

1 **Title**

2 VCFPOP: performing population genetics analyses for polyploids and anisoploids based
3 on next-generation sequencing variant calling dataset

4 **Authors**

5 Kang Huang^{1,2}, Bing Yang¹, Jincuo Ao¹, Yuhang Li¹, Yunxia Cui¹, Yuchen Kong¹, Yifan Wu¹,
6 Derek W. Dunn¹, Baoguo Li¹

7 **Addresses**

8 ¹ Shaanxi Key Laboratory for Animal Conservation, College of Life Sciences, Northwest Uni-
9 versity, Xi'an 710069, China

10 ² Department of Forest and Conservation Sciences, University of British Columbia, Vancou-
11 ver, BC V6T1Z4, Canada

12 **Keywords**

13 Polysomic inheritance, next-generation sequencing data, variant calling format, population
14 genetics.

15 **Corresponding author**

16 Baoguo Li

17 Telephone: +8613572209390; Fax: +86 029 88303304; E-mail: baoguoli@nwu.edu.cn

18 **Running title**

19 Population genetics analysis for autopolyploids
20

21 **Abstract**

22 Polyploids are cells or organisms with a genome consisting of more than two sets of
23 homologous chromosomes. Polyploid plants have important traits that facilitate speciation
24 and are thus often model systems for evolutionary, molecular ecology and agricultural
25 studies. However, due to their unusual mode of inheritance and double-reduction, diploid
26 models of population genetic analysis cannot properly be applied to polyploids. To over-
27 come this problem, we developed a software package entitled VCFPOP to perform a variety
28 of population genetic analyses for autopolyploids, such as parentage analysis, analysis of
29 molecular variance, principal coordinates analysis, hierarchical clustering analysis and
30 Bayesian clustering. We make this software freely available, downloadable from
31 <http://github.com/huangkang1987/vcfpop>.

32 **Keywords:** polysomic inheritance, next-generation sequencing data, population genetics,
33 AMOVA, Bayesian clustering, *F*-statistics.

34

35 Introduction

36 Polyploids are cells or organisms with a genome consisting of more than two sets of
37 homologous chromosomes. Polyploids represent a significant portion of all plant species,
38 with from 30-80% of angiosperms showing polyploidy (Burow et al. 2001). Because of their
39 propensity to facilitate speciation, polyploid species have often been used as model sys-
40 tems for evolutionary, molecular ecology and agricultural studies. More recently, poly-
41 ploids have increasingly become the focus of theoretical and experimental work (Avni et
42 al. 2017; Ling et al. 2018).

43 There are two distinct mechanisms of genome duplication that result in polyploidy:
44 allopolyploidy and autopolyploidy. This paper focuses on autopolyploids that display
45 polysomic inheritance. In autopolyploids, more than two homologous chromosomes can
46 pair at meiosis, resulting in the formation of multivalents and polysomic inheritance. A
47 typical feature of polysomic inheritance is the possibility that a gamete inherits a single
48 gene copy twice, termed double-reduction (Butruille & Boiteux 2000). For example, an au-
49 totetraploid individual $ABCD$ produces a gamete AA . The double-reduction will change
50 the frequency of genotypes within a population (Huang et al. 2019), resulting in increased
51 homozygosity and an inflated inbreeding coefficient (Hardy 2016).

52 Due to differences in data format and modes of inheritance between diploids and pol-
53 yploids, population genetics software designed for diploid organisms such as GENEPOP
54 (Rousset 2008) and ARLEQUIN (Excoffier & Lischer 2010) cannot be used for autopolyploids.
55 Some software packages have been developed in order to accommodate polyploid

56 genotype datasets, e.g. POLYSAT (Clark & Jasieniuk 2011), SPAGEDI (Hardy & Vekemans
57 2002), POLYRELATEDNESS (Huang *et al.* 2014), GENODIVE (Meirmans & Tienderen 2004) and
58 STRUCTURE (Pritchard *et al.* 2000), etc. However, such software includes the assumption of
59 a disomic mode of inheritance and genotypic frequencies in accordance with the *Hardy-*
60 *Weinberg Equilibrium* (HWE), whereas alleles of the same genotype are assumed to be in-
61 dependent.

62 Huang *et al.* (2019) derived the genotypic frequency for various double-reduction
63 models: the *random chromosome segregation* (RCS) (Muller 1914), the *pure random chromatid*
64 *segregation* (PRCS) (Haldane 1930), the *complete equational segregation* (CES) (Mather 1935)
65 and the *partial equational segregation* (PES) (Huang *et al.* 2019). The software package
66 POLYGENE is able to use all of these double-reduction models and can perform population
67 genetics analyses for both allelic phenotypic and genotypic data (Huang *et al.* 2020). How-
68 ever, POLYGENE cannot accommodate the large datasets that result from the use of next-
69 generation sequencing (NGS) methods. To solve this problem, we developed a new soft-
70 ware package entitled VCFPOP.

71 **The new software package VCFPOP**

72 VCFPOP is a free software developed with C++. It works only via command-line mode
73 and will run on Windows, Linux and Mac OS X. To ensure free copying, distribution and
74 modifications of the software and its source code, VCFPOP is distributed under a GNU Gen-
75 eral Public License (GPL, version 3).

76 VCFPOP has been optimized for memory allocation and calculation speed to analyze

77 large genotype datasets. It can analyze the genotypes of haploids, diploids, polyploids and
78 anisoploids. For polyploids, VCFPOP applies the RCS, PRCS, CES and PES double-reduc-
79 tion models and supports a maximum ploidy level of 10. VCFPOP also supports multi-level
80 region definition, so as to analyze the variance component of different hierarchies via *anal-*
81 *ysis of molecular variance* (AMOVA) (Huang *et al.* 2021).

82 **Input format**

83 VCFPOP can handle a multiple genotype format, consisting of *variant call format* (VCF)
84 V4.x (compressed or uncompressed) (Danecek *et al.* 2011), *binary call format* (BCF) V2.x,
85 GENEPOP V4.3 (Rousset 2008), STRUCTURE V2.3 (Pritchard *et al.* 2000), CERVUS V3.0
86 (Kalinowski *et al.* 2007), ARLEQUIN V3.6 (Excoffier & Lischer 2010), POLYGENE V1.4 (Huang
87 *et al.* 2020), and POLYRELATEDNESS V1.7 (Huang *et al.* 2015a).

88 Multiple VCF/BCF files for the same samples sequenced at different variants (vertical
89 concatenation, separated by '|' in `-g_input`), or different samples sequenced at the same
90 variants (horizontal concatenation, separated by '&' in `-g_input`) can be analyzed together
91 without additional concatenation.

92 Because both VCF and BCF formats do not contain population or regionally differen-
93 tiation, populations and regions should be additionally defined using the arguments, in
94 which '`-g_indfile`' reads the content from a file and '`-g_indtext`' reads from the com-
95 mand. Because the command-line mode does not allow linebreaks, '#n' or space is used as
96 the escape character in `-g_indtext`. Multi-level region definition is supported in all cor-
97 responding analysis (e.g., AMOVA and population assignment) (Huang *et al.* AMOVA

98 paper). The example format of `-g_indfile` is shown as follows:

```
99     -g_indtext="pop1:ind1,ind2,ind3  pop2:#4,#5-#6  pop3:ind7,ind8,ind9
100     #REG A1:#1,#2 A2:pop3 #REG B1:#1-#2"
```

101 The identifier of a population is followed by a semi-colon, then the individual identifiers,
102 ordinations or ordination ranges separated by commas. The separator ‘#REG’ is used to
103 separate regions or populations at different levels. Similarly, the identifier of a region is
104 followed by a semi-colon, then the sub-region or population identifiers.

105 Usage

106 After installation, the user should open the terminal to launch VCFPOP.

107 To view the help for all functions, execute

```
108 ./vcfpop -h
```

109 To view the detail help information of some specific functions, execute

```
110 ./vcfpop -h -func1 -func2
```

111 To use specific functions, execute

```
112 ./vcfpop -func1 -func1_parameters -func2 -func2_parameters ...
```

113 To use a parameter file in the same format as the program arguments but allow line
114 breaks, execute

```
115 ./vcfpop -p=parameter_file
```

116 After calculation, the results are saved as ‘*.func.txt’. The specific functions are as
117 follows:

```
118     -g                               General settings
```

119	-f	Filter for individual, locus or genotype
120	-haplotype	Haplotype extraction
121	-convert	File conversion
122	-diversity	Genetic diversity indices
123	-indstat	Individual statistics
124	-fst	Genetic differentiation
125	-gdist	Genetic distance
126	-amova	Analysis of molecular variance
127	-popas	Population assignment
128	-relatedness	Relatedness coefficient
129	-kinship	Kinship coefficient
130	-pcoa	Principal coordinate analysis
131	-cluster	Hierarchical clustering
132	-structure	Bayesian clustering

133 **Functions**

134 **General settings:** configuration of input and output files, output format (e.g., scien-
135 tific notation, decimal places), temporary directory, sampling population for individuals,
136 region definitions, number of threads, *single instruction multiple data* (SIMD) instructions
137 (e.g., SSE, AVX, AVX512), and random number generator seed.

138 **Filter:** exclusion of some variants of low quality, genotyping ratio, or polymorphism,
139 individuals with poor genotyping ratio, and genotypes of low quality or read depth before

140 file conversion and analyses. There are four types of filters: (i) variant information filters
141 that exclude variants by their quality, type (e.g., single nucleotide polymorphism or indel)
142 and original filter; (ii) genotype filters that exclude genotypes by their read depths, geno-
143 type qualities and ploidy levels. The excluded genotype will be set as the missing genotype;
144 (iii) individual filters that exclude individuals by the number of variants typed and ploidy
145 levels; and (iv) diversity filters that exclude variants based on the estimated genetic diver-
146 sity indices (e.g., minor allele frequency, number of alleles, number of individuals geno-
147 typed, heterozygosity, significance of genotypic equilibrium test) in a specified reference
148 population.

149 **Haplotype extraction:** several adjacent variants are combined into a highly polymor-
150 phic locus and the haplotypes are subsequently extracted. The haplotypes are used as al-
151 leles in subsequent analyses.

152 **Conversion:** genotypes can be converted into another genotype format, either
153 GENEPOP (Rousset 2008), SPAGEDI (Hardy & Vekemans 2002), CERVUS (Kalinowski et al.
154 2007) , ARLEQUIN (Excoffier & Lischer 2010), STRUCTURE (Pritchard et al. 2000), POLYGENE
155 (Huang et al. 2020), or POLYRELATEDNESS (Huang et al. 2015a).

156 **Genetic diversity indices:** this estimates the genetic diversity indices (e.g., observed
157 and expected heterozygosity, effective number of alleles and polymorphic information
158 content, Shannon's information index, inbreeding coefficient). A genotypic distribution
159 test (i.e., the HWE test in polyploids) is performed using a Fisher's G-test, with the null
160 hypothesis being the genotypic frequencies are in accordance with the prediction of a spe-
161 cific double-redouble model (e.g., RCS, PRCS, CES, PES).

162 **Individual statistics:** this enables the estimation of individual heterozygosity, the
163 genotype likelihood, inbreeding coefficient and kinship coefficient. In polyploids, the het-
164 erozygosity of a genotype is the probability of randomly sampling two different *identical-*
165 *by-state* (IBS) alleles without replacement (e.g., 2/3 for a tetraploid genotype *AABB*). The
166 individual heterozygosity is the arithmetic average of the heterozygosity of genotypes of
167 an individual across loci. The genotypic likelihood is the product of genotypic frequencies
168 across loci of an individual. Three method-of-moment estimators (Loiselle *et al.* 1995;
169 Ritland 1996; Weir 1996) are employed to estimate the individual kinship coefficient and
170 inbreeding coefficient.

171 **Genetic differentiation:** estimates the differentiation index F_{ST} and tests for genetic
172 differentiation between/among populations/regions. The F_{ST} estimators available are:
173 Nei's (1973) G_{ST} , Weir & Cockerham's (1984)'s θ , Hudson *et al.*'s (1992), Slatkin's (1995)
174 R_{ST} , Hedrick's (2005) G'_{ST} and Jost's (2008) D and Huang *et al.*'s (2021) variance decom-
175 position. All estimators can be applied to both polyploids and anisoploids with the excep-
176 tion of Weir & Cockerham's (1984)'s θ . Differentiation is tested using Fisher's G -test for all
177 loci or for each locus based on the distribution of genotypes or alleles.

178 **Genetic distance:** this enables the calculation of a variety of genetic distance indices
179 between individuals, populations or regions: Nei's (1972) standard genetic distance,
180 Cavalli-Sforza's (1967) chord distance, Reynolds *et al.*'s (1983) θ_w , Nei's (1983) D_A dis-
181 tance, Euclidean distance, Goldstein's (1995) distance, Nei's (1974) minimum genetic dis-
182 tance, Roger's (1972) distance, and two F_{ST} transformation-based genetic distances: (i)
183 Reynolds *et al.*'s (1983) D and (ii) Slatkin's (1995) linearized F_{ST} . Based on the estimated

184 genetic distance matrices, the *principal coordinate analysis* (PCoA) and the *hierarchical*
185 *clustering analysis* (HCA) can be performed.

186 **Analysis of molecular variance:** Classical AMOVA only supports data for haploids
187 and diploids, and will only support from one to four hierarchies. Based on the generalized
188 framework (Huang *et al.* 2021), we extended AMOVA to accommodate any ploidy level
189 and any number of hierarchies. The generalized framework models the symbolic expres-
190 sion of the expected *Sum of Squares* (SS) \mathbf{S} by variance components $\mathbf{\Sigma}$, which can be writ-
191 ten as $\mathbf{S} = \mathbf{C}\mathbf{\Sigma}$. The method-of-moment estimate of variance components can be solved by
192 $\hat{\mathbf{\Sigma}} = \mathbf{C}^{-1}\hat{\mathbf{S}}$. Three methods are provided: (i) the homoploid method, (ii) the anisoploid
193 method and (iii) the likelihood method. The homoploid method uses the dummy haplo-
194 type method as in GENALEX (Peakall & Smouse 2006), and combines all loci into one
195 dummy locus, then calculate $\hat{\mathbf{S}}$ and \mathbf{C} using the dummy haplotypes. This method can
196 only be applied to homoploids and is slightly biased when there are missing data. The
197 anisoploid model can be applied to both homoploids and anisoploids, which calculates $\hat{\mathbf{S}}$
198 and \mathbf{C} using the alleles for each locus. The matrices $\hat{\mathbf{S}}$ and \mathbf{C} are summed over loci, and
199 the variance component matrix $\hat{\mathbf{\Sigma}}$ is solved at once. Because the anisoploid method per-
200 mutes alleles at each locus to test the significance, it takes increased calculation time com-
201 pared with the homoploid method. VCFPOP uses a pseudo-permutation method to solve
202 this problem, which first perform a small number of permutations (e.g., 100) for each locus,
203 then subsamples one permutation at each locus to generate results for each pseudo-per-
204 mutation. For the likelihood method, the F -statistics are first estimated by maximizing gen-
205 otypic likelihood under differentiation or subdivision, and the variance components and

206 other statistics are subsequently solved.

207 **Population assignment:** enables the calculation of the likelihood for each individual
208 of being a member of a particular population (Paetkau *et al.* 2004). Each individual is as-
209 signed to the population with the maximum likelihood. Such assignment can help identify
210 the natal population for an individual.

211 **Kinship coefficient:** this calculates the kinship coefficient (θ) between two individu-
212 als, the probability that two randomly sampled alleles, each from one individual, are *iden-*
213 *tical-by-descent* (IBD). The same estimators (Loiselle *et al.* 1995; Ritland 1996; Weir 1996) are
214 employed to estimate the kinship coefficient between individuals.

215 **Relatedness coefficient:** this estimates the relatedness coefficient between individuals.
216 Two native polyploid relatedness estimators that supports a maximum ploidy level of
217 eight are provided: (i) the method-of-moment estimator (Huang *et al.* 2014), and (ii) the
218 maximum-likelihood estimator (Huang *et al.* 2015a). Moreover, the relatedness coefficient
219 can also be transformed from the kinship coefficient, with VCFPOP providing two transfor-
220 mations. The original transformation is used for outbred populations and is performed by

$$221 \hat{r}_{HL} = v_{\min} \hat{\theta}_{xy},$$

222 Where \hat{r}_{HL} is the estimate of the relatedness coefficient from a higher ploidy individual to
223 a lower ploidy individual, $\hat{\theta}_{xy}$ is the kinship coefficient between these two individuals,
224 v_{\min} is the ploidy level of the lower ploidy individual (Huang *et al.* 2015b). The modified
225 transformation accommodates inbreeding and is performed by

$$226 \hat{r}_{HL} = \frac{v_{\min}}{v_{\min} + v_{\max}} \hat{\theta}_{xy} \left(\frac{1}{\hat{\theta}_{xx}} + \frac{1}{\hat{\theta}_{yy}} \right),$$

227 where v_{\max} is the ploidy level of the higher ploidy individual, and $\hat{\theta}_{xx}$ (or $\hat{\theta}_{yy}$) is the

228 kinship coefficient within the individual x (or y) (Huang *et al.* 2015a).

229 **Bayesian clustering:** This estimates ancestral proportions of each individual by the
230 Markov Chain Monte Carlo (MCMC) method. VCFPOP follows the software STRUCTURE and
231 implements three models of Bayesian clustering: (i) the ADMIXTURE model (Pritchard *et al.*
232 2000), (ii) the LOCPRIORI model (Hubisz *et al.* 2009) and (iii) the F model (Falush *et al.* 2003).

233 Optimization

234 NGS datasets are usually large, with a single VCF file often reaching hundreds of gi-
235 gabytes in size. Workstations or computer clusters are usually required to analyze such
236 data. VCFPOP uses various of methods to reduce the requirement and exploit the capacity
237 of the computer:

- 238 (i) Optimized algorithm (to accelerate the calculation);
- 239 (ii) Advanced instruction set (e.g., SSE, LZCNT, POPCNT, AVX, FMA, AVX512);
- 240 (iii) Lock-free technology (to reduce access conflicts among threads);
- 241 (iv) Virtual memory allocation (to avoid re-allocation and memory move);
- 242 (v) Local memory management class (to allocate millions of small pieces of memory);
- 243 (vi) Variable length array (to place temporary array on stack memory);
- 244 (vii) Fast hash algorithm (to detect identical genotypes);
- 245 (viii) Fast hash table (to access genotypes by either hash value or index);
- 246 (ix) Indexing alleles and genotype (to share instances and reduce memory expense);
- 247 (x) Memory cache (to avoid frequent disk I/O).

248 VCFPOP is optimized for memory expense and calculation speed. For Intel® processors

249 later than SkyLake (released in 2017) and CannonLake (released in 2018), the AVX-512
250 SIMD instructions can be used to accelerate the calculation speed, which enables the pro-
251 cessing of 512 bits simultaneously. For Intel ® processors later than Haswell hardware (re-
252 leased in 2013) and AMD ® processors after Excavator hardware (released in 2015), the
253 AVX instructions can be used, which enables the processing of 256 bits simultaneously.
254 These SIMD instruction-sets can be flexibly switched without additional compilation.

255 For memory usage, VCFPOP indexes the individual genotypes and uses bitwise storage
256 to save the genotype index. The memory usage for a genotype at a biallelic locus (e.g.,
257 single nucleotide polymorphism) is reduced from 4 bytes (e.g., '0/0 ') to 2 bits (4 possible
258 states: *AA*, *AB*, *BB* and missing data) for VCF format (16-folds compression). The detailed
259 information (e.g., ploidy level, alleles) of genotypes is saved in an additional table, with
260 each genotype using approximately 12 bytes. Because additional memory is required for
261 information at the locus, individual and population levels, the typical compression ratio is
262 14.5-fold evaluated by the Chr 22 phased 3 data of the 1000 genome project, which uses
263 734 Mib memory to load 10.4Gib data.

264 Although the compression ratio can reach 50-fold for some professional compression
265 algorithms (e.g., deflate, lzma, z-std), the random access of genotypes requires a longer
266 decompression time. Our method can compress data at a considerable compression ratio,
267 and can also access the genotypes without additional cost. In other words, there is a trade-
268 off between the accession speed and memory usage.

269 Therefore, a typical laptop with 16 GiB of memory can process 100 GiB VCF files with-
270 out considering the calculation speed. The subsequent analysis methods may also require

271 additional memory. For example, the homoploid AMOVA method saves the genetic dis-
272 tance between dummy haplotypes and requires an additional $8H^2$ bytes (about 745 MiB
273 in a dataset with 5000 diploids), where H is the number of dummy haplotypes.

274 The loading speed of VCFPOP is also optimized, and uses a single thread to read the
275 data from the disk and multiple threads to process the data. With a sample benchmark test
276 of a 10.4 GiB uncompressed VCF file, VCFPOP can load data at 320 MiB/s and 560 MiB/s on
277 a laptop (Intel i7-8750H CPU with 2.2GHz and 6 cores, 16 GiB memory, 256 Gib nvme SSD)
278 and a workstation (Intel Xeon E5-2696 V4 CPU with 2.2GHz and 44 cores, 64 GiB memory
279 and 1 TiB nvme SSD), respectively. Restricted by the additional decompression process,
280 the loading speed is reduced to 255 and 460 MiB/s on the laptop and the workstation for
281 the compressed format (`vcf.gz`), respectively.

282 Acknowledgements

283 This work was supported by the Strategic Priority Research Program of the Chinese Acad-
284 emy of Sciences (XDB31020302), the National Natural Science Foundation of China (31730104,
285 32170515, 31770411, 32070453), and the Innovation Capability Support Program of Shaanxi
286 (2021KJXX-027). KH is supported by a scholarship from China Scholarship Council. KH would
287 like to thank Prof. Kermit Ritland for providing a visiting professor position at the University
288 of British Columbia.

289 References

290 Avni R, Nave M, Barad O, *et al.* (2017) Wild emmer genome architecture and diversity elucidate

291 wheat evolution and domestication. *Science* **357**, 93-97.

292 Burow MD, Simpson CE, Starr JL, Paterson AH (2001) Transmission genetics of chromatin from
 293 a synthetic amphidiploid to cultivated peanut (*Arachis hypogaea* L.): broadening the gene
 294 pool of a monophyletic polyploid species. *Genetics* **159**, 823-837.

295 Butruille DV, Boiteux LS (2000) Selection–mutation balance in polysomic tetraploids: impact of
 296 double reduction and gametophytic selection on the frequency and subchromosomal
 297 localization of deleterious mutations. *Proceedings of the National Academy of Sciences* **97**,
 298 6608-6613.

299 Cavalli-Sforza LL, Edwards AW (1967) Phylogenetic analysis. Models and estimation
 300 procedures. *American Journal of Human Genetics* **19**, 233.

301 Clark LV, Jasieniuk M (2011) POLYSAT: an R package for polyploid microsatellite analysis.
 302 *Molecular Ecology Resources* **11**, 562-566.

303 Danecek P, Auton A, Abecasis G, *et al.* (2011) The variant call format and VCFtools.
 304 *Bioinformatics* **27**, 2156-2158.

305 Excoffier L, Lischer HE (2010) Arlequin suite ver 3.5: a new series of programs to perform
 306 population genetics analyses under Linux and Windows. *Molecular Ecology Resources* **10**,
 307 564-567.

308 Falush D, Stephens M, Pritchard JK (2003) Inference of population structure using multilocus
 309 genotype data: linked loci and correlated allele frequencies. *Genetics* **164**, 1567-1587.

310 Goldstein DB, Ruiz LA, Cavallisforza LL, Feldman MW (1995) Genetic absolute dating based
 311 on microsatellites and the origin of modern humans. *Proceedings of the National Academy of*
 312 *Sciences of the United States of America* **92**, 6723-6727.

313 Haldane JB (1930) Theoretical genetics of autopolyploids. *Journal of Genetics* **22**, 359-372.

314 Hardy OJ (2016) Population genetics of autopolyploids under a mixed mating model and the
 315 estimation of selfing rate. *Molecular Ecology Resources* **16**, 103-117.

316 Hardy OJ, Vekemans X (2002) SPAGeDi: a versatile computer program to analyse spatial
 317 genetic structure at the individual or population levels. *Molecular Ecology Notes* **2**, 618-620.

318 Hedrick PW (2005) A standardized genetic differentiation measure. *Evolution* **59**, 1633-1638.

319 Huang K, Dunn DW, Ritland K, Li B (2020) polygene: Population genetics analyses for
 320 autopolyploids based on allelic phenotypes. *Methods in Ecology and Evolution* **11**, 448-456.

321 Huang K, Guo ST, Shattuck MR, *et al.* (2015a) A maximum-likelihood estimation of pairwise
 322 relatedness for autopolyploids. *Heredity* **114**, 133-142.

323 Huang K, Ritland K, Guo S, Shattuck M, Li B (2014) A pairwise relatedness estimator for
 324 polyploids. *Molecular Ecology Resources* **14**, 734-744.

325 Huang K, Ritland K, Guo ST, *et al.* (2015b) Estimating pairwise relatedness between individuals
 326 with different levels of ploidy. *Molecular Ecology Resources* **15**, 772–784.

327 Huang K, Wang T, Dunn DW, *et al.* (2021) A generalized framework for AMOVA with multiple
 328 hierarchies and ploidies. *Integrative Zoology* **16**, 33-52.

329 Huang K, Wang TC, Dunn DW, *et al.* (2019) Genotypic frequencies at equilibrium for polysomic
 330 inheritance under double-reduction. *G3: Genes, Genomes, Genetics* **9**, 1693-1706.

331 Hubisz MJ, Falush D, Stephens M, Pritchard JK (2009) Inferring weak population structure with
 332 the assistance of sample group information. *Molecular Ecology Resources* **9**, 1322-1332.

333 Hudson RR, Slatkin M, Maddison W (1992) Estimation of levels of gene flow from DNA
 334 sequence data. *Genetics* **132**, 583-589.

335 Jost L (2008) G_{ST} and its relatives do not measure differentiation. *Molecular Ecology* **17**, 4015-4026.

336 Kalinowski ST, Taper ML, Marshall TC (2007) Revising how the computer program CERVUS
337 accommodates genotyping error increases success in paternity assignment. *Molecular*
338 *Ecology* **16**, 1099-1106.

339 Ling HQ, Ma B, Shi XL, *et al.* (2018) Genome sequence of the progenitor of wheat A subgenome
340 *Triticum urartu*. *Nature* **557**, 424.

341 Loiselle BA, Sork VL, Nason J, Graham C (1995) Spatial genetic structure of a tropical
342 understory shrub, *Psychotria officinalis* (Rubiaceae). *American Journal of Botany* **82**, 1420-
343 1425.

344 Mather K (1935) Reductional and equational separation of the chromosomes in bivalents and
345 multivalents. *Journal of Genetics* **30**, 53-78.

346 Meirmans PG, Tienderen PHV (2004) GENOTYPE and GENODIVE : two programs for the
347 analysis of genetic diversity of asexual organisms. *Molecular Ecology Notes* **4**, 792-794.

348 Muller HJ (1914) A new mode of segregation in Gregory's tetraploid primulas. *The American*
349 *Naturalist* **48**, 508-512.

350 Nei M (1972) Genetic distance between populations. *American Naturalist* **106**, 283-292.

351 Nei M (1973) Analysis of gene diversity in subdivided populations. *Proceedings of the National*
352 *Academy of Sciences* **70**, 3321-3323.

353 Nei M, Roychoudhury AK (1974) Genic variation within and between the three major races of
354 man, Caucasoids, Negroids, and Mongoloids. *American Journal of Human Genetics* **26**, 421-
355 443.

356 Nei M, Tajima F, Tateno Y (1983) Accuracy of estimated phylogenetic trees from molecular data
357 II. Gene frequency data. *Journal of Molecular Evolution* **19**, 153-170.

358 Paetkau D, Slade R, Burden M, Estoup A (2004) Genetic assignment methods for the direct,
359 real-time estimation of migration rate: a simulation-based exploration of accuracy and
360 power. *Molecular Ecology* **13**, 55-65.

361 Peakall R, Smouse PE (2006) GENALEX 6: genetic analysis in Excel. Population genetic
362 software for teaching and research. *Molecular Ecology Notes* **6**, 288-295.

363 Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using multilocus
364 genotype data. *Genetics* **155**, 945-959.

365 Reynolds J, Weir BS, Cockerham CC (1983) Estimation of the coancestry coefficient: basis for a
366 short-term genetic distance. *Genetics* **105**, 767-779.

367 Ritland K (1996) Estimators for pairwise relatedness and individual inbreeding coefficients.
368 *Genetical Research* **67**, 175-185.

369 Rogers JS (1972) Measures of similarity and genetic distance. In: *Studies in Genetics VII* (ed.
370 Wheeler MR), pp. 145-153. University of Texas Publication, Austin.

371 Rousset F (2008) genepop'007: a complete re-implementation of the genepop software for
372 Windows and Linux. *Molecular Ecology Resources* **8**, 103-106.

373 Slatkin M (1995) A measure of population subdivision based on microsatellite allele frequencies.
374 *Genetics* **139**, 457-462.

375 Weir BS (1996) *Genetic data analysis II: methods for discrete population genetic data* Sinauer
376 Associates, Sunderland.

377 Weir BS, Cockerham CC (1984) Estimating F -statistics for the analysis of population structure.
378 *Evolution* **38**, 1358-1370.

379

380 **Data Accessibility and Benefit-Sharing** 381 **Section**

382 The source code, binary executables, user manual and example files (input files, parame-
383 ters, and commands) are available from GitHub (<http://github.com/huangkang1987/vcfpop>).

384 Benefits from this research accrue from the sharing of our software and source code on
385 public databases as described above.

386

387 **Author Contributions**

388 KH and BGL designed the project, KH and BY designed the software and wrote the draft,
389 YFW and YXC reviewed the code, JCA, YHL and YCK performed simulations and tests, and
390 DWD checked the model and helped write the manuscript.

391