistical Society, Series C (Applied Statistics)*, **28**, 100–108.
- 601 Nei M, Chesser RK (1983) Estimation of fixation indices and gene diversities. *Annals of*
602 *Human Genetics*, **47**, 253–259.

- 603 Obbard DJ, Harris SA, Pannell JR (2006) Simple allelic-phenotype diversity and differ-
604 entiation statistics for allopolyploids. *Heredity*, **97**, 296–303.
- 605 Peakall R, Smouse PE, Huff DR (1995) Evolutionary implications of allozyme and RAPD
606 variation in diploid populations of dioecious buffalograss *Buchlo dactyloides*. *Molecular*
607 *Ecology*, **4**, 135–148.
- 608 Rodzen JA, May B (2002) Inheritance of microsatellite loci in white sturgeon (*Acipenser*
609 *transmontanus*). *Genome*, **45**, 1064–1076.
- 610 Rousseau-Gueutin M, Lerceteau-Köhler E, Barrot L, *et al* (2008) Comparative genetic
611 mapping between octoploid and diploid *Fragaria* species reveals a high level of col-
612 linearity between their genomes and the essentially disomic behavior of the cultivated
613 octoploid strawberry. *Genetics*, **179**, 2045–2060.
- 614 Swaminathan K, Chae WB, Mitros T, *et al* (2012) A framework genetic map for *Miscant-*
615 *hus sinensis* from RNAseq-based markers shows recent tetraploidy. *BMC Genomics*,
616 **13**, 142.
- 617 Udall JA, Wendel JF (2006) Polyploidy and crop improvement. *Crop Science*, **46**, S3–S14.
- 618 Waples RS (1988) Estimation of allele frequencies at isoloci. *Genetics*, **118**, 371–384.

619 **Data Accessibility**

620 POLYSAT is available from CRAN (<http://cran.r-project.org>). All datasets and scripts
621 used in this manuscript are provided as Supporting Information.

622 **Author Contributions**

623 LC wrote and tested the software and drafted the manuscript. AS provided the white
624 sturgeon data and gave critical feedback on the manuscript.

625 **Tables and Figures**

Table 1: Percentages of simulated datasets with correct allele assignments under different levels of population structure. Two populations of 50 allotetraploid individuals were simulated under different allele frequencies, then merged into one dataset that was then used for making allele assignments. The value shown in the leftmost column was randomly added or subtracted from the frequency of each allele in the first population to generate the allele frequencies of the second population. For isoloci with odd numbers of alleles, one allele had the same frequency in both populations. For each difference in allele frequency, 1000 simulations were performed (5000 total). F_{ST} was calculated from allele frequencies as $(H_T - H_S)/H_T$, where H_S is the expected heterozytosity in each subpopulation, averaged across the two subpopulations, and H_T is the expected heterozytosity if the two subpopulations were combined into one population with random mating. Means and standard deviations across 1000 simulations are shown for F_{ST} . The third column shows the percentages of datasets in which significant positive correlations were detected between any pair of alleles; positive correlations can be used as an indication that there is population structure in the dataset. The fourth and fifth columns indicate the percentages of datasets with correct allele assignments, using our method and that of Catalán *et al* (2006). 95% confidence intervals are given for percentages.

Difference in allele frequency	F_{ST}	Significant positive correlations	K-means + UPGMA + swap ≤ 0.50	Catalán
0.0	0.000 \pm 0.000	0% \pm 0%	94% \pm 1%	84% \pm 2%
0.1	0.016 \pm 0.004	2% \pm 1%	99% \pm 1%	89% \pm 2%
0.2	0.062 \pm 0.013	21% \pm 3%	93% \pm 2%	94% \pm 1%
0.3	0.117 \pm 0.021	62% \pm 3%	88% \pm 2%	99% \pm 1%
0.4	0.176 \pm 0.026	82% \pm 2%	88% \pm 2%	100% \pm 0%

Table 2: For datasets from Fig. 4 with correct allele assignments at `rare.al.check = 0` (no swapping), percentages of genotypes that could be unambiguously resolved. Means and standard deviations are shown.

Freq. of homoplasious allele	Mean percentage of genotypes that could be resolved
0.1	85.6% \pm 5.6%
0.2	71.2% \pm 8.3%
0.3	59.4% \pm 9.4%
0.4	51.5% \pm 9.0%
0.5	48.5% \pm 7.0%

Table 3: Assignment of alleles from eight microsatellite markers to two tetrasomic genomes in octoploid white sturgeon (*Acipenser transmontanus*). Alleles were assigned using negative correlations, with the exception of Atr117 in Pop 1 due to a fixed allele in that locus and population. Assignments were performed without allele swapping (“No swapping”, `rare.al.check = 0` in `testAlGroups`) and with allele swapping (“Swap ≤ 0.5 ”, `rare.al.check = 0.5`). In testing for homoplasmy `testAlGroups` was run with the defaults of `tolerance = 0.05` to allow for 5% of genotypes to disagree with allele assignments, and `null.weight=0.5` to allow for the possibility of null alleles. Assignments were performed using the whole dataset of 249 individuals (“whole set”) or a subset of 183 individuals based on population structure (“Pop 1”, Supplementary Table 3 and Supplementary Fig. 1). The assignments from Pop 1 with `Swap ≤ 0.5` were then used to split the dataset into isoloci using the `recodeAllopolypoly` function. Genotypes that could not be unambiguously determined were coded as missing data; percentages of missing data in each of two isoloci in Pop 1 and Pop 2 are shown.

Marker	Number of alleles	Number of homoplasious alleles				Percent missing data in recoded dataset	
		Whole set used for assignment		Pop 1 used for assignment		Pop 1	Pop 2
		No swapping	Swap ≤ 0.5	No swapping	Swap ≤ 0.5		
AciG110	20	3	1	0	0	0%, 1%	29%, 29%
As015	18	3	1	2	1	57%, 82%	62%, 80%
AciG35	18	2	0	1	0	0%, 1%	26%, 26%
Atr109	25	6	3	4	3	73%, 74%	70%, 65%
Atr117	22	1	1	0	0	0%, 0%	36%, 36%
AciG52	22	4	1	0	0	0%, 1%	32%, 33%
Atr107	24	3	1	0	0	0%, 1%	45%, 45%
Atr1173	18	3	2	3	2	62%, 77%	64%, 91%

Table 4: Accuracy of allele assignment in mapping populations. Percentages of datasets with accurate allele assignments are shown. 95% confidence intervals are indicated. 1000 loci were simulated, each with 200 individuals.

Generation	Allotetraploid	Allohexaploid	Allo-octoploid
F ₂	100% ± 0%	100% ± 0%	13.6% ± 2.1%
F ₃	100% ± 0%	100% ± 0%	91.4% ± 1.7%
F ₄	100% ± 0%	100% ± 0%	100% ± 0%
F ₅	100% ± 0%	100% ± 0%	100% ± 0%
F ₆	100% ± 0%	100% ± 0%	100% ± 0%
F ₇	100% ± 0%	100% ± 0%	100% ± 0%
F ₈	100% ± 0%	100% ± 0%	100% ± 0%

Figure 1: Negative correlation between two allelic variables at a locus. (A) Qualitative reasoning for the expectation of negative correlation between two allelic variables at the same isolocus. (B) Use of Fisher's exact test to identify negative correlation between a pair of allelic variables. Ten individuals are shown for the sake of illustration, but an ideal dataset would have 100 or more individuals. In the allelic variables, presence of an allele in an individual is indicated by 1, and absence is indicated by 0.

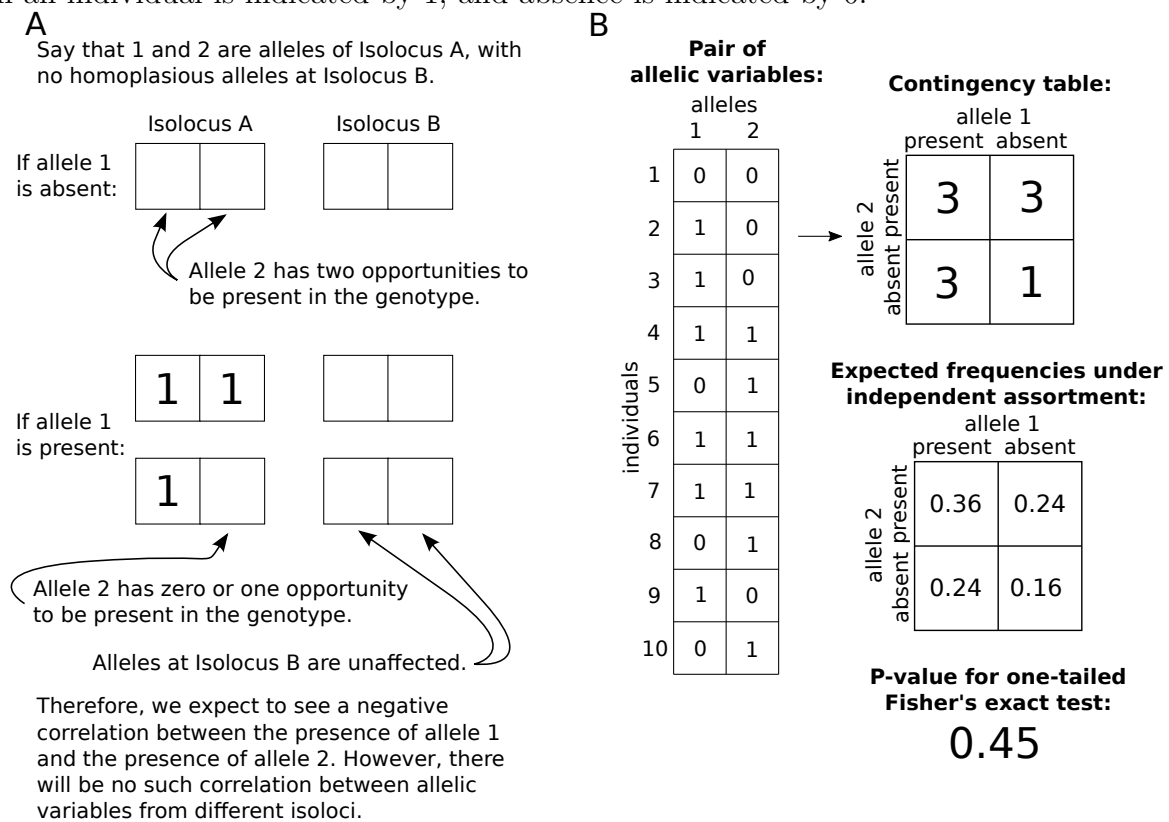


Figure 2: Overview of functions in POLYSAT 1.5 for processing allopolyploid and diploidized autopolyploid datasets. Additionally, the `processDatasetAllo` function can be used to automatically run `alleleCorrelations` and `testAlGroups` on every locus in a dataset. In the box representing the `alleleCorrelations` function, all alleles belonging to the locus on the left are variable in the dataset, so Fisher's exact test is used to find correlations between allelic variables, then K-means and UPGMA are used to perform clustering. The locus on the right has one allele (4) that is present in all individuals, making it impossible to assign alleles to isoloci using Fisher's exact test. In the box representing the `testAlGroups` function, all steps are performed on all loci regardless of whether or not fixed alleles are present.

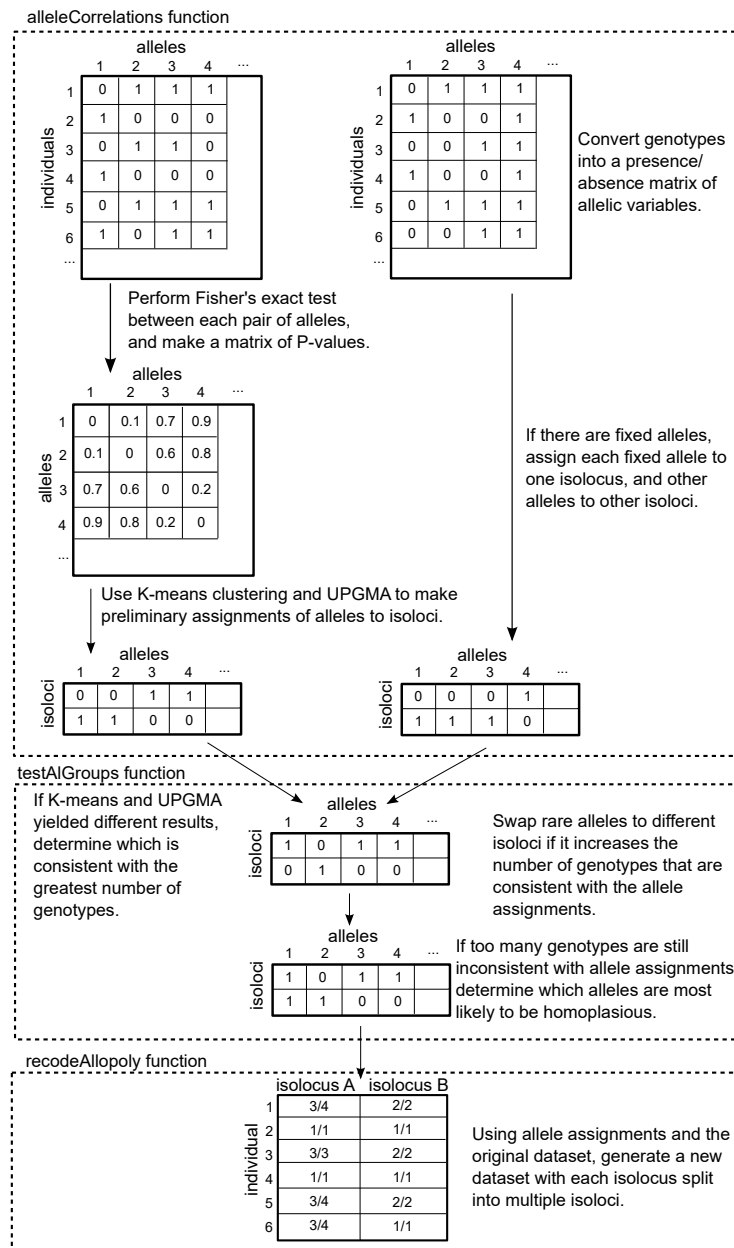


Figure 3: Accuracy of allele assignments with different sample sizes. For each ploidy and sample size, 1000 simulations were performed. Octoploids were simulated with two tetraploid genomes. Whiskers indicate 95% confidence intervals. “Swap ≤ 0.5 ” indicates that `testAllGroups` was used with `rare.al.check = 0.5`. The y-axis indicates the percentage of datasets for which allele assignments were completely correct.

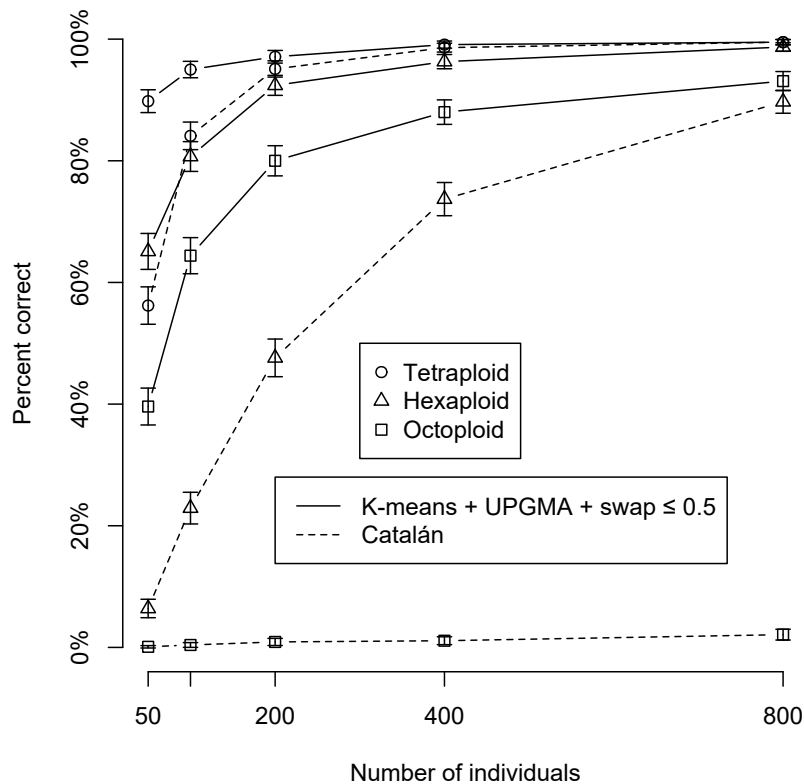


Figure 4: Percentages of simulated datasets with correct allele assignments when homoplasious alleles are present. Whiskers indicate 95% confidence intervals. The y-axis indicates the percentage of datasets for which allele assignments were completely correct. Allotetraploid datasets were simulated with one pair of homoplasious alleles (alleles from two different isoloci, but with identical amplicon size) for each locus. The frequency of homoplasious alleles was identical at both isoloci in each dataset, and was set at five different levels (0.1 through 0.5). Five different sample sizes were tested (50, 100, 200, 400, and 800). For each homoplasious allele frequency and sample size, 1000 datasets were simulated. Allele assignments were made using three methods: K-means + UPGMA (A; `rare.al.check = 0`), K-means + UPGMA + swap ≤ 0.25 (B; `rare.al.check = 0.25`), or K-means + UPGMA + swap ≤ 0.50 (C; `rare.al.check = 0.5`); plus an algorithm in the function `testAlGroups` that identifies the alleles most likely to be homoplasious, and assigns alleles as homoplasious until all genotypes are consistent with allele assignments.

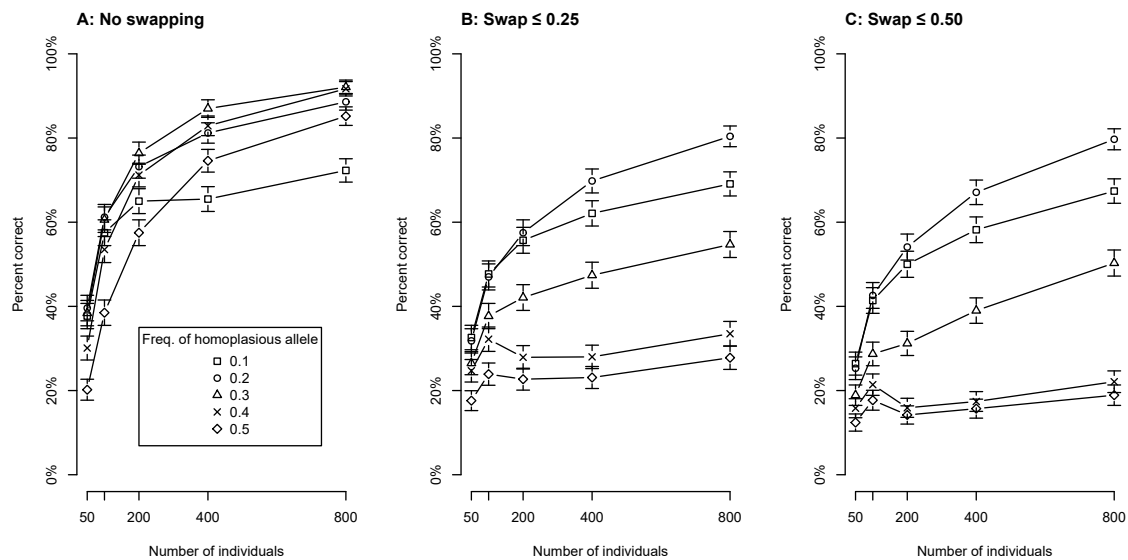


Figure 5: Percentages of simulated datasets with correct allele assignments when one isolocus has a null allele. Whiskers indicate 95% confidence intervals. The y-axis indicates the percentage of datasets for which allele assignments were completely correct. Allotetraploid datasets were simulated, and frequency of the null allele was set at one of five levels (x-axis). 1000 datasets were simulated at each null allele frequency. Two parameters for `testAlGroups` were adjusted: `rare.al.check` at values of zero, 0.25, and 0.5 (corresponding to the methods K-means + UPGMA, K-means + UPGMA + swap ≤ 0.25 , and K-means + UPGMA + swap ≤ 0.50 , respectively); and `null.weight` at values of zero (null alleles are allowed when checking for evidence of homoplasy) and 0.5 (genotypes lacking alleles belonging to a given isolocus are taken as evidence that their other alleles are homoplasious).

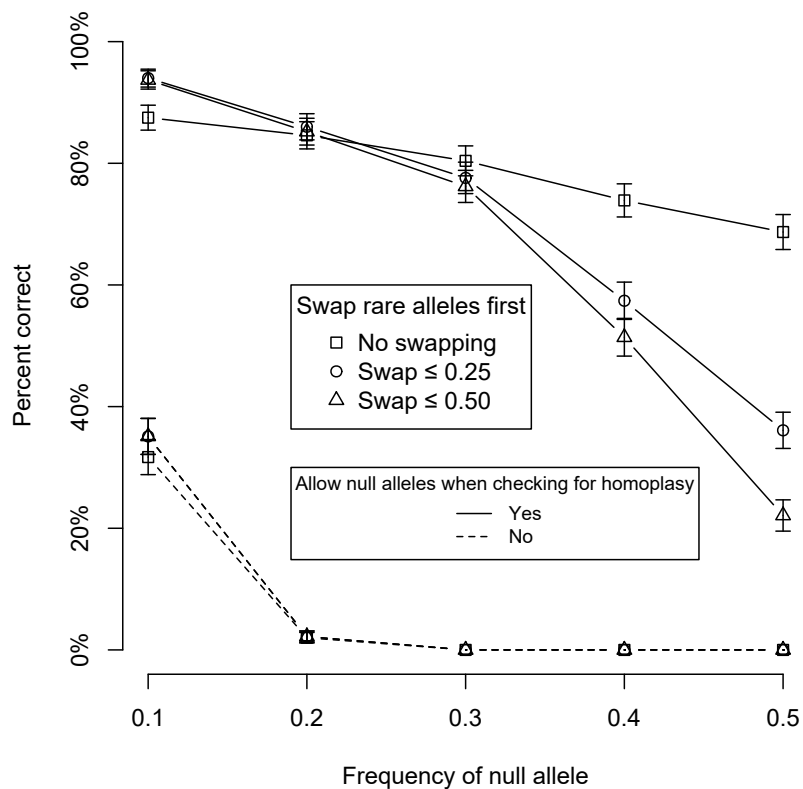


Figure 6: Percentages of simulated datasets with correct allele assignments when meiotic error causes compensated aneuploidy. Whiskers indicate 95% confidence intervals. The y-axis indicates the percentage of datasets for which allele assignments were completely correct. Meiotic error was simulated in the `simAllopolypoly` function on a per-gamete basis, with each error causing an allele from one isocus to be substituted with an allele from the other isocus. Each dataset was otherwise simulated for an allotetraploid organism with 100 individuals. Meiotic error rate, as shown in the x-axis, was controlled using the `meiotic.error.rate` argument of `simAllopolypoly`. For each error rate, 1000 datasets were simulated. For the `testAlGroups` function, the `tolerance` argument was set to 1 to prevent the function from checking for homoplasmy, and `rare.al.check` was set to zero, 0.25, or 0.5 (corresponding to the methods K-means + UPGMA, K-means + UPGMA + swap ≤ 0.25 , and K-means + UPGMA + swap ≤ 0.50 , respectively). Each dataset was tested for all three values of `rare.al.check`.

