

1 Resolving microsatellite genotype ambiguity in
2 populations of allopolyploid and diploidized
3 autopolyploid organisms using negative
4 correlations between alleles

5 Lindsay V. Clark
current email: lvclark@illinois.edu
permanent forwarding email: lindsay.v.clark.04@dartmouth.edu
Department of Crop Sciences
University of Illinois, Urbana-Champaign

6 Andrea Drauch Schreier
Department of Animal Science
University of California, Davis

7 May 25, 2015

8 **1 Abstract**

9 A major limitation in the analysis of genetic marker data from polyploid
10 organisms is non-Mendelian segregation, particularly when a single marker
11 yields allelic signals from multiple, independently segregating loci (isoloci).
12 However, with markers such as microsatellites that detect more than two al-
13 leles, it is sometimes possible to deduce which alleles belong to which isoloci.
14 Here we describe a novel mathematical property of codominant marker data
15 when it is recoded as binary (presence/absence) allelic variables: under ran-
16 dom mating in an infinite population, two allelic variables will be negatively
17 correlated if they belong to the same locus, but uncorrelated if they belong to
18 different loci. We present an algorithm to take advantage of this mathemat-
19 ical property, sorting alleles into isoloci based on correlations, then refining

20 the allele assignments after checking for consistency with individual geno-
21 types. We demonstrate the utility of our method on simulated data, as well
22 as a real microsatellite dataset from a natural population of octoploid white
23 sturgeon (*Acipenser transmontanus*). Our methodology is implemented in
24 the R package POLYSAT version 1.4.

25 **2 Introduction**

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

44 The situation is further complicated when not all copies of a locus pair
45 with each other with equal probability at meiosis (referred to as “disomic in-
46 heritance” assuming that the locus behaves as multiple independent diploid
47 loci [Obbard et al., 2006]; similarly, one could refer to an octoploid locus as
48 having “tetrasomic inheritance” if it behaved as two tetrasomic loci). In this
49 manuscript we will refer to duplicated loci that do not pair with each other
50 at meiosis (or pair infrequently) as “isoloci” after Obbard et al. [2006]. When
51 a genetic marker consists of multiple isoloci, it is not appropriate to analyze
52 that marker under the assumption of polysomic inheritance; for example, if
53 allele A can be found at both isoloci but allele B is only found at one isolocus
54 in a population, the genotypes AAAB and AABB are possible but ABBB is

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

76 Given that marker data do not follow straightforward Mendelian laws in
77 polyploid organisms, they are often recoded as a matrix of ones and zeros
78 reflecting the presence and absence of alleles (sometimes referred to as “al-
79 lelic phenotypes”) [Obbard et al., 2006]. In mapping populations such binary
80 data are useful if one parent is heterozygous for a particular allele and the
81 other parent lacks that allele, in which case segregation follows a 1:1 ratio and
82 can be analyzed under the diploid testcross model [Swaminathan et al., 2012,
83 Rousseau-Gueutin et al., 2008]. However, in natural populations, inheritance
84 of dominant (presence/absence) markers typically remains ambiguous, and
85 such markers are treated as binary variables that can be used to assess sim-
86 ilarity among individuals and populations but are inappropriate for many
87 population genetic analyses, *e.g.* tests that look for departures from or make
88 assumptions of Hardy-Weinberg Equilibrium [Clark and Jasieniuk, 2011].

89 Microsatellites are a special case given that they have multiple alleles,
90 allowing for the possibility of assigning alleles to isoloci, which would dras-
91 tically reduce the complexity of interpreting genotypes in allopolyploids and
92 diploidized autopolyploids. For example, if an allotetraploid individual has

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

119 In this manuscript, we present a novel methodology for assigning mi-
120 crosatellite alleles to isoloci based on the distribution of alleles among geno-
121 types in the dataset. Our method is appropriate for both mapping pop-
122 ulations and natural populations, as long as the dataset can be split into
123 reasonably-sized groups of individuals (~ 100 individuals or more) lacking
124 strong population structure. Negative correlations between alleles are used to
125 cluster alleles into putative isolocus groups, which are then checked against
126 individual genotypes. If necessary, alleles are swapped between clusters or
127 declared homoplasious so that the clusters agree with the observed genotypes
128 within a certain error tolerance. Genotypes can then be recoded, with each
129 marker split into two or more isoloci, such that isoloci can then be analyzed
130 as diploid or polysomic markers. Our method still works when there are null

131 alleles, homoplasmy between isoloci, or occasional meiotic recombination be-
132 tween isoloci, albeit with reduced power to find the correct set of allele assign-
133 ments. We test our methodology on simulated allotetraploid, allohexaploid,
134 and allo-octoploid (having two tetrasomic genomes) data, and compare its
135 effectiveness to that of the method of Catalán et al. [2006]. We also demon-
136 strate the utility of our method on a real dataset from a natural population
137 of octoploid white sturgeon (*Acipenser transmontanus*). Our methodology,
138 as well as a modified version of the Catalán et al. [2006] methodology, are
139 implemented in the R package POLYSAT version 1.4.

140 **3 Alleles in an unstructured population are** 141 **negatively correlated if they belong to the** 142 **same locus**

143 Say that a microsatellite dataset is recoded as an “allelic phenotype” matrix,
144 such that each row represents one individual, and each allele becomes a col-
145 umn (or an “allelic variable”) of ones and zeros indicating whether that allele
146 is present in that individual or not. Under Hardy-Weinberg equilibrium and
147 in the absence of linkage disequilibrium, these allelic variables are expected
148 to be uncorrelated if the alleles belong to different loci or different isoloci.
149 However, if two alleles belong to the same locus (or isolocus), the allelic vari-
150 ables should be negatively correlated. This is somewhat intuitive given that
151 the presence of a given allele means that there are fewer locus copies remain-
152 ing in which the other allele might appear (Fig. 1). The negative correlation
153 can also be proved mathematically.

154 Define a set of three or more alleles numbered $1 \cdots n$ that all belong to one
155 isolocus. k and j are alleles in this set. Allele frequencies in the population
156 being sampled are defined as $p_1 \cdots p_n$, where each allele frequency is between
157 one and zero, and all allele frequencies sum to one. G_k means that allele k
158 is present in a given individual’s genotype.

159 In a diploid, the probability that allele k is present in an individual is

$$P(G_k) = p_k^2 + 2p_k * (1 - p_k) = 2p_k - p_k^2 \quad (1)$$

160 and the probability that allele k is absent in an individual is

$$P(-G_k) = (1 - p_k)^2 \quad (2)$$

161 Given that allele j is present, the conditional probability that allele k is
162 also present is

$$P(G_k|G_j) = \frac{P(G_k \cap G_j)}{P(G_j)} = \frac{2p_k p_j}{2p_j - p_j^2} = \frac{2p_k}{2 - p_j} \quad (3)$$

163 and given that allele j is absent, the conditional probability that allele k
164 is present is

$$P(G_k|-G_j) = \frac{P(G_k \cap -G_j)}{P(-G_j)} = \frac{2p_k - p_k^2 - 2p_k p_j}{(1 - p_j)^2} \quad (4)$$

165 It follows that

$$P(G_k|-G_j) - P(G_k|G_j) = \frac{p_k[2(1 - p_k - p_j) + p_k p_j]}{(1 - p_j)^2(2 - p_j)} \quad (5)$$

166 Given that $0 < p_k < 1$, $0 < p_j < 1$, and $p_k + p_j < 1$, $P(G_k|-G_j) -$
167 $P(G_k|G_j) > 0$ and therefore $P(G_k|-G_j) > P(G_k|G_j)$. The magnitude of
168 difference between these probabilities is dependent on allele frequencies, but
169 the presence of allele j in a genotype always reduces the probability that
170 allele k is also present. A proof of the same principal at a tetrasomic locus
171 is provided in appendix 1.

172 4 Methods for clustering alleles into isoloci 173 using negative correlations

174 As demonstrated in the previous section, the occurrences of two different
175 alleles across individuals in a panmictic population should be negatively cor-
176 related if the two alleles belong to the same isolocus, and uncorrelated if
177 they belong to different isoloci. Therefore, tests for negative correlation can
178 be used to guide the assignment of alleles to isoloci. For any pair of alle-
179 les, a two-by-two contingency table can be generated, containing counts of
180 individuals that have both alleles, neither allele, the first allele but not the
181 second, and the second allele but not the first. A statistical test is then per-
182 formed to check for independence of the two allelic variables; we use Fisher's
183 exact test here because it is appropriate for small sample sizes, which are
184 likely to occur in typical population genetics datasets when rare alleles are

185 present. A one-tailed Fisher's exact test is used, with the alternative hy-
186 pothesis being that a disproportionately high number of individuals will just
187 have one allele of the pair, as opposed to both alleles or neither allele. This
188 alternative hypothesis corresponds to the two alleles belonging to the same
189 isolocus, whereas the null hypothesis is that they belong to different isoloci
190 and therefore assort independently. The P-values from Fisher's exact test
191 on each pair of allelic variables from a single microsatellite marker are then
192 stored in a symmetric square matrix. We expect to see clusters of alleles
193 with low P-values between them; alleles within a cluster putatively belong
194 to the same isolocus. For clustering algorithms, zeros are inserted along the
195 diagonal of the matrix, since the P-values are used as a dissimilarity statistic.
196 The function `alleleCorrelations` in POLYSAT 1.4 produces such a matrix
197 of P-values for a single microsatellite marker.

198 Multiple clustering algorithms exist for square matrices. Here we test K-
199 means and hierarchical clustering for their ability to correctly assign alleles
200 to isoloci based on P-values produced by Fisher's exact test. K-means clus-
201 tering places rows of a matrix into clusters, where each row is more similar
202 to the mean of its own cluster than to the mean of any other cluster. We
203 use the Hartigan and Wong [1979] method of K-means clustering as it is the
204 default method implemented in R, and we found that if we set the number of
205 randomly chosen starting centroids high enough (`nstart = 50` in the `kmeans`
206 function) it converged on the same answer as all other methods of K-means
207 clustering (data not shown). One potential issue with K-means clustering is
208 that since it only seeks to group similar rows together, it can group pairs of
209 alleles with high P-values rather than low P-values. Hierarchical clustering,
210 on the other hand, treats the values in the square matrix as a dissimilarity
211 statistic, and seeks to cluster objects that are most similar. We used three
212 methods of hierarchical clustering (complete linkage, single linkage, and UP-
213 GMA, which differ in how they use distances between elements of clusters
214 to determine distances between clusters) on simulated tetraploid, hexaploid,
215 and octoploid datasets, and found that UPGMA was the most accurate for
216 assigning alleles to isoloci for all ploidies (Table 1). K-means was more accu-
217 rate than UPGMA for all ploidies, and the method of Catalán et al. [2006]
218 had similar accuracy to K-means for allotetraploid datasets, but performed
219 much more poorly at higher ploidies (Table 1). Although K-means was more
220 accurate overall than UPGMA, UPGMA sometimes found the correct an-
221 swer when K-means found the incorrect answer, and therefore both results
222 are output by the `alleleCorrelations` function in POLYSAT. To choose be-

223 tween K-means and UPGMA when they give different results, the function
224 `testAllGroups` in POLYSAT checks every genotype in the dataset against both
225 results. A genotype is consistent with a set of allele assignments if it has at
226 least one allele belonging to each isolocus, and no more alleles belonging to
227 each isolocus than the ploidy of that isolocus (*e.g.* two in an allotetraploid).
228 The set of results that is consistent with the greatest number of genotypes
229 is selected, or K-means in the event of a tie. Selecting the best results out
230 of K-means and UPGMA improved the accuracy of allele assignments at all
231 ploidies, particularly hexaploids (Table 1).

232 When the best set of allele assignments was chosen from K-means and
233 UPGMA, 13% of tetraploid datasets, 31% of hexaploid datasets, and 41%
234 of octoploid datasets still had incorrect allele assignments (Table 1). We
235 expected that rare alleles would be the most likely to be assigned incorrectly,
236 given that they would be present in the fewest genotypes and therefore there
237 would be the least statistical power to detect correlations between them and
238 other alleles. To correct the allele assignments, an algorithm was added to
239 the `testAllGroups` function that individually swaps the assignment of each
240 rare allele to the other isolocus (or isoloci) and then checks whether the new
241 set of assignments is consistent with a greater number of genotypes than
242 the old set of assignments. If an allele is successfully swapped, then every
243 other rare allele is checked once again, until no more swaps are made. The
244 maximum number of genotypes in which an allele must be present to be
245 considered a rare allele is adjusted using the `rare.al.check` argument to
246 the `testAllGroups` function. On the same set of datasets that were used to
247 compare K-means and hierarchical clustering methods to the Catalán et al.
248 [2006] method, we tested the accuracy of allele assignments when alleles
249 present in $\leq 25\%$ or $\leq 50\%$ of genotypes were subjected to the swapping
250 algorithm (Table 1). Note that the frequency of genotypes with a given
251 allele will always be higher than the allele frequency itself, although a 50%
252 threshold is still much higher than the cutoff for considering an allele to be
253 “rare” in most population genetic analyses. For all ploidies, swapping rare
254 alleles between isoloci resulted in a considerable improvement in accuracy.

255 The accuracy of allele assignment was dependent on the number of alleles
256 belonging to the two isoloci in the simulated datasets. The assignments for
257 tetraploids were very accurate (98-100%) when both isoloci had the same
258 number of alleles, but decreased in accuracy as the difference in number of
259 alleles between the two isoloci increased (Table 2). When one isolocus had
260 two alleles and the other had eight alleles, the accuracy was 80% (Table 2).

261 The accuracy of allele assignments would be expected to increase with
262 sample size for our method, given that power to detect correlation is very
263 dependent on sample size. The Catalán et al. [2006] method should also
264 improve with increasing sample size, given that with more individuals in the
265 dataset, there is a greater probability of producing the double homozygotes
266 that are needed to resolve the allele assignments. For all ploidies, we found
267 that the accuracy of both our method and the Catalán et al. [2006] method
268 was dependent on sample size, and that our method performed better than
269 the Catalán et al. [2006] method at all sample sizes (Fig. 2). For tetraploids
270 and hexaploids, the effect of sample size was greater on the Catalán et al.
271 [2006] method than on our method, particularly at small sample sizes (Fig.
272 2). For octoploids, the success of the Catalán et al. [2006] method was near
273 zero even with 800 individuals in the dataset (due to the low probability of
274 producing fully homozygous genotypes at tetrasomic isoloci), whereas our
275 method had an accuracy of 93% with 800 octoploid individuals.

276 5 Caveats of the method

277 *Population structure.* Both negative and positive correlations between alleles
278 at different loci (or isoloci) can occur when the assumption of random mating
279 is violated by population structure, confounding the use of negative corre-
280 lations for assigning alleles to isoloci. We simulated allotetraploid datasets
281 consisting of two populations of fifty individuals differing in allele frequency
282 by a predetermined amount, and found that accuracy of our method remained
283 high ($\sim 90\%$) even at moderate levels of F_{ST} (~ 0.2 ; Table 3). Interestingly,
284 low levels of population structure ($F_{ST} \approx 0.02$) improved the accuracy of
285 our method to 99%, compared to 94% when $F_{ST} = 0$ (Table 3), probably
286 as a result of an increase in the number of double homozygous genotypes,
287 which would have been informative during the allele swapping step. For this
288 same reason, the Catalán et al. [2006] method, which depends on double ho-
289 mozygous genotypes, had an improved success rate as population structure
290 increased, and exceeded our method in accuracy at moderate levels of F_{ST}
291 (Table 3). However, accuracy of our method decreased with increasing F_{ST}
292 when $F_{ST} > 0.02$ (Table 3), likely because correlations between alleles caused
293 by population structure outweighed the benefits of increased homozygosity.
294 Because our method can be negatively impacted by population structure, the
295 `alleleCorrelations` function checks for significant positive correlations be-

296 tween alleles, which could only be caused by population structure or scoring
297 error (such as stutter peaks being mis-called as alleles, and therefore tend-
298 ing to be present in the same genotypes as their corresponding alleles), and
299 prints a warning if such correlations are found. In our simulations, signif-
300 icant positive correlations between alleles were found in most datasets that
301 had moderate population structure (Table 3).

302 Because the user may want to split highly structured datasets into multi-
303 ple populations for making allele assignments (to avoid the issue of allele cor-
304 relations caused by population structure), the function `mergeAlleleAssignments`
305 is included in POLYSAT 1.4 in order to merge several sets of allele assignments
306 into one. This is particularly useful if some alleles are found in some pop-
307 ulations but not others, or if alleles with identical amplicon sizes are found
308 belonging to different isoloci in different populations.

309 *Homoplasmy.* Although our algorithm attempts primarily to sort alleles
310 into non-overlapping groups, there is always a possibility that different isoloci
311 have some alleles with identical amplicon sizes. Therefore, we introduced an
312 algorithm to the `testAlGroups` function to check whether any genotypes
313 were still inconsistent with the allele assignments after the allele swapping
314 step, and assign alleles to multiple isoloci until all genotypes (or a particular
315 proportion that can be adjusted with the `threshold` argument) are consistent
316 with the allele assignments. The allele that could correct the greatest number
317 of inconsistent genotypes (or in the event of a tie, the one with the lowest P-
318 values from Fisher’s exact test between it and the alleles in the other isolocus)
319 is made homoplasious first, then all genotypes are re-checked and the cycle is
320 repeated until the desired level of agreement between allele assignments and
321 genotypes is met.

322 We tested the accuracy of allele assignments across several sample sizes
323 and frequencies of homoplasious alleles, with and without allele swapping be-
324 forehand (Fig. 3). Allele assignments were most accurate when allele swap-
325 ping was not performed before testing for homoplasious alleles, and when the
326 homoplasious allele was at a frequency of 0.3 in both isoloci. In addition to
327 the issue of assignment accuracy, there is the fact that many genotypes may
328 not be unambiguously resolvable when they contain a homoplasious allele.
329 For example, if alleles A and B belong to different isoloci, and C belongs to
330 both, the genotype ABC could be AA BC, AC BB, or AC BC. In such cases,
331 the POLYSAT function `recodeAlloPoly`, which recodes genotypes from single
332 loci into multiple isoloci, marks the entire genotype as missing. When allele
333 assignments were correct, we tested the mean proportion of genotypes that

334 were resolvable, given several frequencies of a homoplasious allele (Table 4).
335 Although accuracy of assignment had been highest with a homoplasious allele
336 frequency of 0.3, only 57% of genotypes could be resolved in such datasets
337 (Table 4).

338 Homoplasmy between alleles within an isolocus is also possible, meaning
339 that two or more alleles belonging to one isolocus are identical in amplicon
340 size but not identical by descent. Although such homoplasmy is an important
341 consideration for analyses that determine similarity between individuals and
342 populations, homoplasmy within isoloci does not affect the allele assignment
343 methods described in this manuscript.

344 *Null alleles.* Mutations in primer annealing sites are a common occur-
345 rence with microsatellite markers, and result in alleles that produce no PCR
346 product, known as null alleles. To test the effect of null alleles on the accuracy
347 of our allele assignment method, we simulated datasets in which one isolocus
348 had a null allele (Fig. 4). One potential issue with null alleles is that, when
349 homozygous, they can result in genotypes that do not appear to have any alle-
350 les from one isolocus. Such genotypes are used by the `testAlGroups` function
351 as an indicator that alleles should be swapped or made homoplasious, which
352 would be incorrect actions if the genotype resulted from a null allele rather
353 than inaccuracy of allele assignment. We therefore added an argument to
354 the `testAlGroups` function, `null.weight`, to indicate how genotypes with
355 no apparent alleles for one isolocus should be prioritized for determining
356 which alleles to make homoplasious. If null alleles are expected to be com-
357 mon, `null.weight` can be set to zero so that genotypes with no apparent
358 alleles for one isolocus are not used for assigning homoplasmy. The default
359 value of 0.5 for `null.weight` will cause `testAlGroups` to use genotypes with
360 no apparent alleles for one isolocus as evidence of homoplasmy, but with lower
361 priority than genotypes with too many alleles per isolocus. (No argument
362 was added to adjust the allele swapping algorithm, since it only swaps alleles
363 if the overall agreement with the dataset is improved.) We found that, when
364 null alleles were present, the accuracy of the algorithm was greatly improved
365 when genotypes lacking alleles for one isolocus were not used as evidence of
366 homoplasmy (Fig. 4). We also found that the allele swapping algorithm still
367 improved the accuracy of allele assignments, particularly when the null allele
368 was at a frequency greater than 0.25 in the population.

369 It is also possible for an entire isolocus to be null. This is often ap-
370 parent when a marker has fewer alleles per genotype than expected, *e.g.* a
371 maximum of two alleles per individual in a tetraploid. Such loci should be

372 excluded from the allele assignment analysis described in this manuscript.
373 If they are included in analysis accidentally, they can be identified by weak
374 K-means/UPGMA clustering of alleles (which can be evaluated from the
375 `betweenss`, `totss`, and `heatmap.dist` outputs from `alleleCorrelations`)
376 and by a high proportion of alleles appearing to be homoplasious.

377 *Meiotic error due to intermediate inheritance patterns.* As mentioned in
378 the introduction, although homologous chromosomes are likely to preferen-
379 tially pair with each other, occasional pairing between homeologous (in an
380 allopolyploid) or paralogous (in an autopolyploid) chromosomes may occur
381 during meiosis. As a result, offspring may be aneuploid, having too many or
382 two few chromosomes from either homologous pair, or may have translocati-
383 ons between homeologous or paralogous chromosomes. Most commonly, the
384 aneuploidy or translocations will occur in a compensated manner [Chester
385 et al., 2015], meaning that for a given pair of isoloci, the total number of
386 copies will be the same as in a non-aneuploid, but one isolocus will have
387 more copies than expected and the other isolocus will have fewer (*e.g.* three
388 copies of one isolocus and one copy of the other isolocus in an allotetraploid).

389 We simulated datasets in which gametes resulting in compensated ane-
390 uploidy (meiotic error) occurred at a range of frequencies from 0.01 to 0.2
391 (Fig. 5). At all meiotic error rates, the allele swapping algorithm from
392 `testAlGroups` improved the accuracy of allele assignment (Fig. 5). Meiotic
393 error did not have a large impact on the success of our method; even at a
394 meiotic error rate of 0.2 (where 0.5 would be fully autopolyploid), our algo-
395 rithm still had an accuracy of 62% on datasets of 100 individuals with no
396 homoplasmy, null alleles, or population structure (Fig. 5). Although the allele
397 swapping algorithm assumes fully disomic (or tetrasomic in the case of an
398 octoploid with two subgenomes) inheritance in the dataset, allele swapping
399 still improves accuracy of allele assignments when inheritance is not fully
400 disomic (Fig. 5).

401 To avoid calling alleles as homoplasious when meiotic error or intermedi-
402 ate inheritance causes some genotypes to have too many alleles per isolocus
403 (or none, as with null alleles), the `tolerance` argument was included for the
404 `testAlGroups` function. This argument indicates the proportion of geno-
405 types that can still be in disagreement with the allele assignments when the
406 algorithm stops adding homoplasious alleles. (For example, if `tolerance`
407 is zero, all genotypes must be consistent with the allele assignments, if
408 `tolerance` is one the algorithm will not check for homoplasmy at all, and
409 at the default of 0.05, the algorithm will keep adding homoplasious alleles

410 until no more than 5% of genotypes disagree with the allele assignments.)
411 Allowing for a few genotypes to disagree with the allele assignments is also
412 expected to make the algorithm more robust to occasional scoring error. The
413 function `recodeAllopolypoly` also has an argument, `allowAneuploidy`, to allow
414 for meiotic error when splitting markers into multiple isoloci. For genotypes
415 with too many alleles for one isolocus, the function will adjust the recorded
416 ploidy for the relevant samples and isoloci. (Ploidy is used by other POLYSAT
417 functions, such as those that estimate allele frequency.)

418 *Fixed alleles.* If one or more alleles are present in all genotypes in a
419 dataset, it is not possible to perform Fisher's exact test to look for correlations
420 between those fixed alleles and any others. The function `alleleCorrelations`
421 therefore checks for fixed alleles before performing Fisher's exact test. Each
422 fixed allele is assigned to its own isolocus. If only one isolocus remains, all
423 remaining alleles are assigned to it. If no isoloci remain (*e.g.* in an allote-
424 traploid with two fixed alleles and several variable alleles), then all remaining
425 alleles are assigned as homoplasious to all isoloci. If multiple isoloci remain
426 (*e.g.* in an allohexaploid with one fixed allele), then Fisher's exact test,
427 K-means clustering and UPGMA are performed to assign the alleles to the
428 remaining isoloci. It is possible that an allele with a very high frequency
429 may be present in all genotypes but not truly fixed (*i.e.* some genotypes
430 are heterozygous). However, allele swapping performed by `testAlGroups`
431 can assign alleles to an isolocus even if that isolocus already has an allele
432 assigned to it that is present in all individuals.

433 **6 Assignment of alleles to isoloci in octoploid** 434 **sturgeon**

435 To demonstrate the usefulness of our allele assignment method on a real
436 dataset, we used previously published data from natural populations of oc-
437 toploid white sturgeon (*Acipenser transmontanus*) [Drauch Schreier et al.,
438 2012]. Previous studies of inheritance patterns in this species suggested that
439 it possesses two tetrasomic subgenomes, at least for portions of its genome
440 [Rodzen and May, 2002, Drauch Schreier et al., 2011]. We selected for anal-
441 ysis the eight microsatellite markers that, based on number of alleles per
442 genotype, appeared to be present in eight copies rather than four.

443 Because population structure can impact allele clustering, we first per-

444 formed a preliminary analysis of population structure using a simple genetic
445 dissimilarity statistic and principal coordinates analysis (Fig. 6). The two
446 major genetic groups that we identified (Table 5) were similar to the pop-
447 ulation structure previously observed, which divides the Fraser River white
448 sturgeon population into two genetic groups that exist on either side of a
449 natural obstruction in the river [Drauch Schreier et al., 2012]. The smaller
450 group (Pop 2) consisted of only 66 individuals and, likely due to small sample
451 size, produced poor quality allele assignments with high levels of homoplasmy
452 when analyzed by itself (data not shown). We therefore tested our method
453 on Pop 1 (183 individuals) and on the combined set of 249 individuals. Be-
454 cause simulations indicated that the allele swapping algorithm reduced the
455 accuracy when there was homoplasmy in the dataset (Fig. 3), we tried allele
456 assignment with and without allele swapping. In checking for homoplasmy,
457 we allowed up to 5% of genotypes to disagree with allele assignments in an-
458 ticipation of meiotic error, scoring error, or genotypes homozygous for null
459 alleles (`tolerance = 0.05` in `testAlGroups`), and to allow for null alleles at
460 low frequency we set `null.weight = 0.5` so that genotypes with too many
461 alleles per isolocus would be used for assignment of homoplasmy first, before
462 genotypes with no alleles for one of their isoloci.

463 For four out of eight loci, our algorithm found allele assignments devoid
464 of homoplasmy when only Pop 1 was used for assignment and when the al-
465 lele swapping algorithm was used (Table 6). Eliminating the allele swapping
466 algorithm or using the whole dataset for allele assignment increased the num-
467 ber of apparent homoplasious alleles in most cases, and did not reduce the
468 number of apparent homoplasious alleles for any locus (Table 6). For the four
469 loci with homoplasmy, most genotypes in the dataset could not be assigned un-
470 ambiguously (Table 6). For the four loci with no apparent homoplasmy, nearly
471 all genotypes in Pop 1 could be assigned unambiguously, and approximately
472 three quarters of the genotypes in Pop 2 (which was not used for creating
473 the assignments) could be assigned unambiguously (Table 6). Despite the
474 fact that Pop 1 was previously determined to consist of three subpopula-
475 tions with Phi-PT values ranging from 0.06 to 0.17 [Drauch Schreier et al.,
476 2012], allele correlations resulting from population structure did not appear
477 to prevent us from obtaining reasonable assignments of alleles to isoloci. Sig-
478 nificant positive correlations between alleles were found at one and two out of
479 eight loci when Pop 1 and the whole dataset were used to make assignments,
480 respectively (data not shown).

481 7 Conclusions

482 Here we introduce the R package POLYSAT version 1.4, with several new func-
483 tions applicable to the analysis of allopolyploids and diploidized autopoly-
484 ploids. These include `simAllopoly`, which generates simulated datasets;
485 `catalanAlleles`, which uses the the Catalán et al. [2006] method to as-
486 sign alleles to isoloci; `alleleCorrelations`, which performs Fisher’s exact
487 test between each pair of alleles from a marker, and then uses K-means
488 clustering and UPGMA to make initial assignments of alleles to isoloci;
489 `testAlGroups`, which checks the consistency of allele assignments with indi-
490 vidual genotypes, chooses between the K-means and UPGMA method, swaps
491 alleles to different isoloci if it improves consistency, and identifies homopla-
492 sious alleles; `mergeAlleleAssignments`, which merges the allele assignments
493 from two different populations using the same microsatellite marker; and
494 `recodeAllopoly`, which uses allele assignments to recode the dataset, split-
495 ting each microsatellite marker into multiple isoloci. An overview of the data
496 analysis workflow is given in Fig. 7.

497 We found that, with simulated data, the accuracy of our allele assignment
498 algorithm was impacted by issues such as homoplasmy and null alleles, and that
499 the optimal parameters for the algorithm depended on which of these issues
500 were present in the dataset. This suggests, since most users will not know
501 whether their dataset has homoplasmy or null alleles, that the `testAlGroups`
502 function should initially be run with several different parameter sets, and for
503 each locus, the results with the fewest homoplasious alleles should be chosen.
504 A heatmap of the P-values generated from Fisher’s exact test can also serve
505 as a qualitative visual indicator of how well the alleles can be separated into
506 isolocus groups. We also found that, although our algorithm was negatively
507 impacted by meiotic error (pairing of non-homologous chromosomes during
508 meiosis) and moderate population structure, its accuracy remained fairly high
509 in both cases. Sensitivity to population structure is the biggest drawback of
510 our method in comparison to that of Catalán et al. [2006], which actually
511 has improved results as population structure increases. However, even low
512 frequencies of null alleles, homoplasmy, or meiotic error can prevent the method
513 of Catalán et al. [2006] from working entirely.

514 Using a real microsatellite dataset from natural populations of white stur-
515 geon, we found that our method was useful for recoding half of the markers
516 into two independently segregating loci each. Given that white sturgeon are
517 octoploid with two tetrasomic subgenomes [Drauch Schreier et al., 2011], we

518 expected this dataset to be problematic; having tetrasomic isoloci as opposed
519 to disomic isoloci would reduce the magnitude of the negative correlations
520 between alleles, and was observed in simulations to reduce the accuracy of
521 assignment using our method, although not nearly as severely as the reduction
522 in efficacy of the Catalán et al. [2006] method (Table 1, Fig. 2). In population
523 genetic studies, we expect that microsatellite markers that can be recoded
524 using our method could then be used for analyses requiring polysomic or di-
525 somic inheritance (for example, Structure [Falush et al., 2007], estimation of
526 allele frequency, or tests of Hardy-Weinberg Equilibrium), while the remain-
527 ing markers will still be useful for other analysis (for example, Mantel tests
528 using simple dissimilarity statistics). Additionally, we found that the allele
529 assignments that we made were still fairly useful for recoding genotypes in a
530 population that was not used for making the assignments.

531 Although inappropriate for biallelic marker systems such as single nu-
532 cleotide polymorphisms (SNPs) and dominant marker systems such as AFLPs,
533 the method that we have described could theoretically used to assign alle-
534 les to isoloci in any marker system in which multiple alleles are the norm.
535 Allozymes, although rarely used in modern studies, are one such system. Al-
536 though data from genotyping-by-sequencing (GBS, and the related technique
537 restriction site-associated DNA sequencing, or RAD-seq) are typically pro-
538 cessed to yield biallelic SNP markers, in the future as typical DNA sequencing
539 read lengths increase, it may become common to find multiple SNPs within
540 the physical distance covered by one read. In that case, haplotypes may be
541 treated as alleles, and negative correlations between haplotypes may be used
542 to assign them to isoloci.

543 **8 Obtaining polysat 1.4**

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

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

547 In the “doc” subdirectory of the package installation, PDF tutorials are
548 available for POLYSAT as a whole and for the methodology described in this
549 manuscript.

550 9 Supplementary files

- 551 • polysat_1.4-0.tar.gz: Source package for POLYSAT 1.4.
- 552 • polysat_1.4-0.zip: Microsoft Windows installation file for POLYSAT 1.4.
- 553 • allopolyVignette.pdf: Tutorial for creating and using allele assignments
554 in POLYSAT.
- 555 • tables_figs.R, sturgeon.R: R scripts for reproducing the analyses in
556 this manuscript.
- 557 • sturgeon.csv: White sturgeon microsatellite dataset used in sturgeon-
558 test.R.
- 559 • appendix1.pdf: Proof that alleles belonging to the same tetrasomic
560 locus are negatively correlated in an infinite, unstructured population.

561 10 References

562 References

- 563 P. Catalán, J. G. Segarra-Moragues, M. Palop-Esteban, C. Moreno, and
564 F. González-Candelas. A Bayesian approach for discriminating among
565 alternative inheritance hypotheses in plant polyploids: the allotetraploid
566 origin of genus *Bordera* (Dioscoreaceae). *Genetics*, 172(3):1939–1953, 2006.
- 567 M. Chester, R. K. Riley, P. S. Soltis, and D. E. Soltis. Patterns of chromoso-
568 mal variation in natural populations of the neoallotetraploid *Tragopogon*
569 *mirus* (Asteraceae). *Heredity*, 114(3):309–317, 2015.
- 570 L. V. Clark and M. Jasieniuk. Polysat: an R package for polyploid mi-
571 crosatellite analysis. *Molecular Ecology Resources*, 11(3):562–566, 2011.
- 572 H. N. De Silva, A. J. Hall, E. Rikkerink, M. A. McNeilage, and L. G. Fraser.
573 Estimation of allele frequencies in polyploids under certain patterns of
574 inheritance. *Heredity*, 95(4):327–334, 2005.

- 575 A. Drauch Schreier, D. Gille, B. Mahardja, and B. May. Neutral mark-
576 ers confirm the octoploid origin and reveal spontaneous autopolyploidy in
577 white sturgeon, *Acipenser transmontanus*. *Journal of Applied Ichthyology*,
578 27(Suppl. 2):24–33, 2011.
- 579 A. Drauch Schreier, B. Mahardja, and B. May. Hierarchical patterns of pop-
580 ulation structure in the endangered fraser river white sturgeon (*Acipenser*
581 *transmontanus*) and implications for conservation. *Canadian Journal of*
582 *Fisheries and Aquatic Sciences*, 69(12):1968–1980, 2012.
- 583 F. Dufresne, M. Stift, R. Vergilino, and B. K. Mable. Recent progress and
584 challenges in population genetics of polyploid organisms: an overview of
585 current state-of-the-art molecular and statistical tools. *Molecular Ecology*,
586 23(1):40–69, 2014.
- 587 D. Falush, M. Stephens, and J. K. Pritchard. Inference of population struc-
588 ture using multilocus genotype data: dominant markers and null alleles.
589 *Molecular Ecology Notes*, 7(4):574–578, 2007.
- 590 R. T. Gaeta and J. C. Pires. Homeologous recombination in allopolyploids:
591 the polyploid ratchet. *New Phytologist*, 186(1):18–28, 2010.
- 592 T. R. Gregory and B. K. Mable. Polyploidy in animals. In T. R. Gregory,
593 editor, *The Evolution of the Genome*, chapter 8, pages 427–517. Elsevier,
594 San Diego, 2005.
- 595 J. A. Hartigan and M. A. Wong. A K-means clustering algorithm. *Journal of*
596 *the Royal Statistical Society, Series C (Applied Statistics)*, 28(1):100–108,
597 1979.
- 598 D. J. Obbard, S. A. Harris, and J. R. Pannell. Simple allelic-phenotype
599 diversity and differentiation statistics for allopolyploids. *Heredity*, 97(4):
600 296–303, 2006.
- 601 J. A. Rodzen and B. May. Inheritance of microsatellite loci in white sturgeon
602 (*Acipenser transmontanus*). *Genome*, 45(6):1064–1076, 2002.
- 603 M. Rousseau-Gueutin, E. Lerceteau-Köhler, L. Barrot, D. J. Sargent,
604 A. Monfort, D. Simpson, P. Arús, G. Guérin, and B. Denoyes-Rothan.
605 Comparative genetic mapping between octoploid and diploid *Fragaria*

- 606 species reveals a high level of colinearity between their genomes and the es-
607 sentially disomic behavior of the cultivated octoploid strawberry. *Genetics*,
608 179(4):2045–2060, 2008.
- 609 K. Swaminathan, W. B. Chae, T. Mitros, K. Varala, L. Xie, A. Barling,
610 K. Glowacka, M. Hall, S. Jezowski, R. Ming, M. Hudson, J. A. Juvik,
611 D. S. Rokshar, and S. P. Moose. A framework genetic map for *Miscant-*
612 *hus sinensis* from RNAseq-based markers shows recent tetraploidy. *BMC*
613 *Genomics*, 13:142, 2012.
- 614 J. A. Udall and J. F. Wendel. Polyploidy and crop improvement. *Crop*
615 *Science*, 46(S1):S3–S14, 2006.
- 616 R. S. Waples. Estimation of allele frequencies at isoloci. *Genetics*, 118(2):
617 371–384, 1988.

618 11 Tables and Figures

Clustering method	4x =	6x =	8x =
	2x + 2x	2x + 2x + 2x	4x + 4x
Catalán	83.3%	21.0%	0.4%
K-means	84.0%	59.0%	57.9%
Complete linkage	59.1%	32.9%	25.9%
Single linkage	34.3%	12.5%	6.8%
UPGMA	68.8%	49.7%	33.3%
K-means + UPGMA	86.6%	68.6%	59.3%
K-means + UPGMA + swap ≤ 0.25	95.0%	83.7%	62.8%
K-means + UPGMA + swap ≤ 0.50	95.0%	84.3%	64.4%

Table 1: Percentages of simulated datasets with correct assignments of alleles to isoloci using different clustering methods. For each type of ploidy tested, including tetraploid ($4x = 2x + 2x$), hexaploid ($6x = 2x + 2x + 2x$), and octoploid ($8x = 4x + 4x$), 10,000 datasets were simulated with 100 individuals each, with two to eight alleles at each isolocus and randomly generated allele frequencies. Data were simulated without homoplasmy, null alleles, population structure, or meiotic error. Datasets were simulated using the `simAllopolypoly` function in POLYSAT 1.4. The Catalán et al. [2006] method, as implemented in the `catalanAlleles` function in POLYSAT 1.4, was used to attempt to assign alleles to isoloci by examining fully-homozygous genotypes. Note that if the Catalán method fails to find the correct answer, it will not produce any results, unlike the other methods listed in this table. For all other methods in the table, a square matrix was calculated for each dataset, containing P-values from Fisher’s exact test for negative correlation between each pair of alleles, using the R function `fisher.test` with `alternative = "less"`. K-means clustering was performed using the Hartigan and Wong [1979] method as implemented in the R function `kmeans`, with 50 random sets chosen (`n.start = 50`) to ensure convergence on the most optimal set of clusters. Hierarchical clustering was performed in R with the `hclust` (`method = "complete"` for complete linkage, `method = "single"` for single linkage, and `method = "average"` for UPGMA) and `cutree` functions. For “K-means + UPGMA”, both methods were performed, and if they produced different results, the results that were consistent with the greatest number of genotypes were retained. For the “swap ≤ 0.25 ” and “swap ≤ 0.50 ” methods, alleles occurring in $\leq 25\%$ or $\leq 50\%$, respectively, of individuals were tested to see if agreements between assignments and genotypes were improved when those alleles were individually swapped to different isoloci, in which case the swaps were retained. In POLYSAT 1.4, the Fisher’s exact test, K-means clustering, and UPGMA steps are included in the `alleleCorrelations` function, whereas the comparison of K-means and UPGMA results as well as allele swapping are implemented in the `testAlGroups` function.

	2	3	4	5	6	7	8
2	98.0%	97.0%	94.9%	89.0%	90.8%	83.7%	80.1%
3		100.0%	98.8%	97.5%	94.0%	93.3%	89.3%
4			98.5%	98.2%	96.6%	95.7%	94.6%
5				99.0%	99.0%	96.8%	96.6%
6					98.0%	97.8%	98.3%
7						98.2%	99.1%
8							97.9%

Table 2: Percentages of allele assignments for allotetraploid datasets from Table 1 that were correct using K-means + UPGMA + swap ≤ 0.50 , by the number of alleles at each of two isoloci.

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

Table 3: Percentages of simulated datasets with correct allele assignments under different levels of population structure. Allotetraploid datasets were simulated as in Table 1, but instead of one population of 100 individuals, two populations of 50 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 $(H_T - H_S)/H_T$, and means and standard deviations across 1000 simulations are shown. 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 two methods described in Table 1.

Freq. of homoplasious allele	Mean percentage of genotypes that could be resolved
0.1	87.4%
0.2	73.0%
0.3	57.0%
0.4	43.9%
0.5	40.9%

Table 4: For datasets from Fig. 3 with correct allele assignments, percentages of genotypes that could be unambiguously resolved.

Sampling region	Pop 1	Pop 2
SG-1	3	35
SG-2	16	22
SG-3	39	1
UFR	46	1
NKO	46	4
SL	26	3
FL	7	0

Table 5: Distribution of sampling regions of white sturgeon among two groups of individuals (Pop 1 and Pop 2) determined by principal coordinate analysis (Fig 6). Region codes refer to those used by [Drauch Schreier et al., 2012]. Pop 1 was used for assigning alleles to isoloci.

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	58%, 83%	59%, 77%
AciG35	18	2	0	1	0	0%, 1%	23%, 23%
Atr109	25	6	3	4	2	58%, 63%	55%, 52%
Atr117	22	1	1	0	0	0%, 0%	35%, 35%
AciG52	22	4	1	1	0	0%, 1%	24%, 24%
Atr107	24	3	1	2	1	64%, 66%	53%, 53%
Atr1173	18	3	2	3	2	61%, 76%	61%, 86%

Table 6: Assignment of alleles from eight microsatellite markers to two tetrasomic genomes in octoploid white sturgeon (*Acipenser transmontanus*). Alleles were assigned using the K-means + UPGMA method from Table 1, 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”, Table 5 and Fig. 6). The assignments from Pop 1 with Swap ≤ 0.5 were then used to split the dataset into isoloci using the `recodeAllopoly` function. Genotypes that could not be unambiguously determined were coded as missing data; percentages of missing data in Pop 1 and Pop 2 are shown.

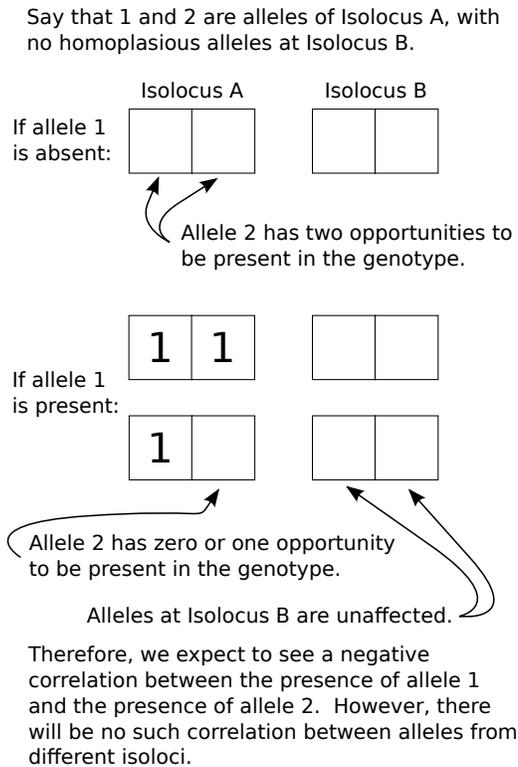


Figure 1: Qualitative reasoning for the expectation of negative correlation between two alleles at the same isolocus.

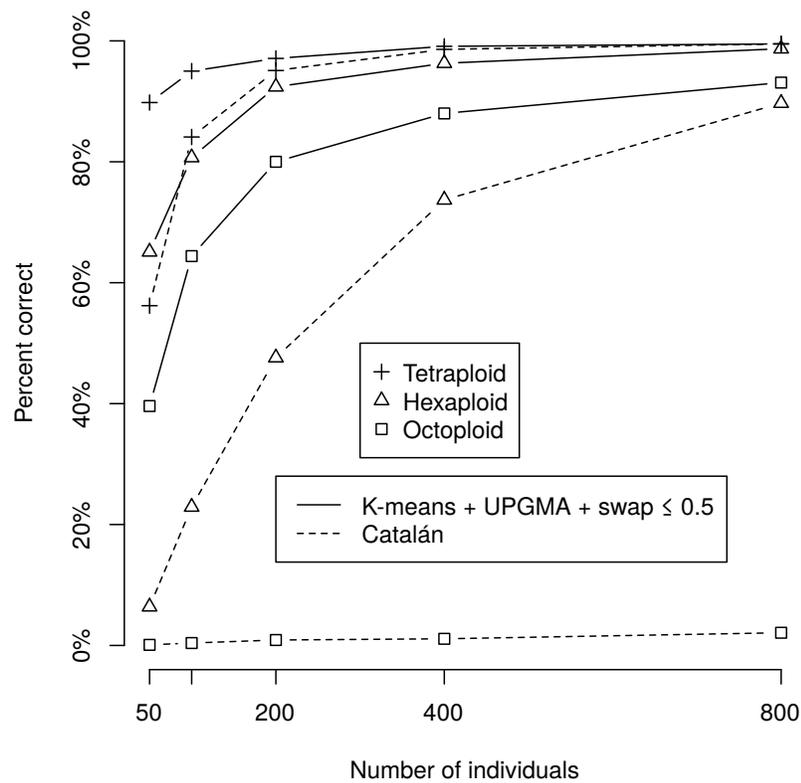


Figure 2: Accuracy of allele assignments with different sample sizes. For each ploidy and sample size, 1000 simulations were performed. Octoploids were simulated with two tetraploid genomes. Simulations and allele assignments were performed as in Table 1.

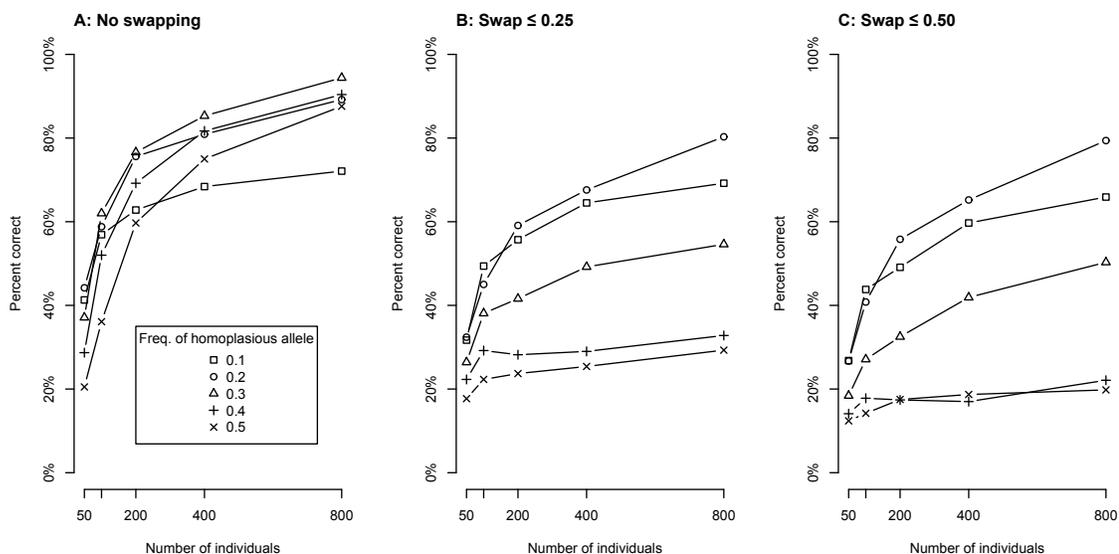


Figure 3: Percentages of simulated datasets with correct allele assignments when homoplasious alleles are present. Allotetraploid datasets were simulated as in Table 1, but in each dataset one pair of homoplasious alleles (alleles from two different isoloci, but with identical amplicon size) was simulated. 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 from Table 1: K-means + UPGMA (A), K-means + UPGMA + swap ≤ 0.25 (B), or K-means + UPGMA + swap ≤ 0.50 (C); plus an algorithm in the function `testAllGroups` that identifies the alleles most likely to be homoplasious, and assigns alleles as homoplasious until all genotypes are consistent with allele assignments.

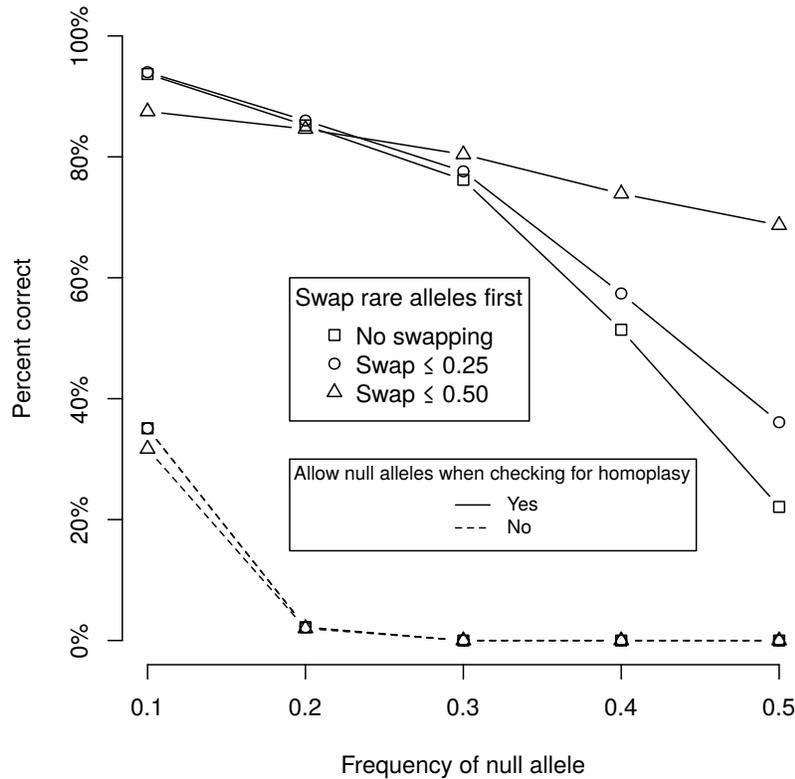


Figure 4: Percentages of simulated datasets with correct allele assignments when one is locus has a null allele. Allotetraploid datasets were simulated as in Table 1, and frequency of the null allele was set at one of five levels (x-axis). 1000 datasets were simulated at each null allele frequency. Two 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 \leq 0.25, and K-means + UPGMA + swap \leq 0.50, respectively); and `null.weight` at values of zero (null alleles are allowed when checking for evidence of homoplasmy) and 0.5 (genotypes lacking alleles belonging to a given is locus are taken as evidence that their other alleles are homoplasious).

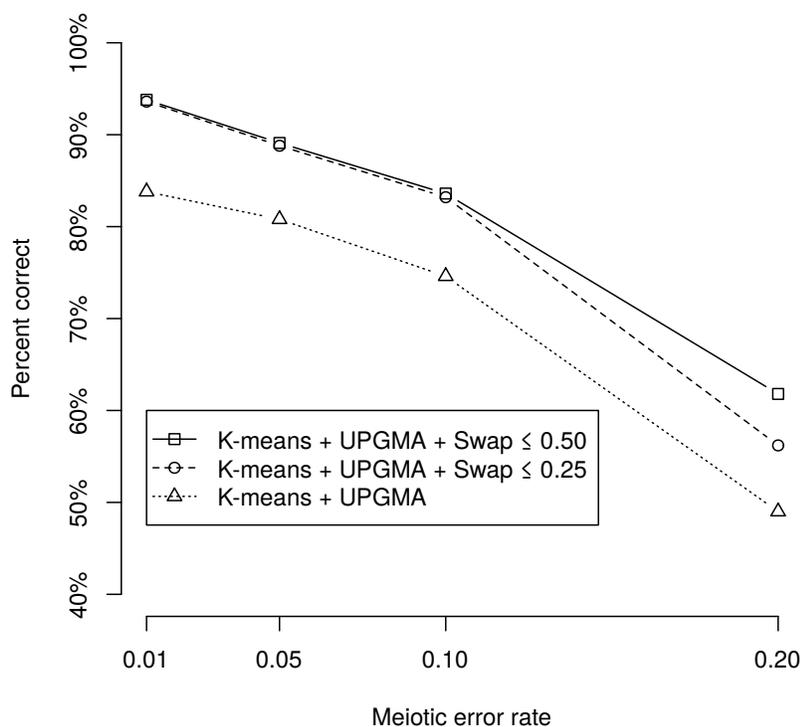


Figure 5: Percentages of simulated datasets with correct allele assignments when meiotic error causes compensated aneuploidy. Meiotic error was simulated in the `simAllopolypoly` function on a per-gamete basis, with each error causing an allele from one is locus to be substituted with an allele from the other is locus. Each dataset was otherwise simulated for an allotetraploid organism with 100 individuals as in Table 1. 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`.

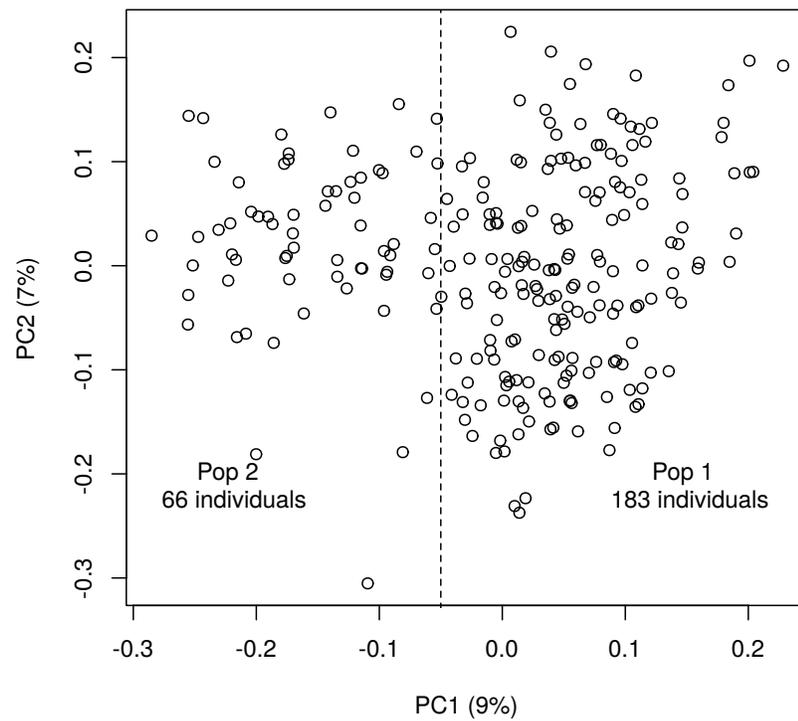


Figure 6: Principal coordinates analysis of 249 white sturgeon individuals, based on genotypes at eight microsatellite loci from Drauch Schreier et al. [2012]. Inter-individual distances were calculated using the `Lynch.distance` function in `POLYSAT`. Percentages of variation explained by the first two axes are shown. The dashed line indicates the cutoff for dividing the set into two groups (Pop 1 and Pop 2) based on population structure.

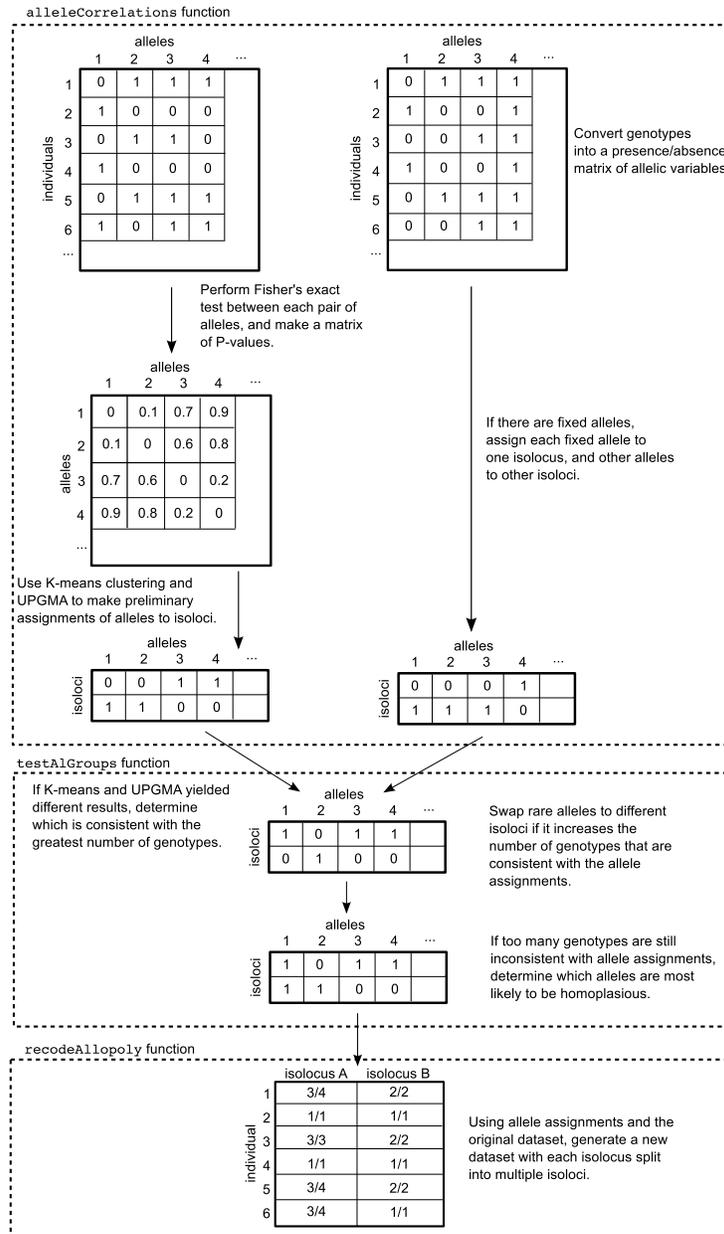


Figure 7: Overview of functions in POLYSAT 1.4 for processing allopolyploid and diploidized autopolyploid datasets.